

# Operating System

*Dr. GuoJun LIU*

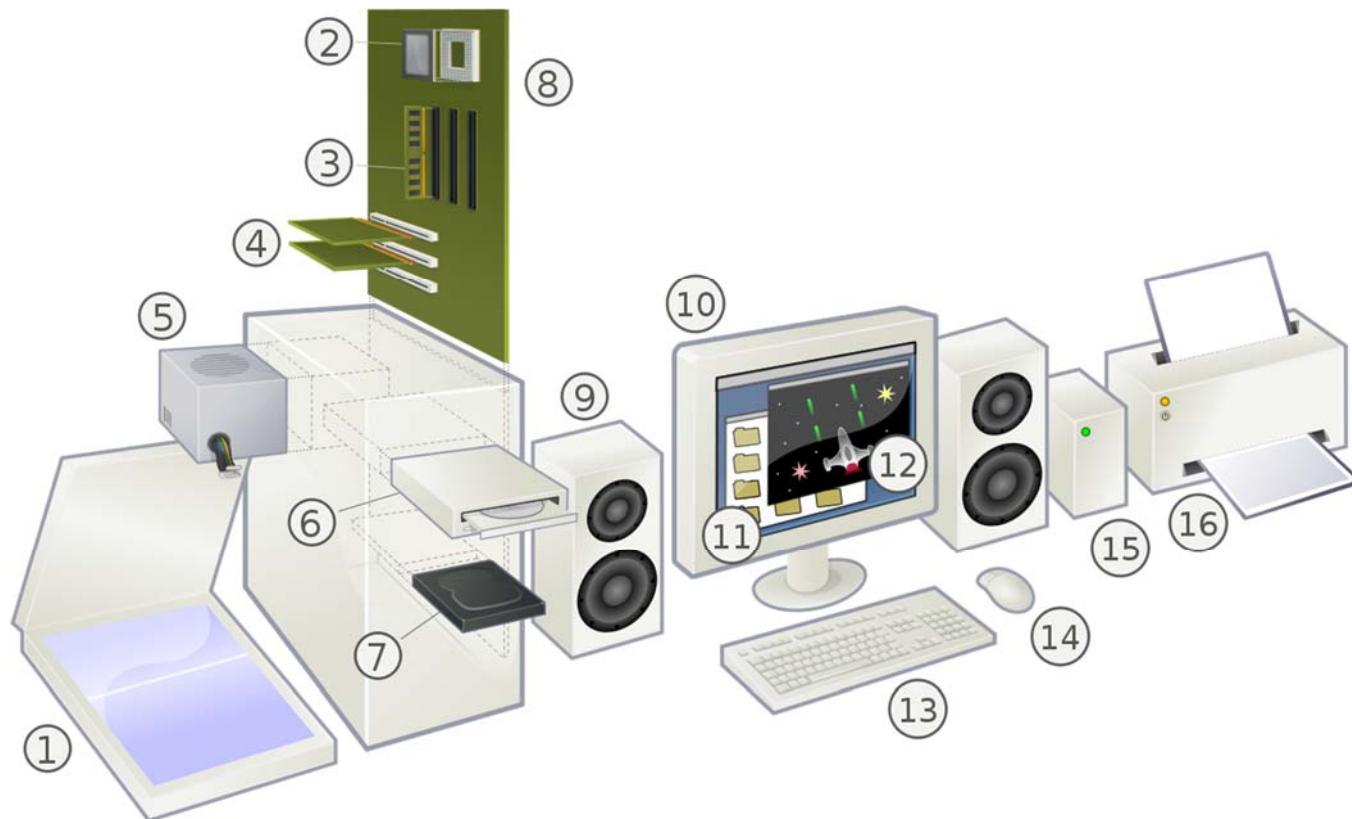
Harbin Institute of Technology

<http://guojun.hit.edu.cn/os/>

# *My Summary of OS*

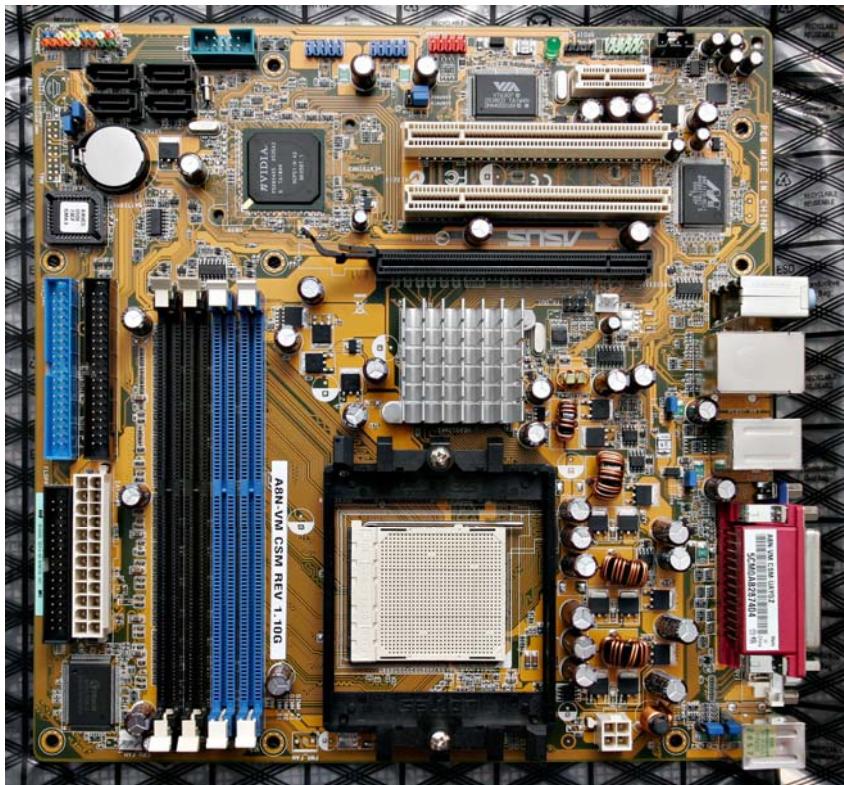
... ...

# A modern PC and peripherals



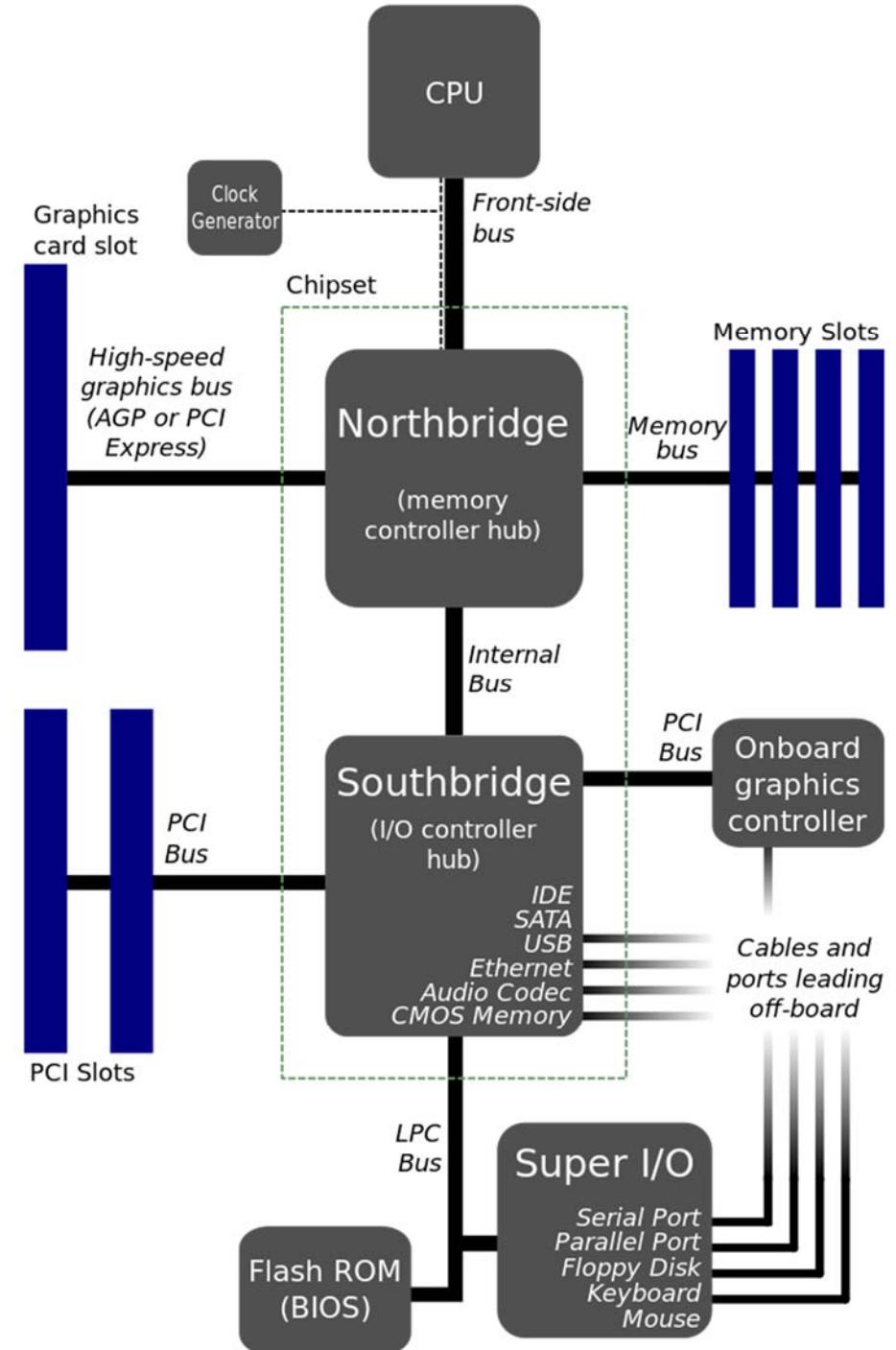
An exploded view of a modern personal computer and peripherals

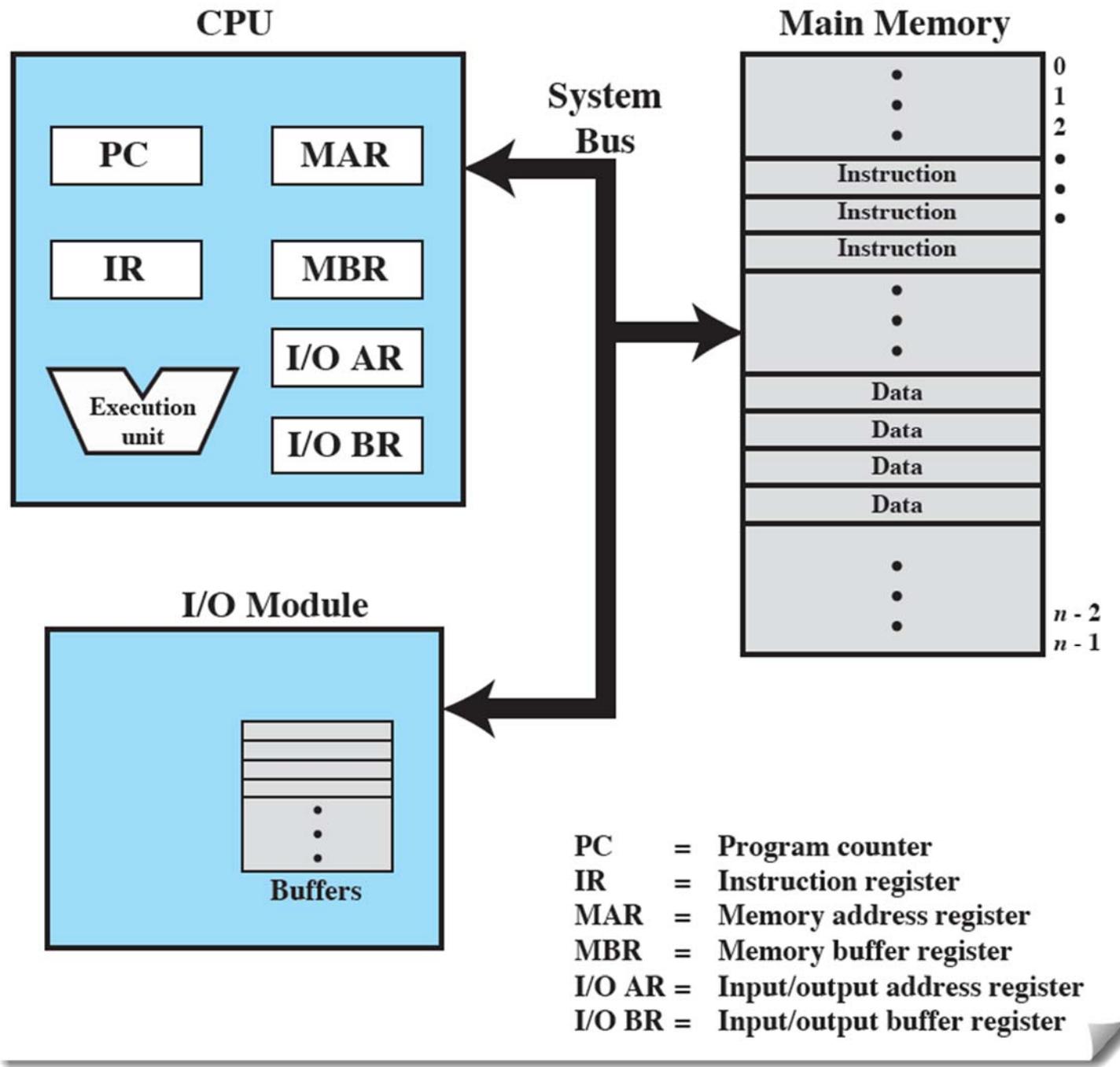
1. Scanner
2. CPU
3. Memory (RAM)
4. Expansion cards  
(graphics cards, etc.)
5. Power supply
6. Optical disc drive
7. Storage (Hard disk)
8. Motherboard
9. Speakers
10. Monitor
11. System software
12. Application software
13. Keyboard
14. Mouse
15. External hard disk
16. Printer



Motherboard

*Photo from wikipedia.org*





Computer Components: Top-Level View

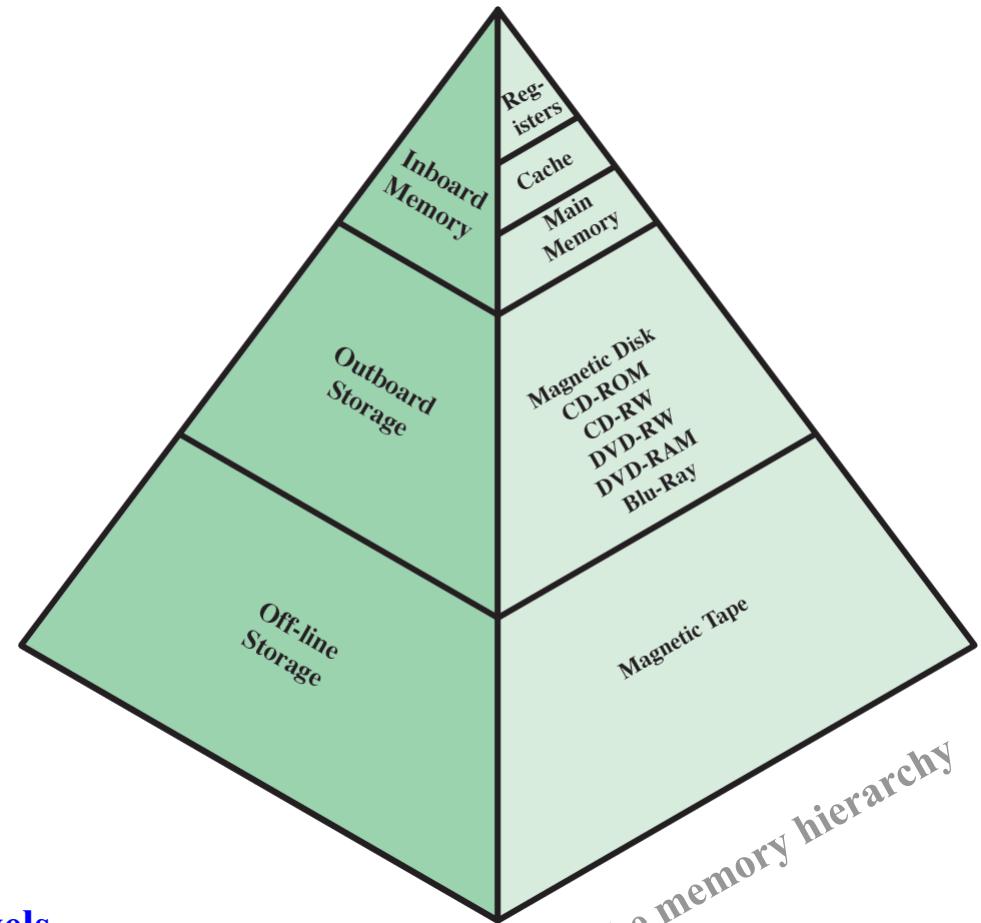
# The Memory Hierarchy

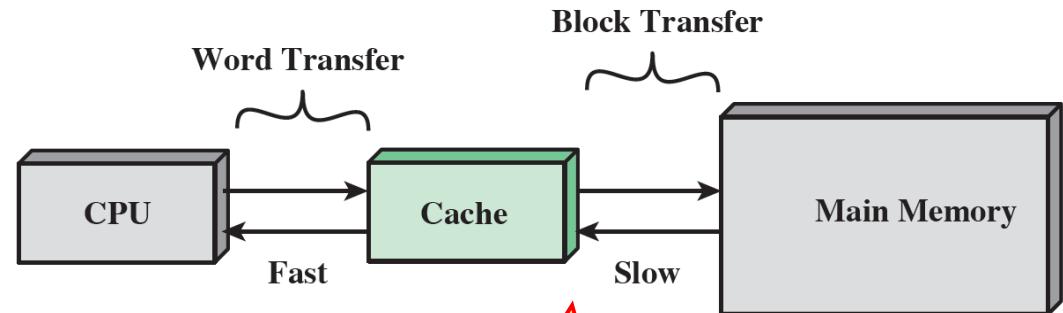
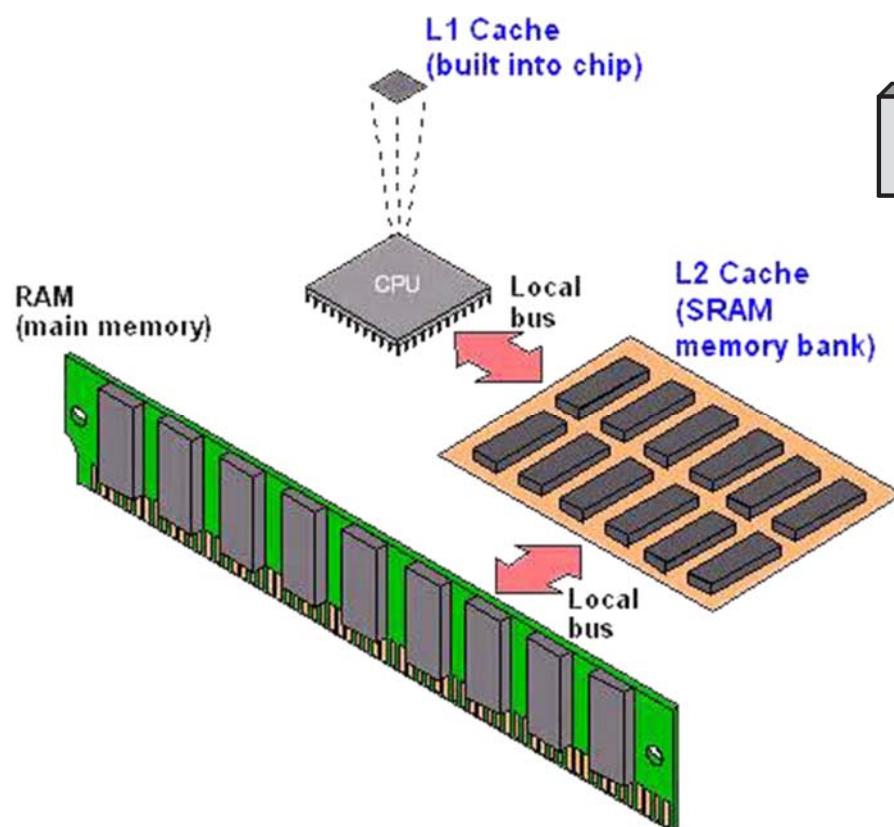
- Going down the hierarchy:

- a. decreasing cost per bit
- b. increasing capacity
- c. increasing access time
- d. **decreasing frequency** of access to the **memory** by the **processor**

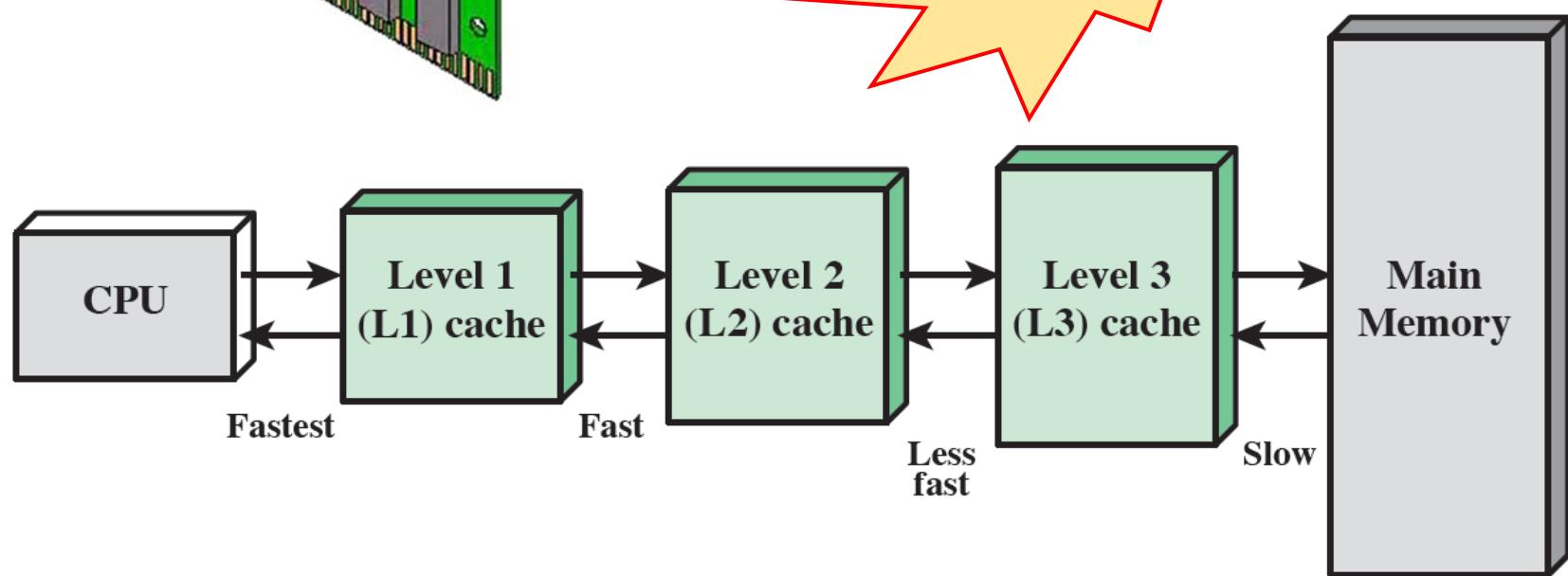


**Decreasing frequency of access at lower levels**

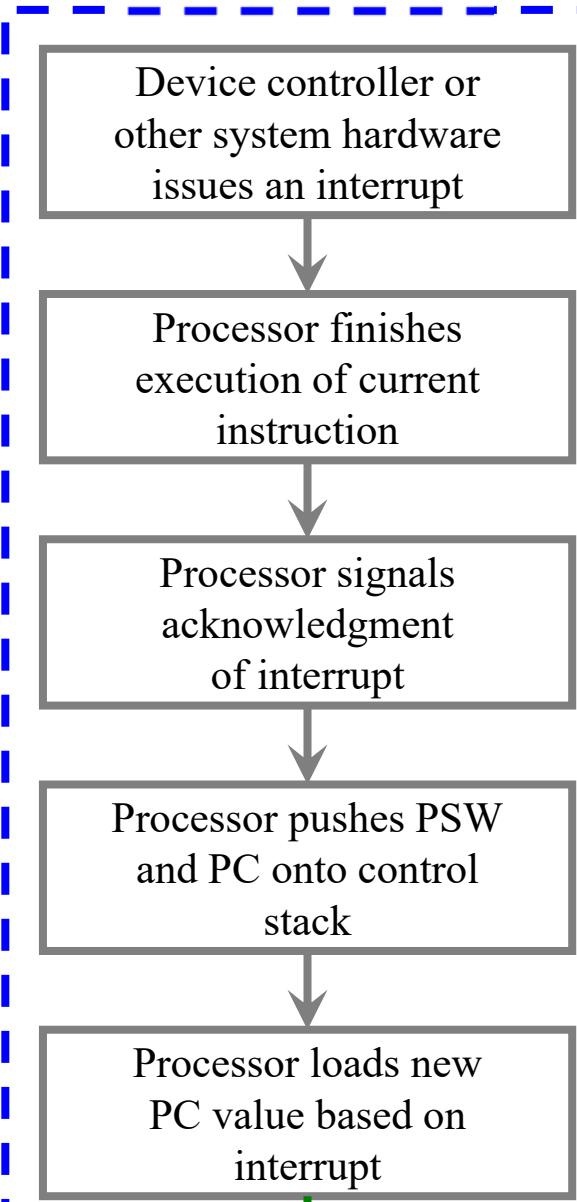




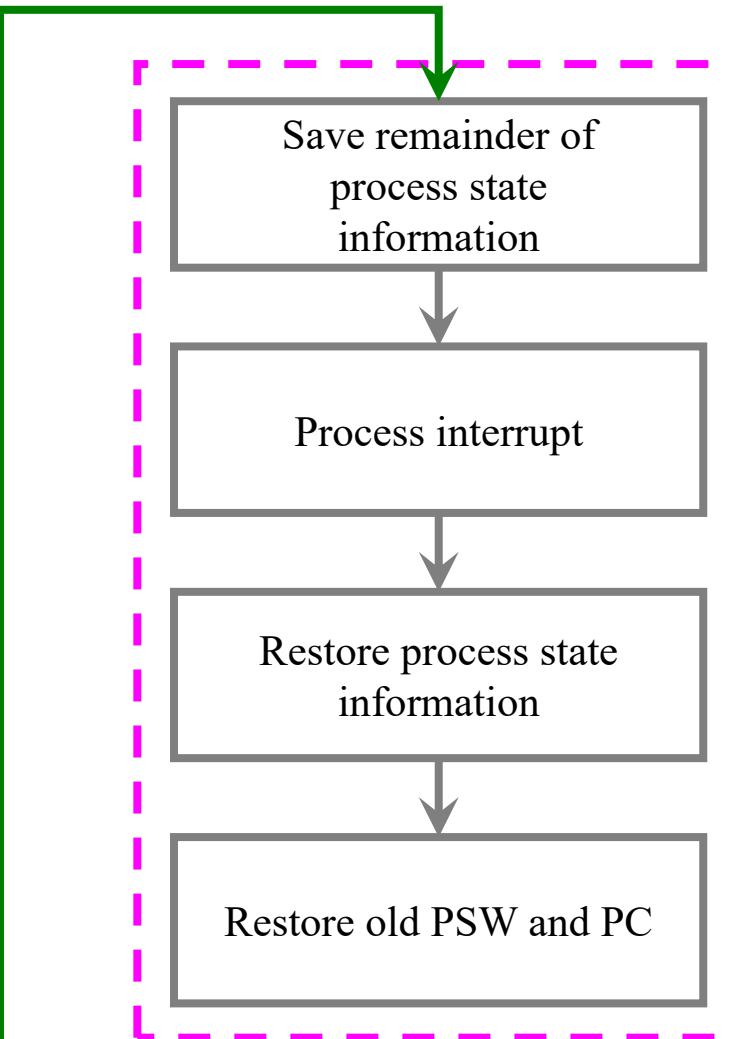
The cache is not  
usually visible to  
the programmer  
or, indeed,  
to the processor



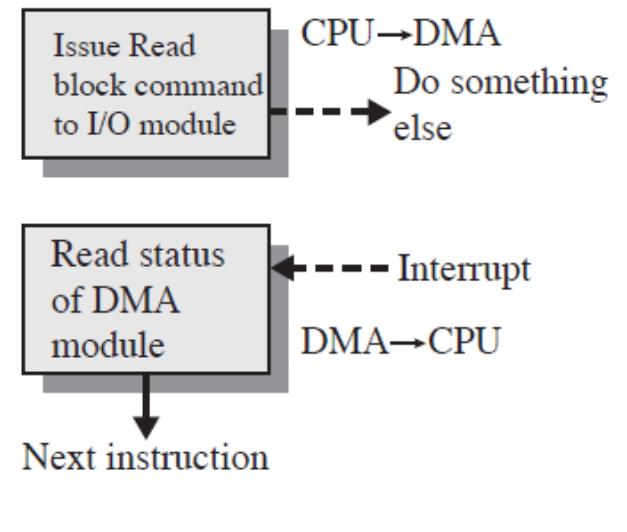
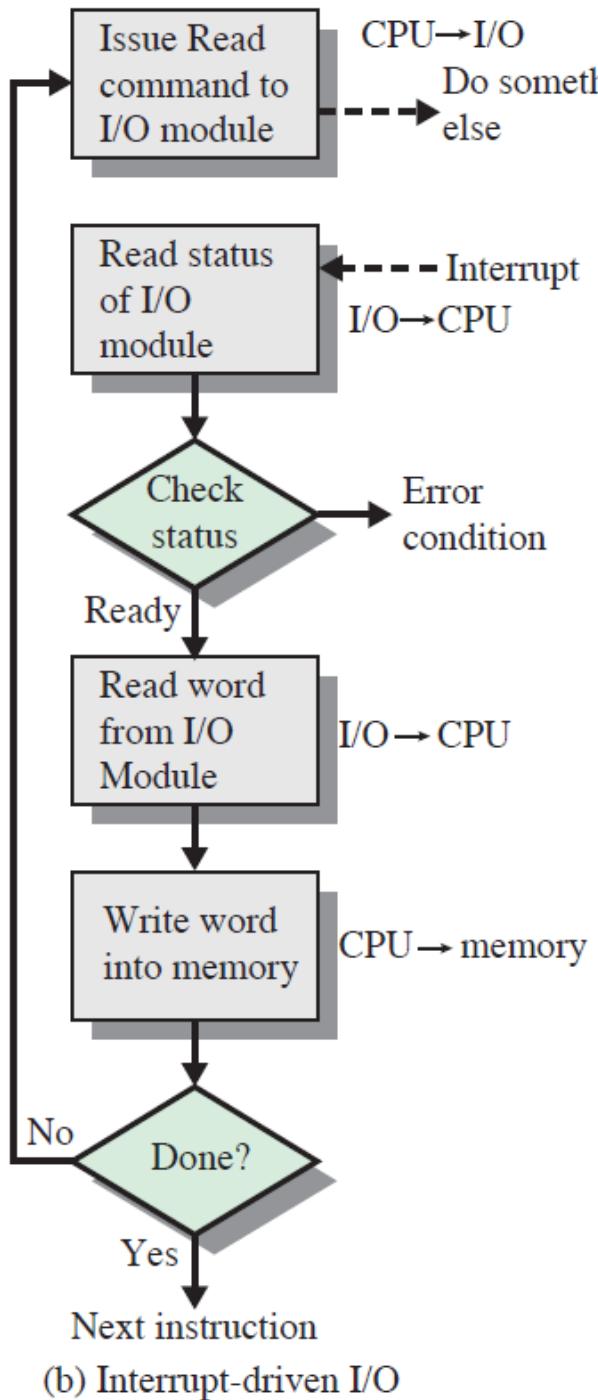
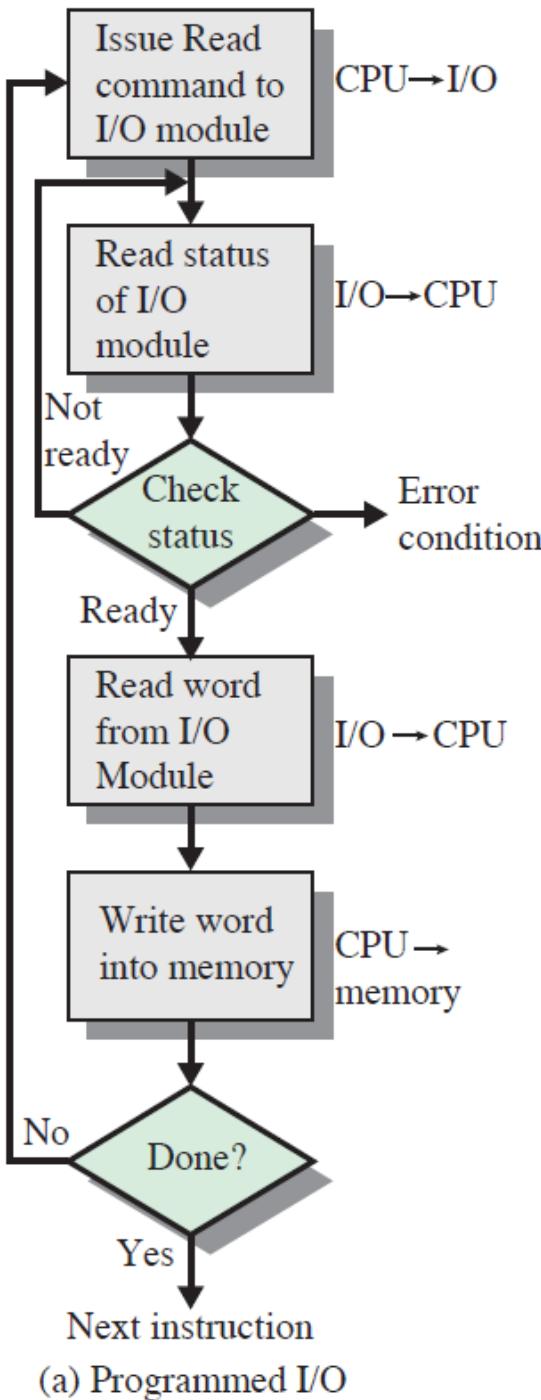
## Hardware



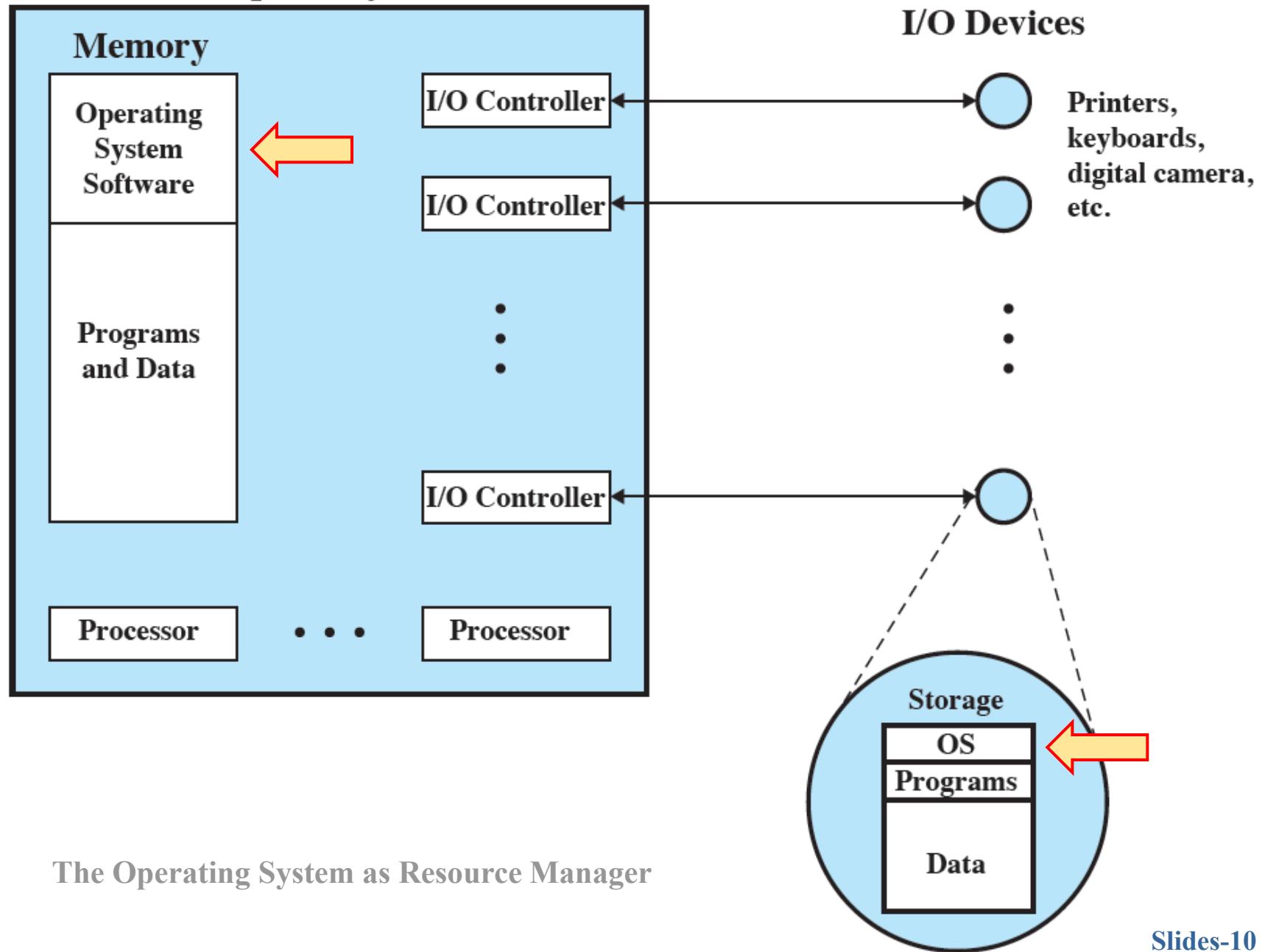
## Software



Simple Interrupt Processing

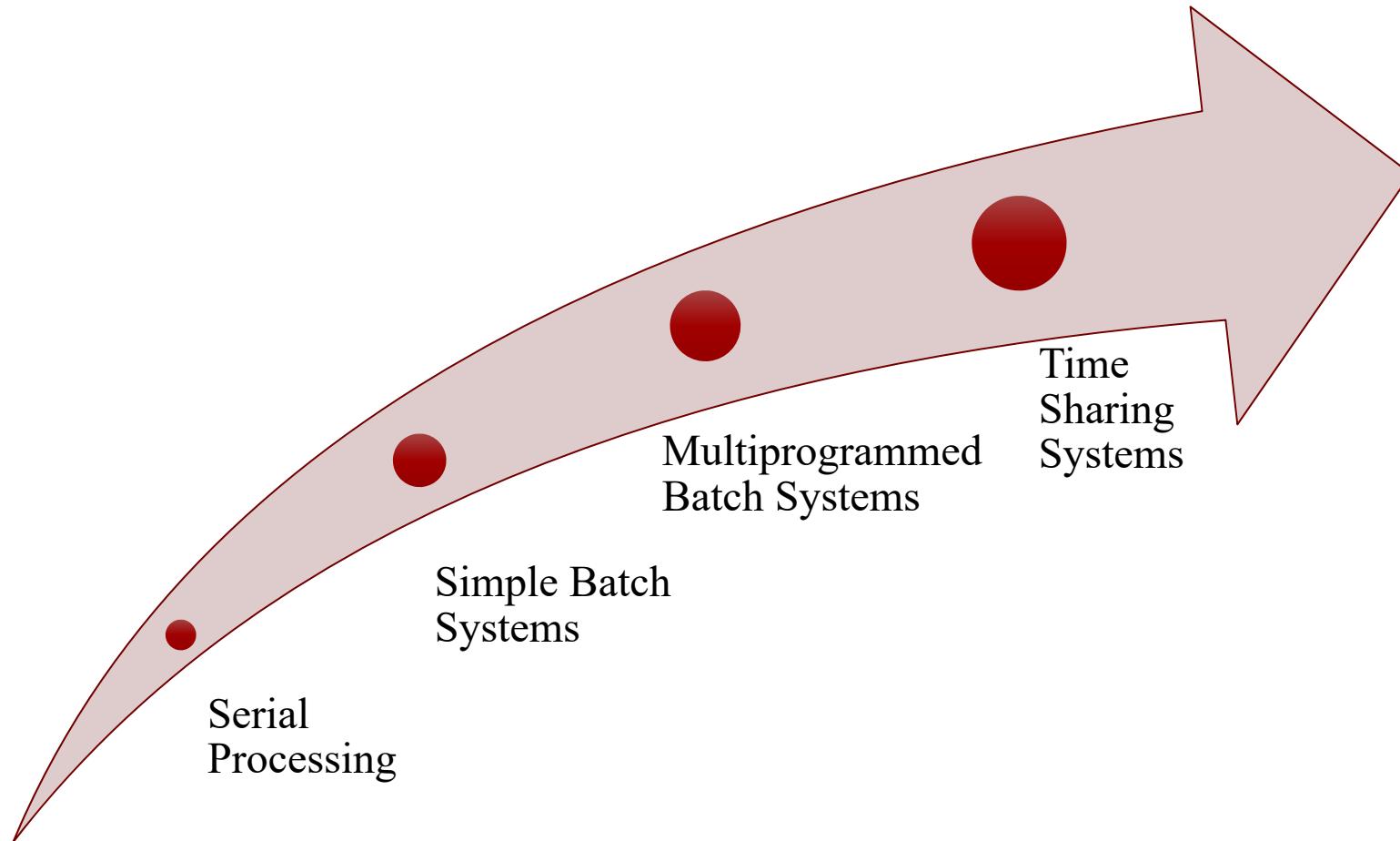


# Computer System

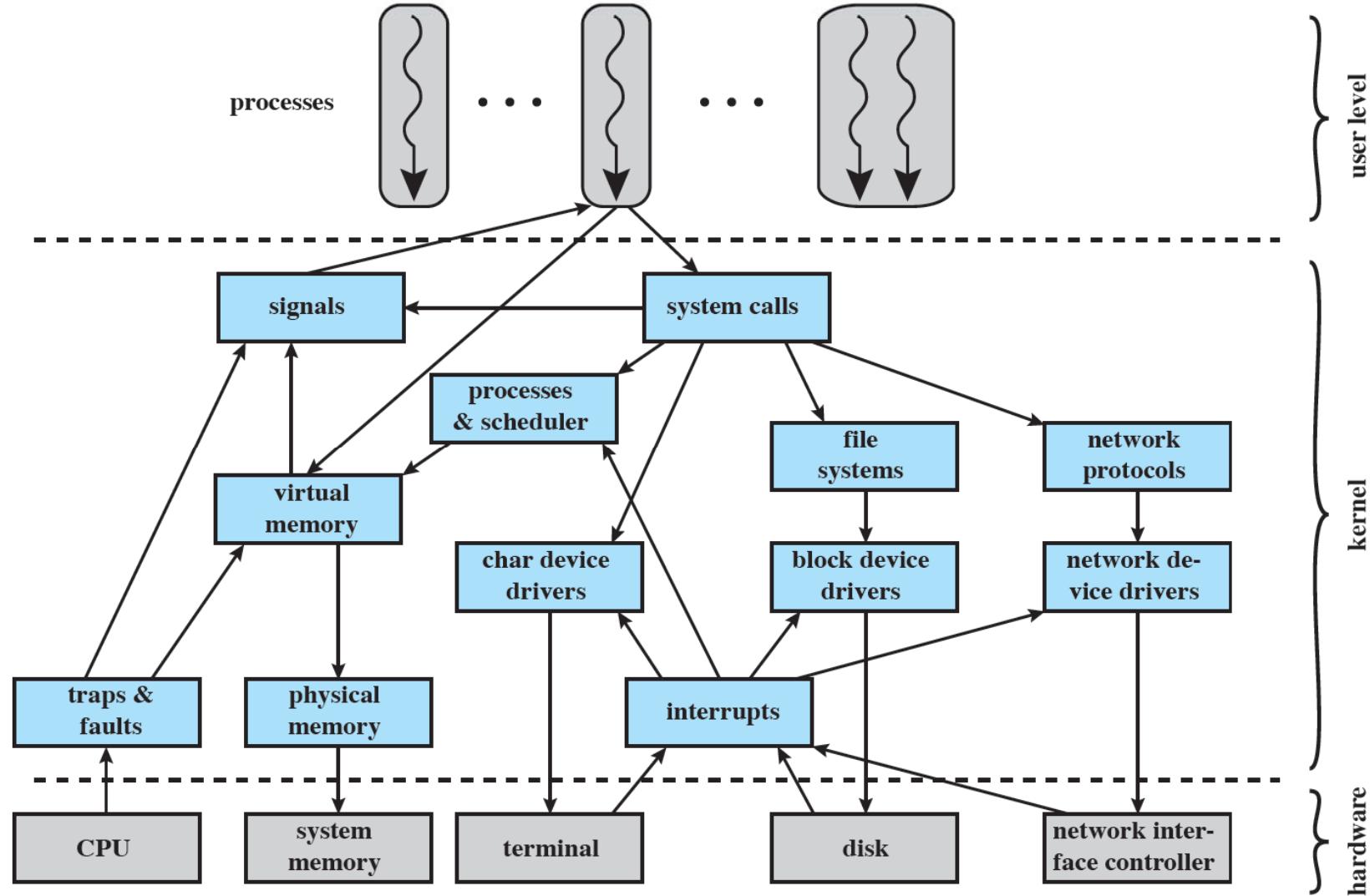


# Evolution of Operating Systems

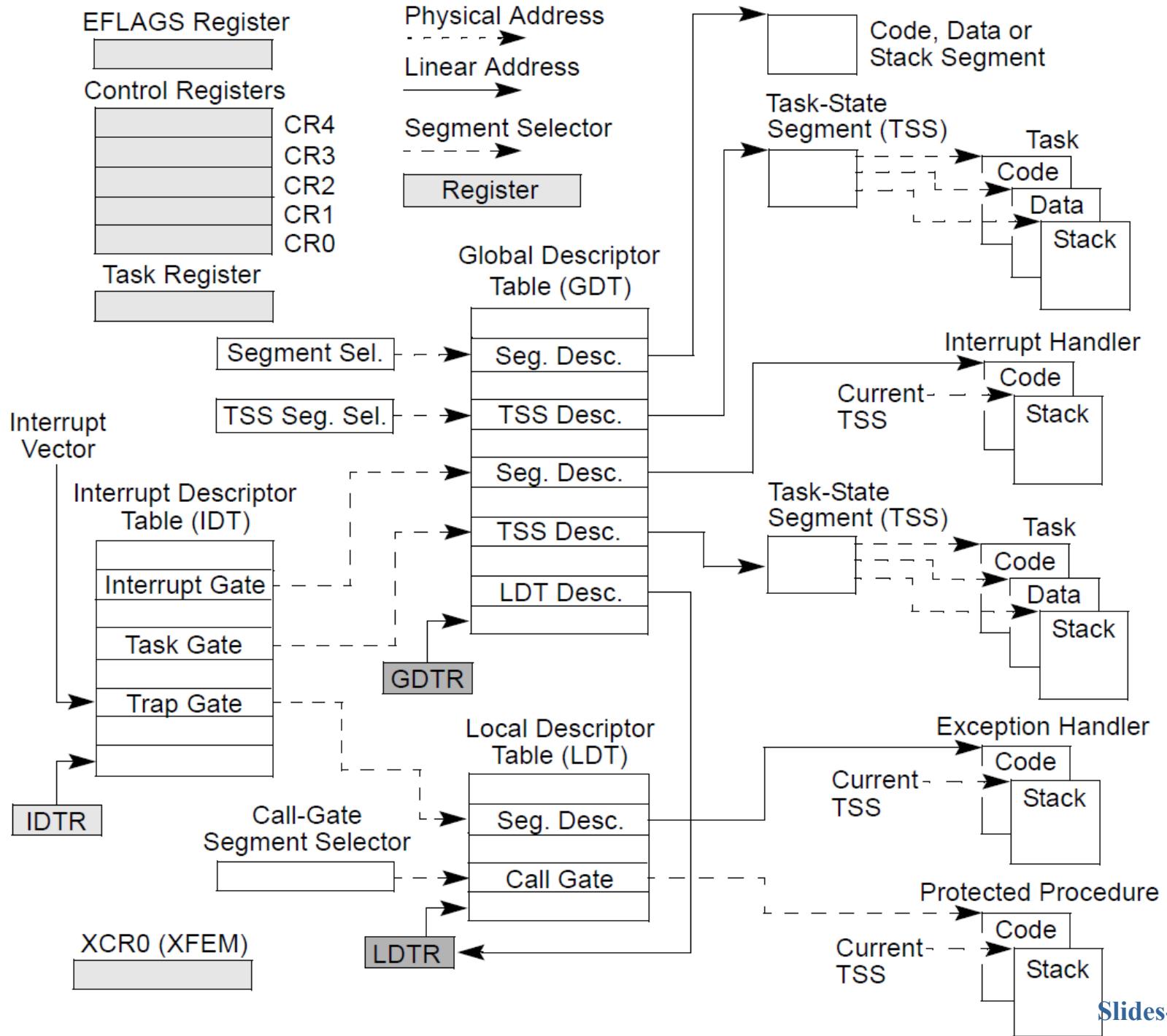
---



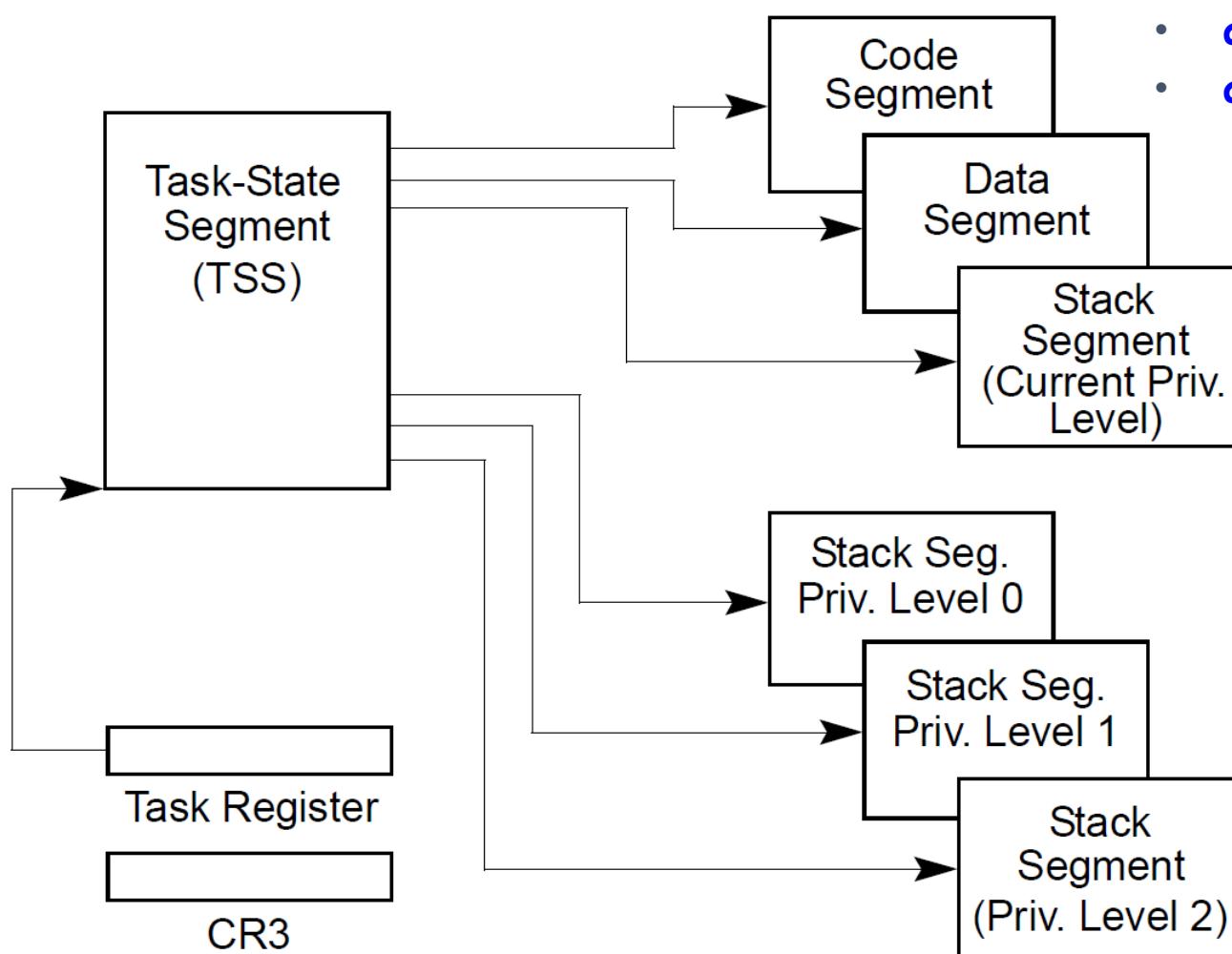
# Linux Kernel Components



# IA-32 System-Level Registers and Data Structures

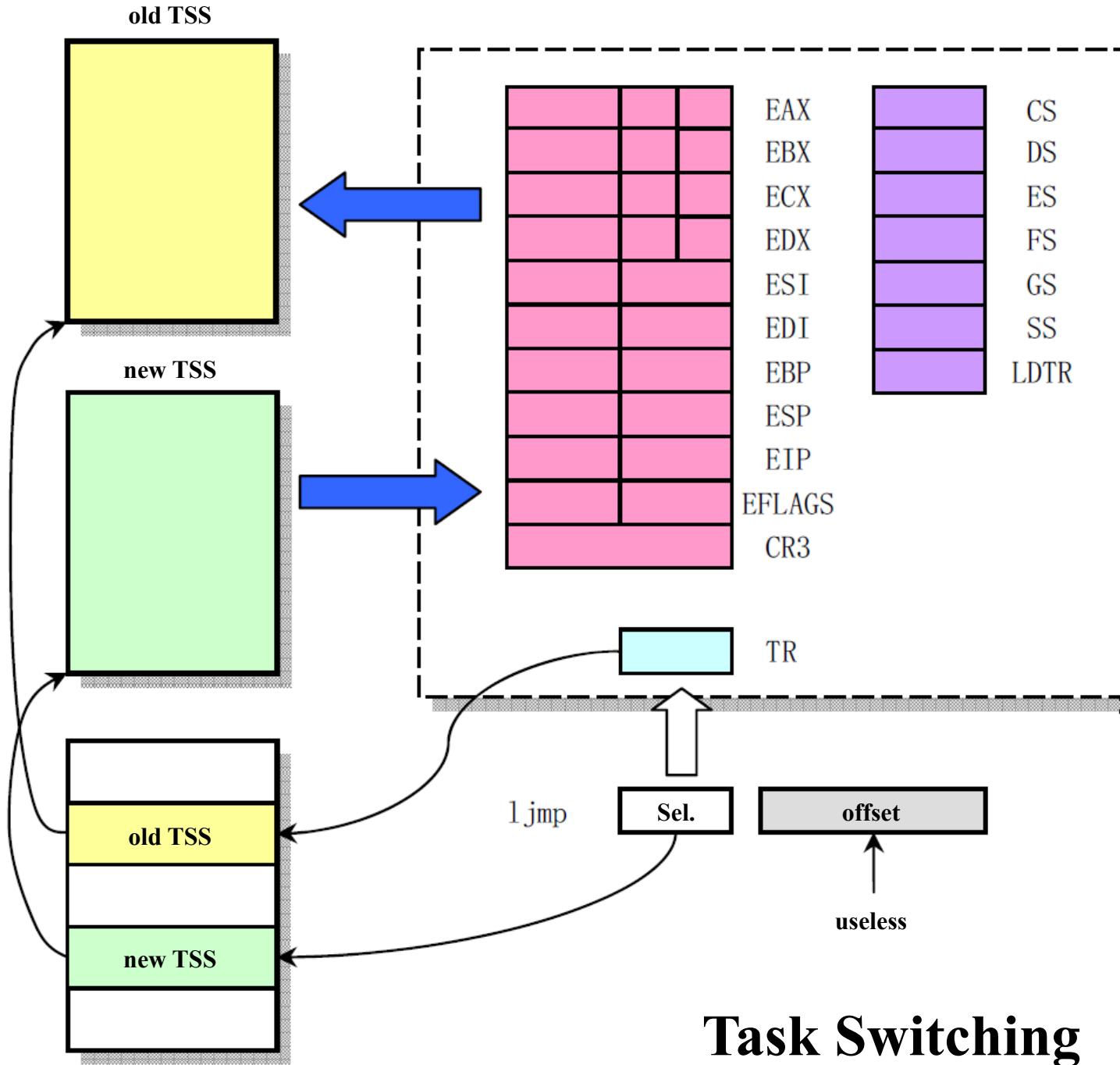


# Structure of a Task



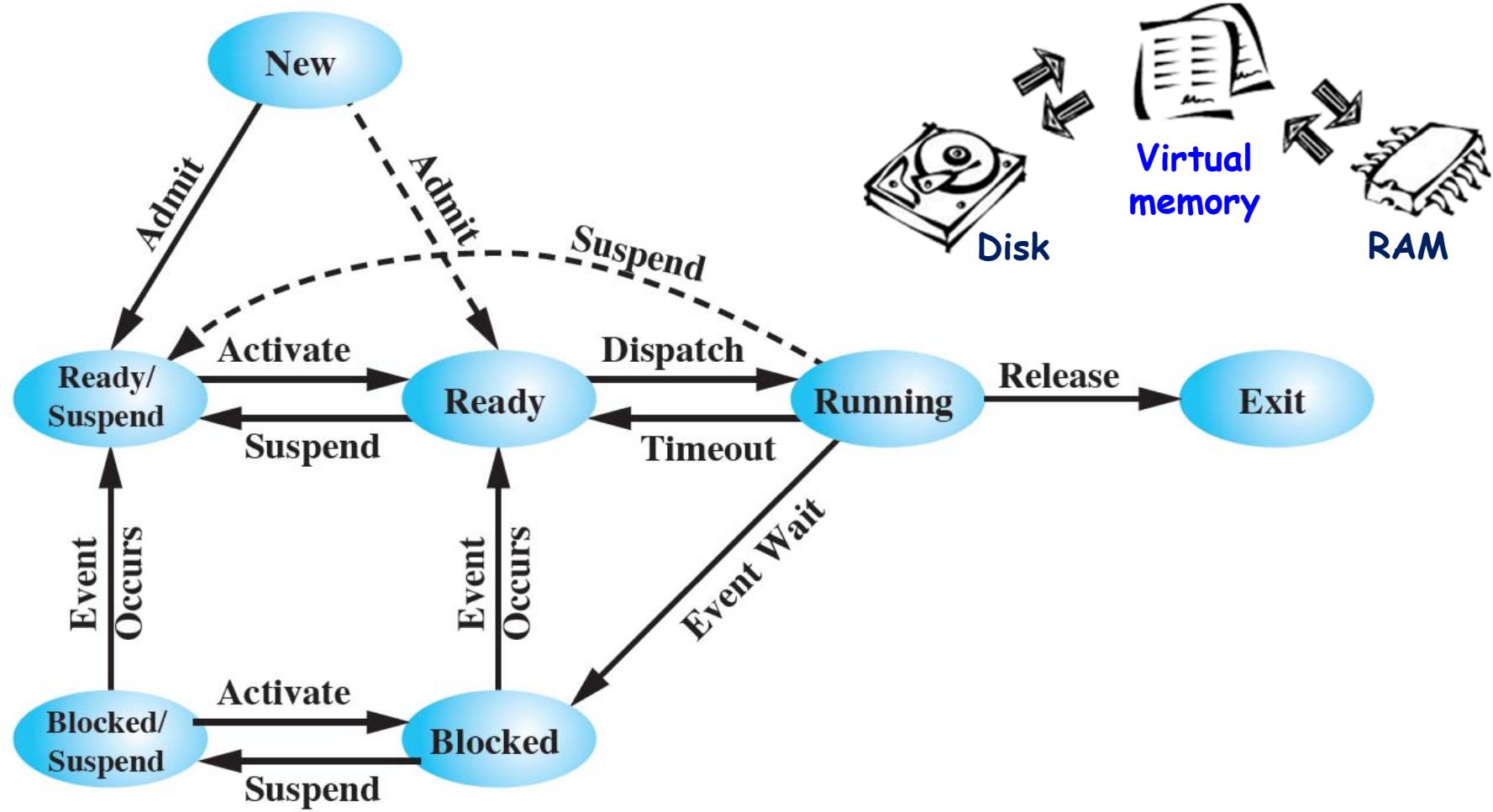
A task is made up of two parts:

- a task execution space
- a task-state segment (TSS)

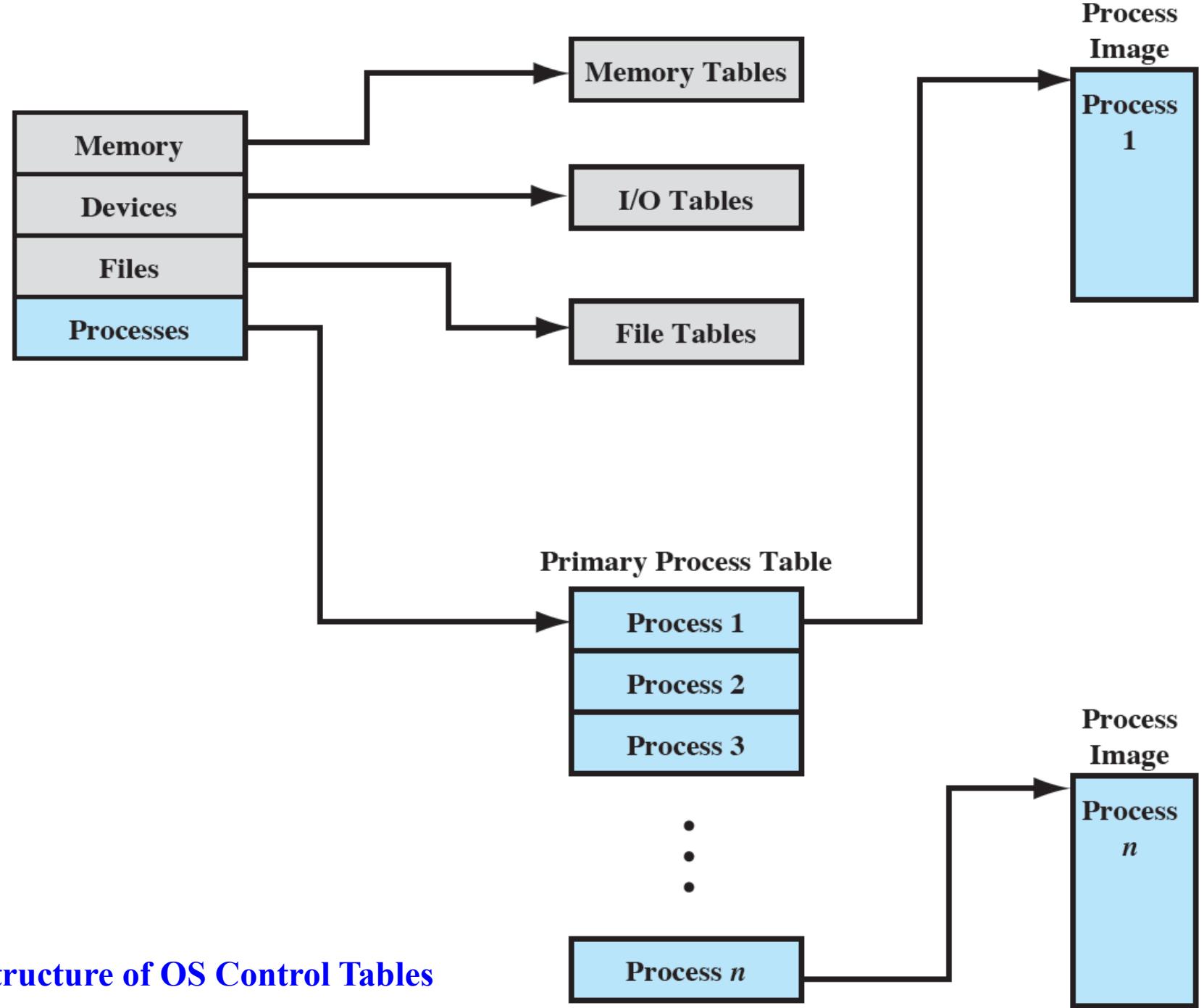


## Task Switching

# Two Suspend States

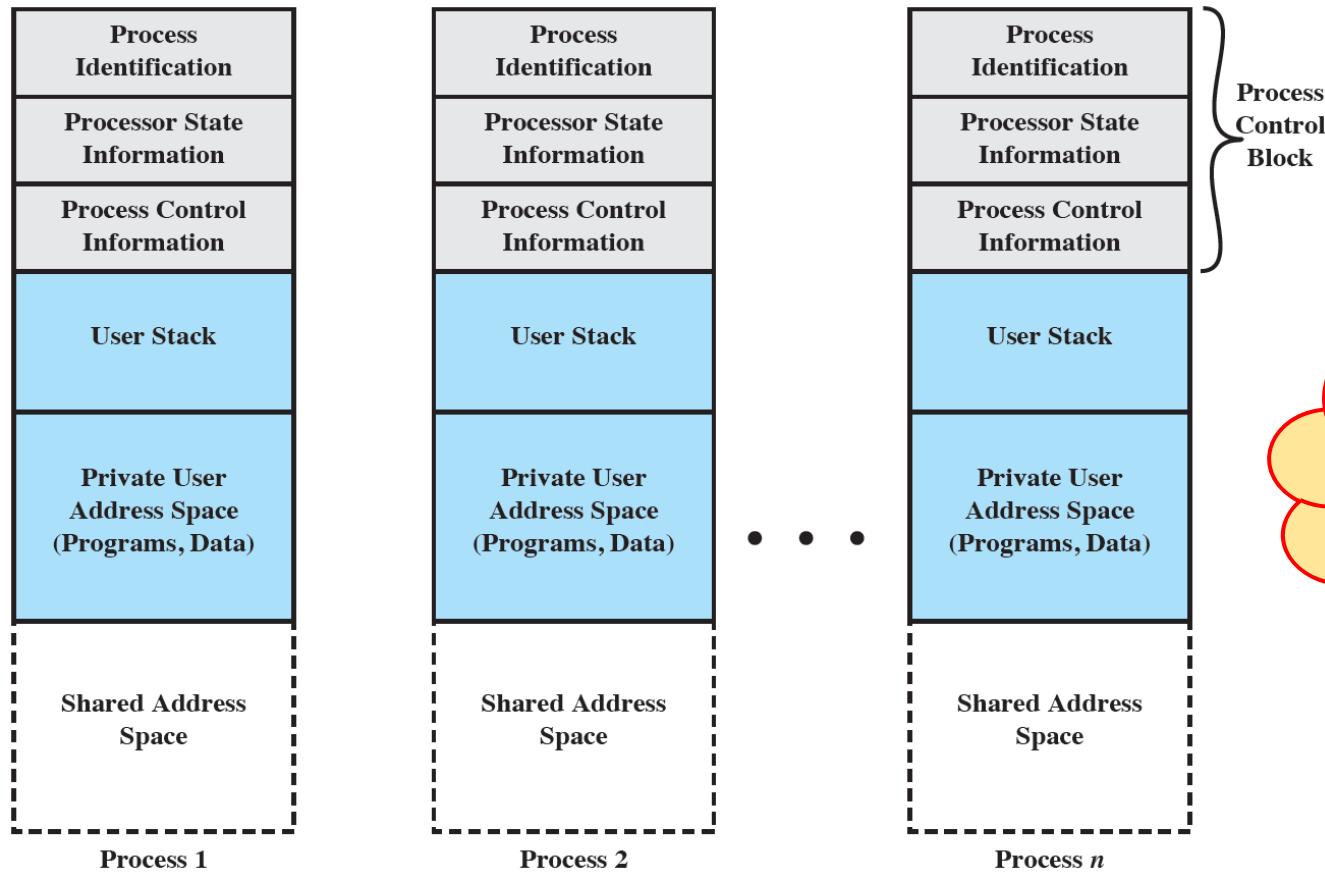


(b) With Two Suspend States



General Structure of OS Control Tables

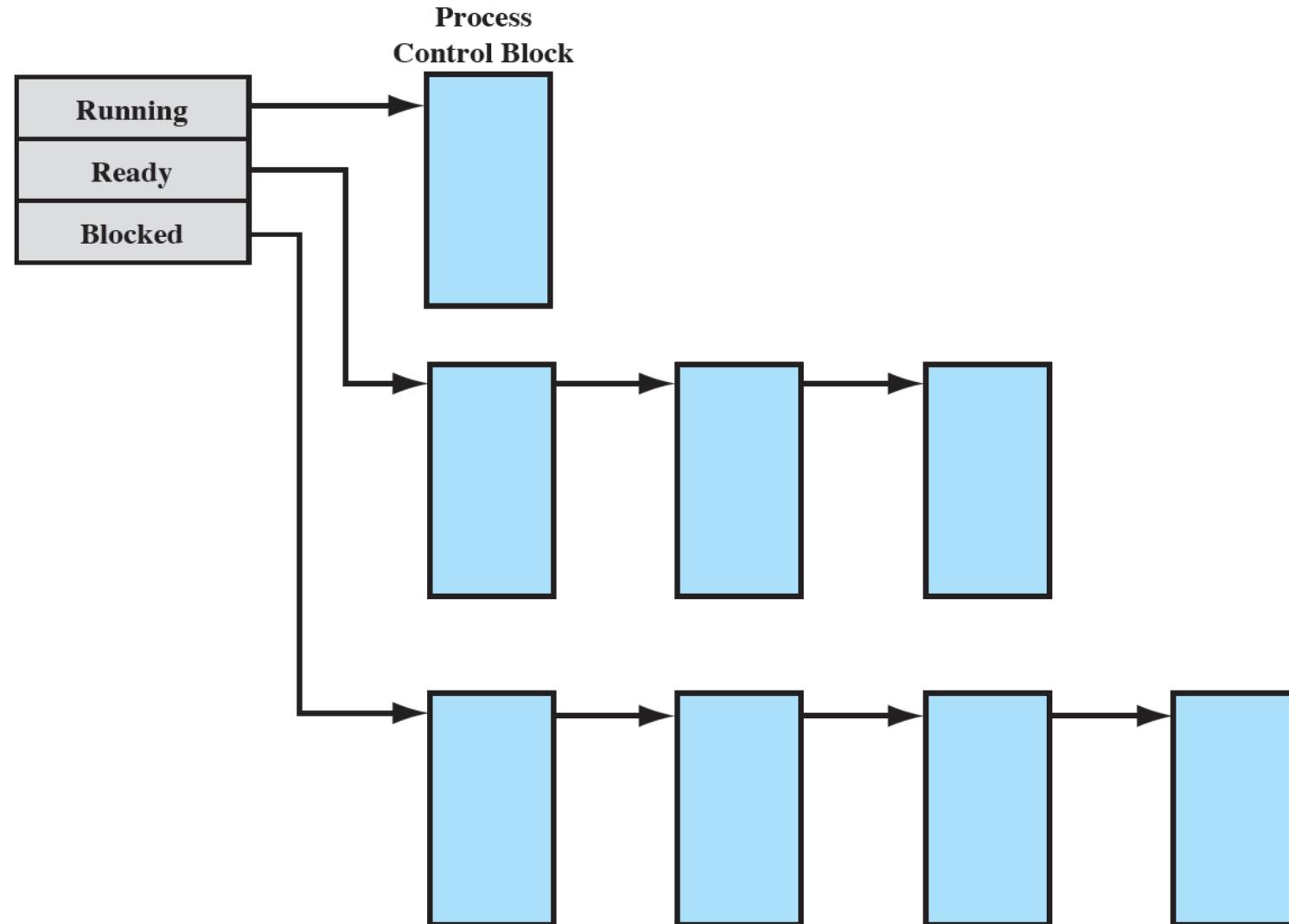
# Structure of Process Images



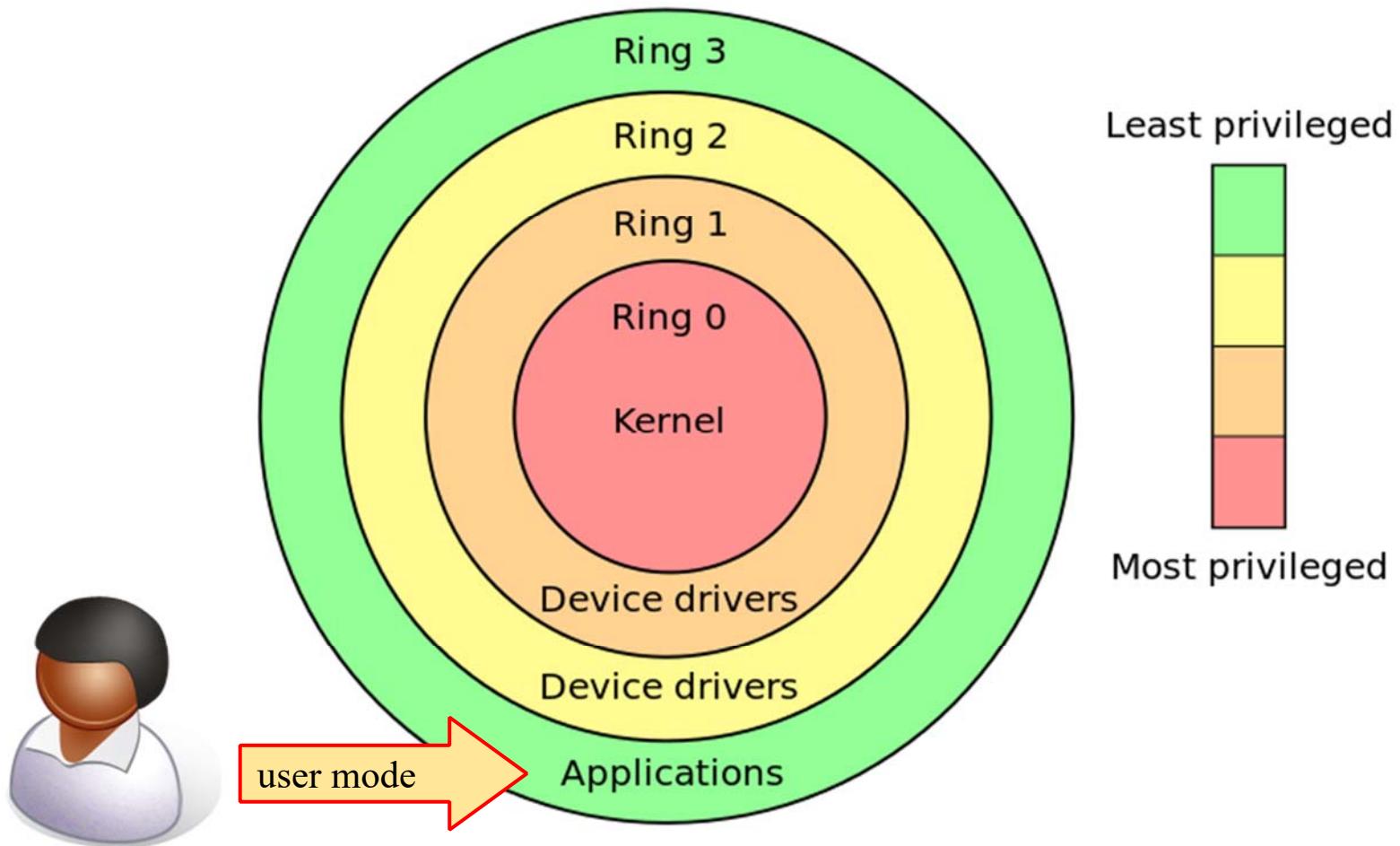
**User Processes in Virtual Memory**

# Process List Structures

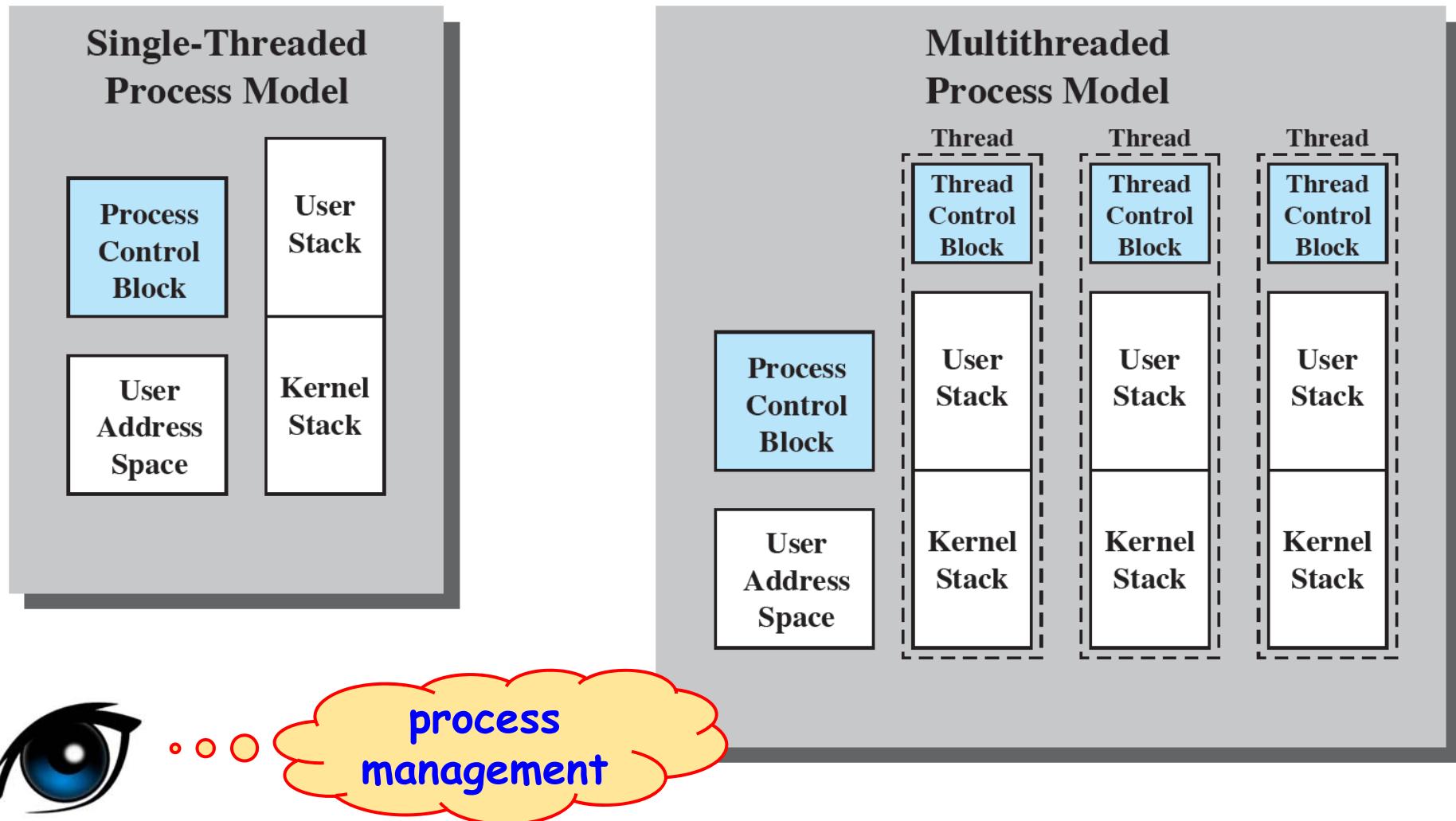
---



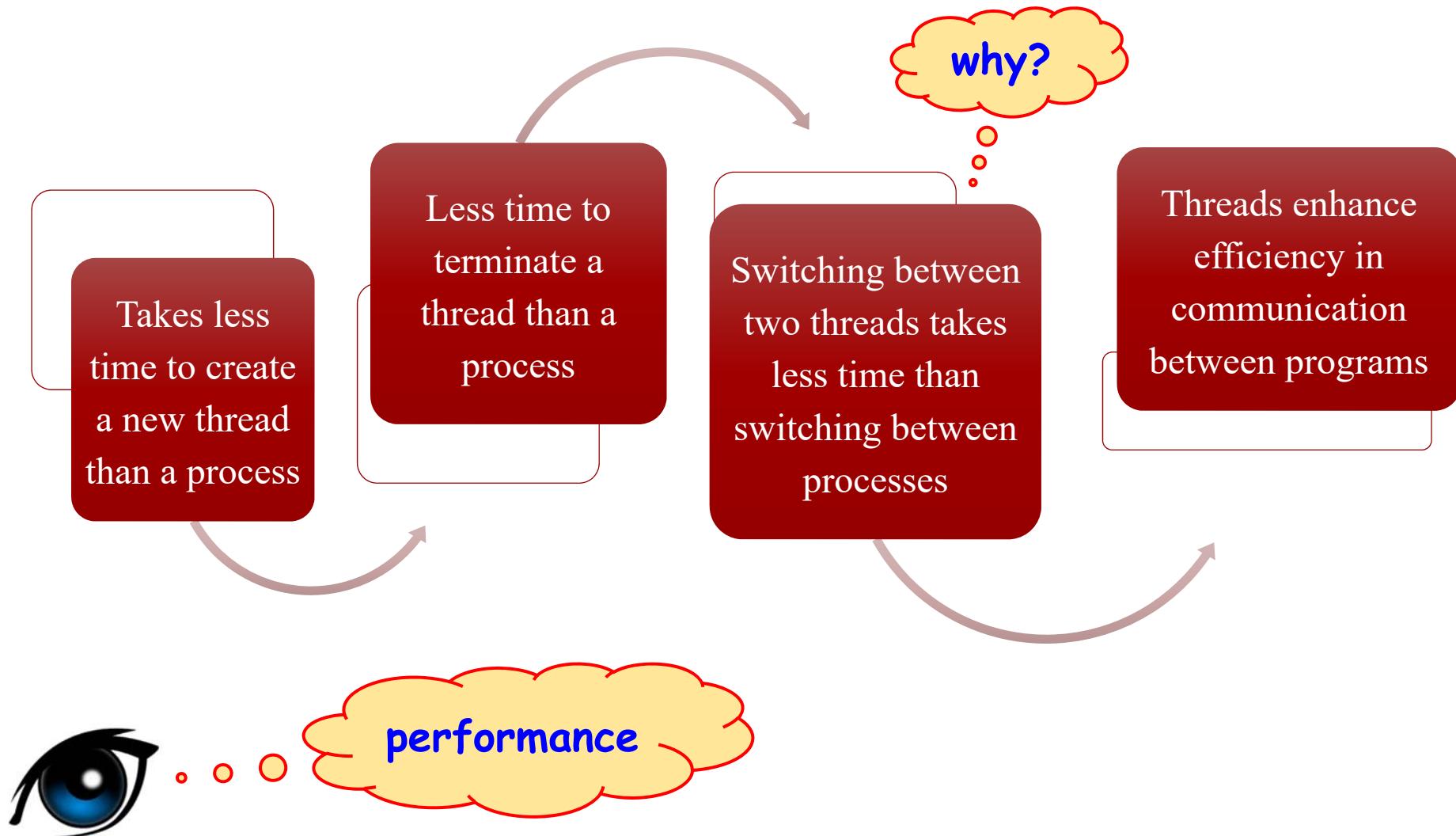
# Modes of Execution



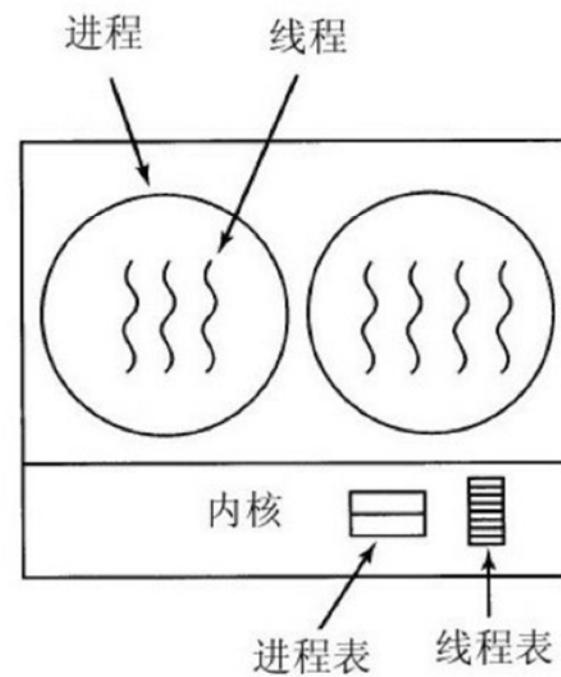
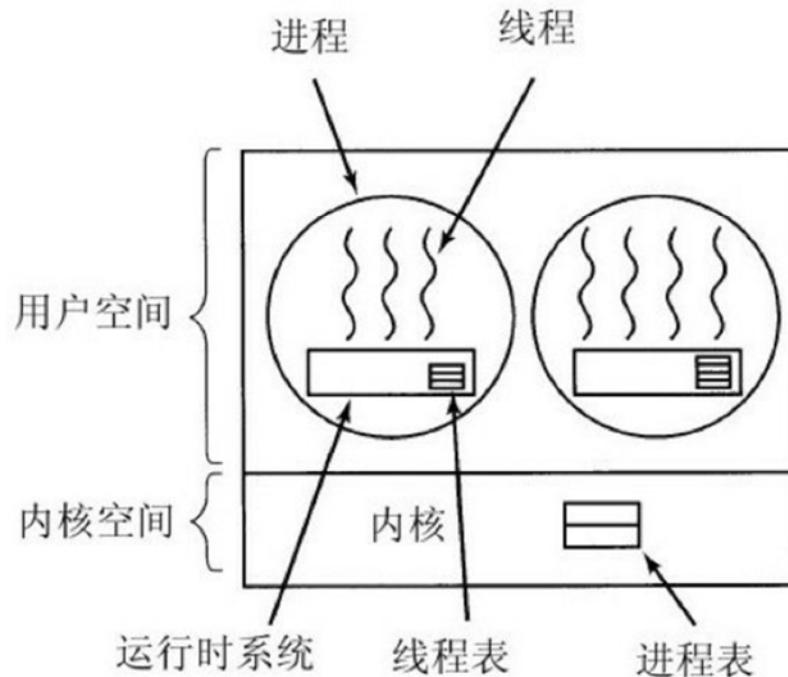
# Threads vs. Processes



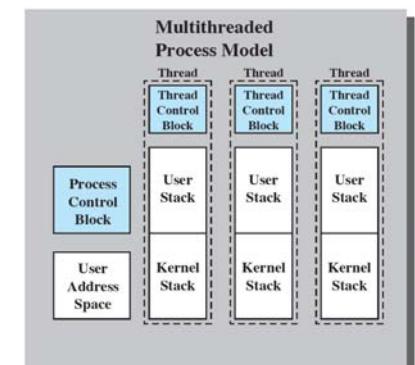
# Key Benefits of Threads



# ULT vs. KLT



线程表位置不同



## A Solution to the Infinite-Buffer Producer/Consumer Problem Using Semaphores

```
/* program producerconsumer */
semaphore n = 0, s = 1;
void producer()
{
    while (true) {
        produce();
        semWait(s);
        append();
        semSignal(s);
        semSignal(n); Red arrow here
    }
}
void consumer()
{
    while (true) {
        semWait(n); Red arrow here
        semWait(s);
        take();
        semSignal(s);
        consume();
    }
}
void main()
{
    parbegin (producer, consumer);
}
```

Which one  
produces a  
serious flaw ?

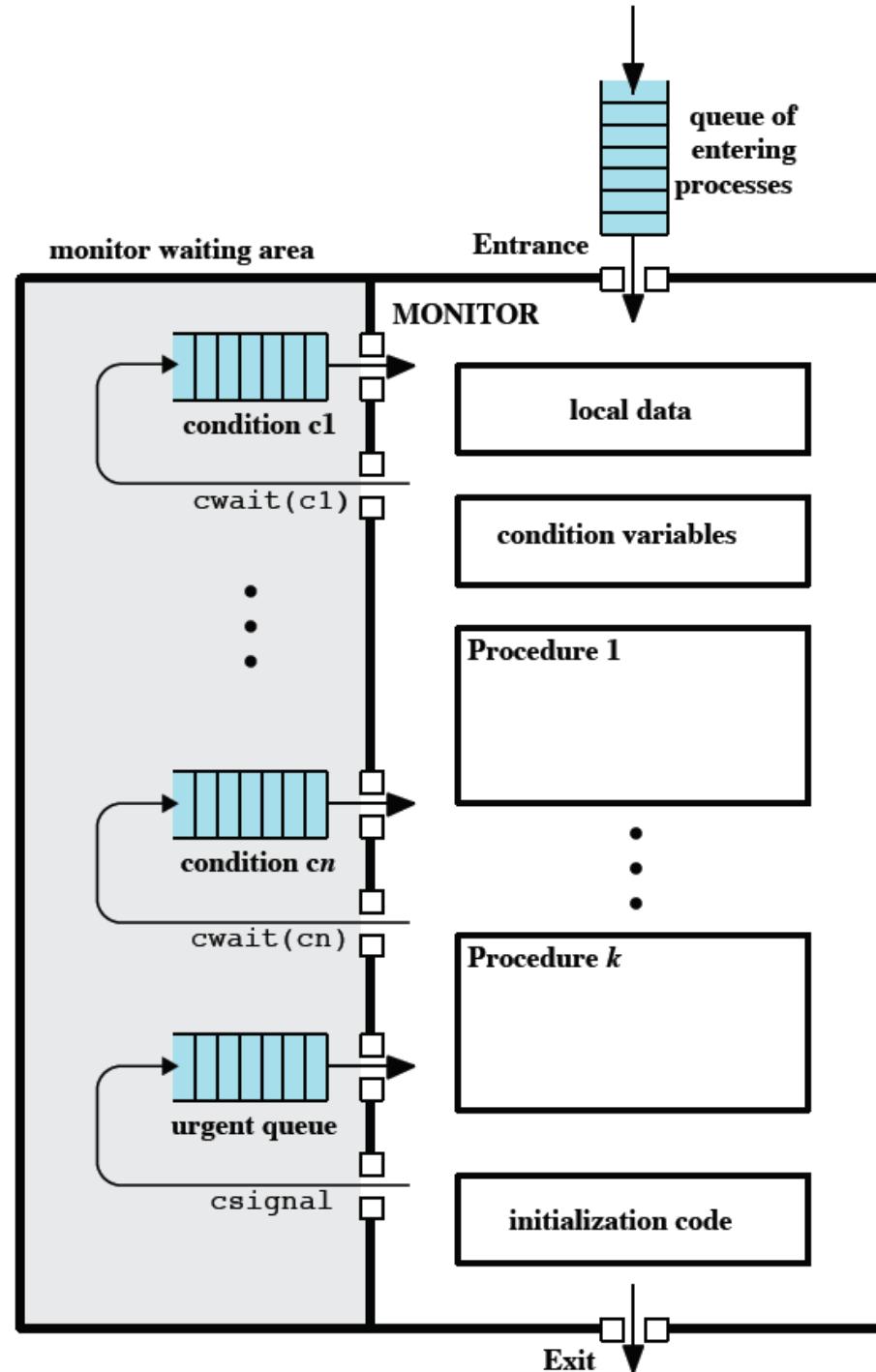


# 信号量总结

---

- 信号量总是成对出现的
- 互斥信号量在一个进程内，同步信号量分散在不同的进程
- 同步信号量 `semWait` 操作在互斥信号量之前，`semSignal` 顺序不重要
- 不要把同步信号量加在互斥信号量之间，否则容易死锁
- 要求会写伪代码

# Structure of a Monitor



# Deadlock

---

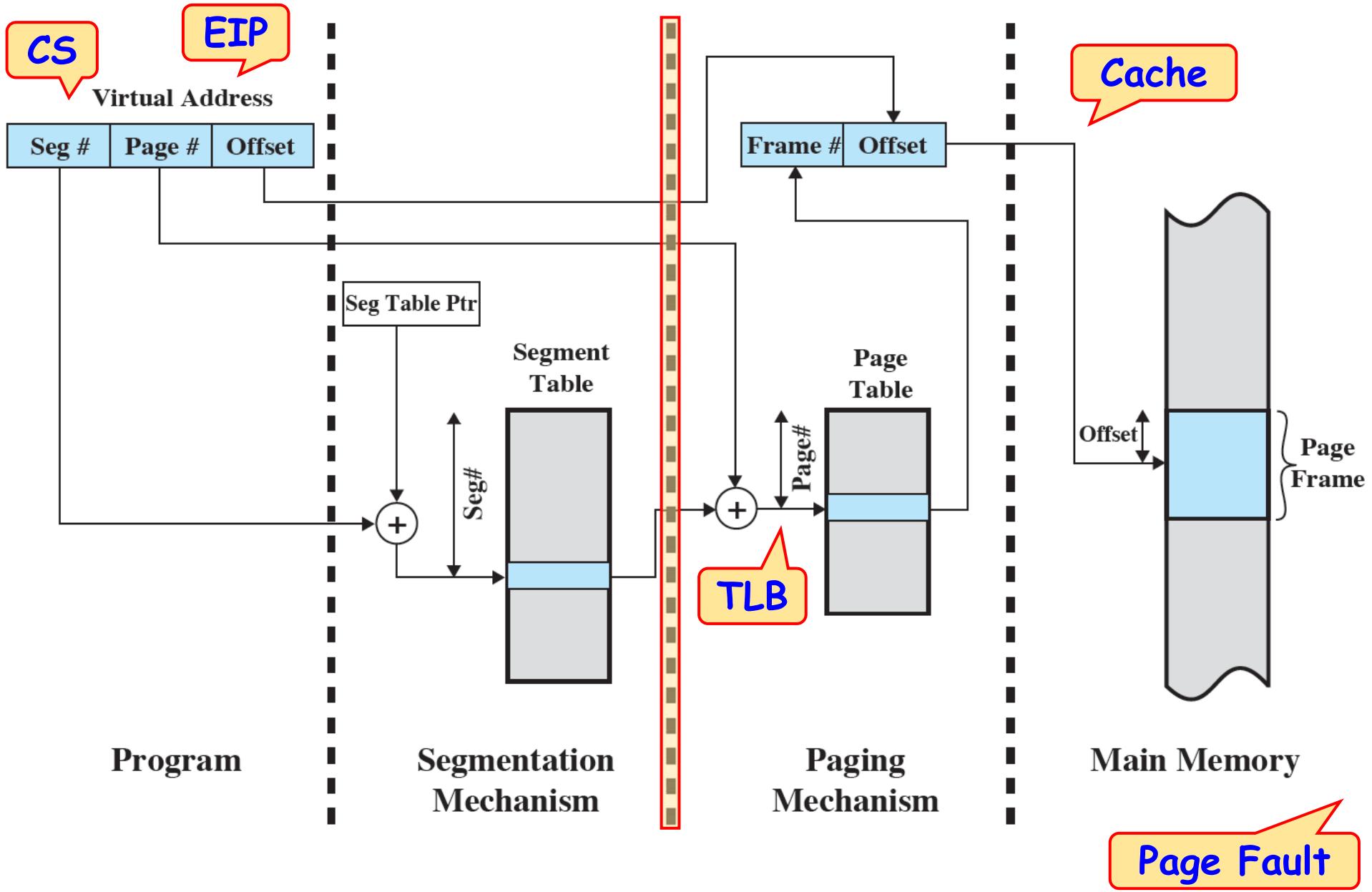


(a) 遇到危险的鸵鸟



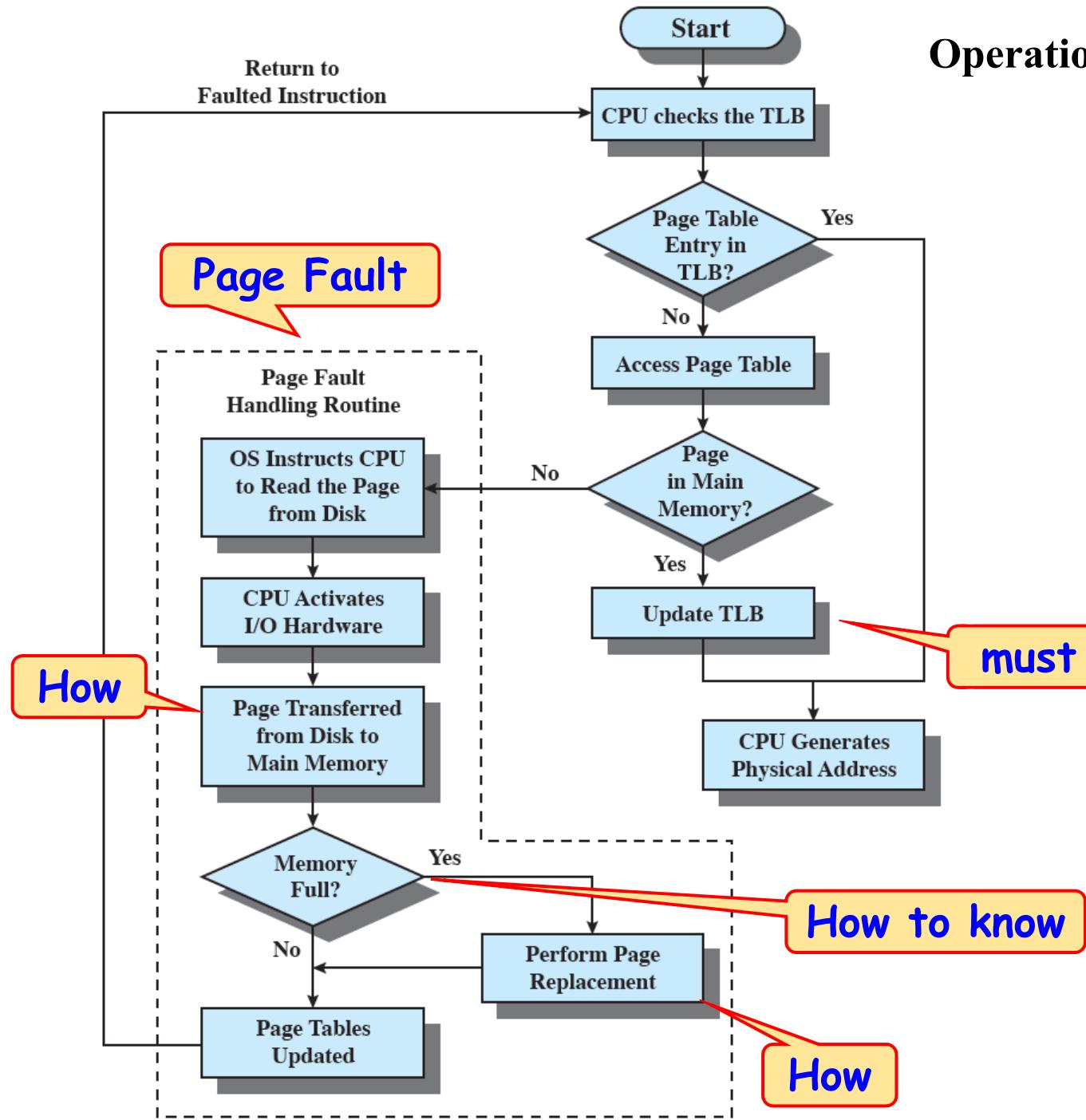
(b) 假装什么都没看见

图 6.1 鸵鸟算法

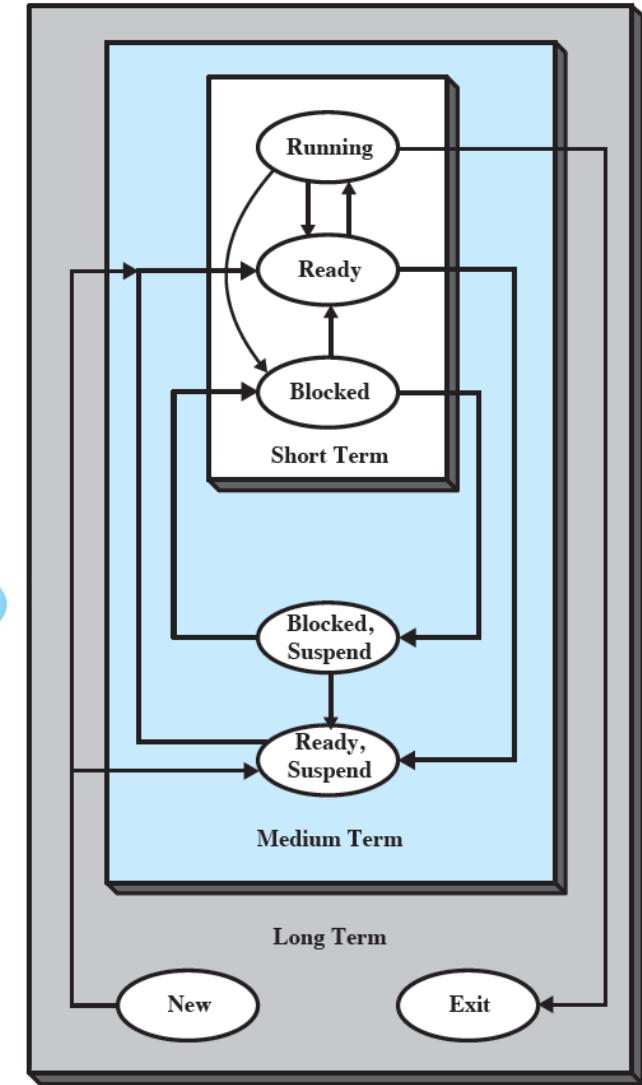
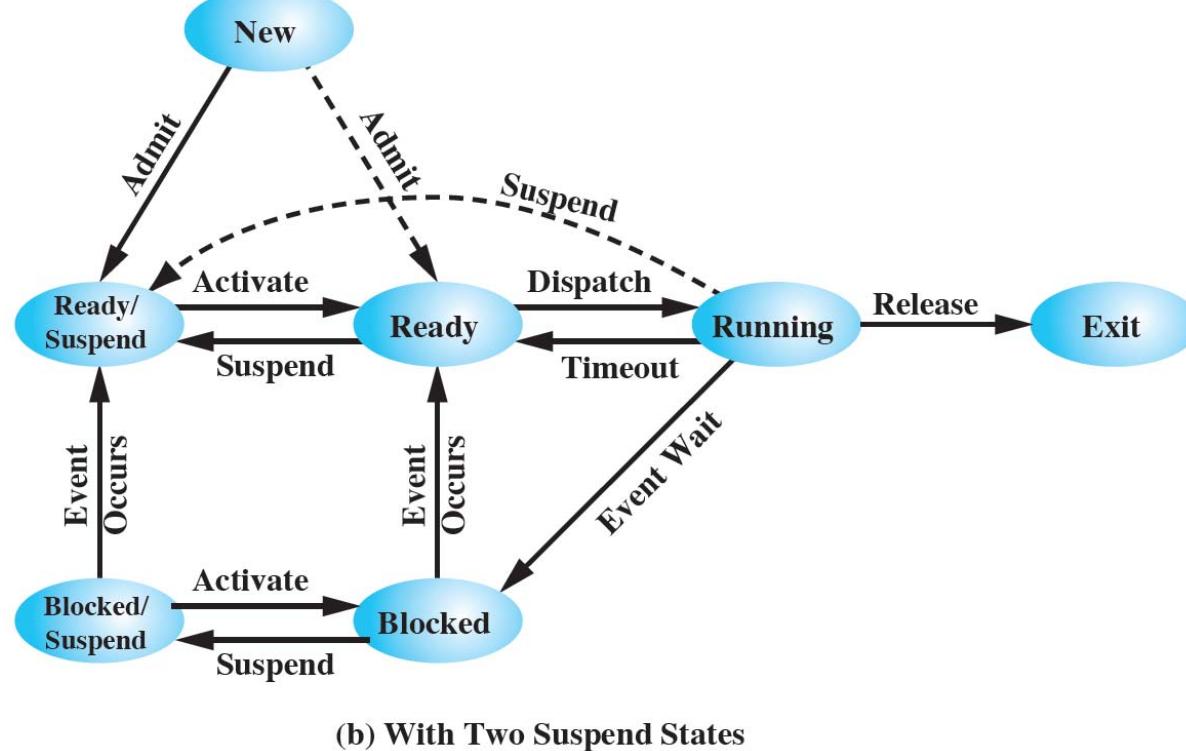


## Address Translation in a Segmentation/Paging System

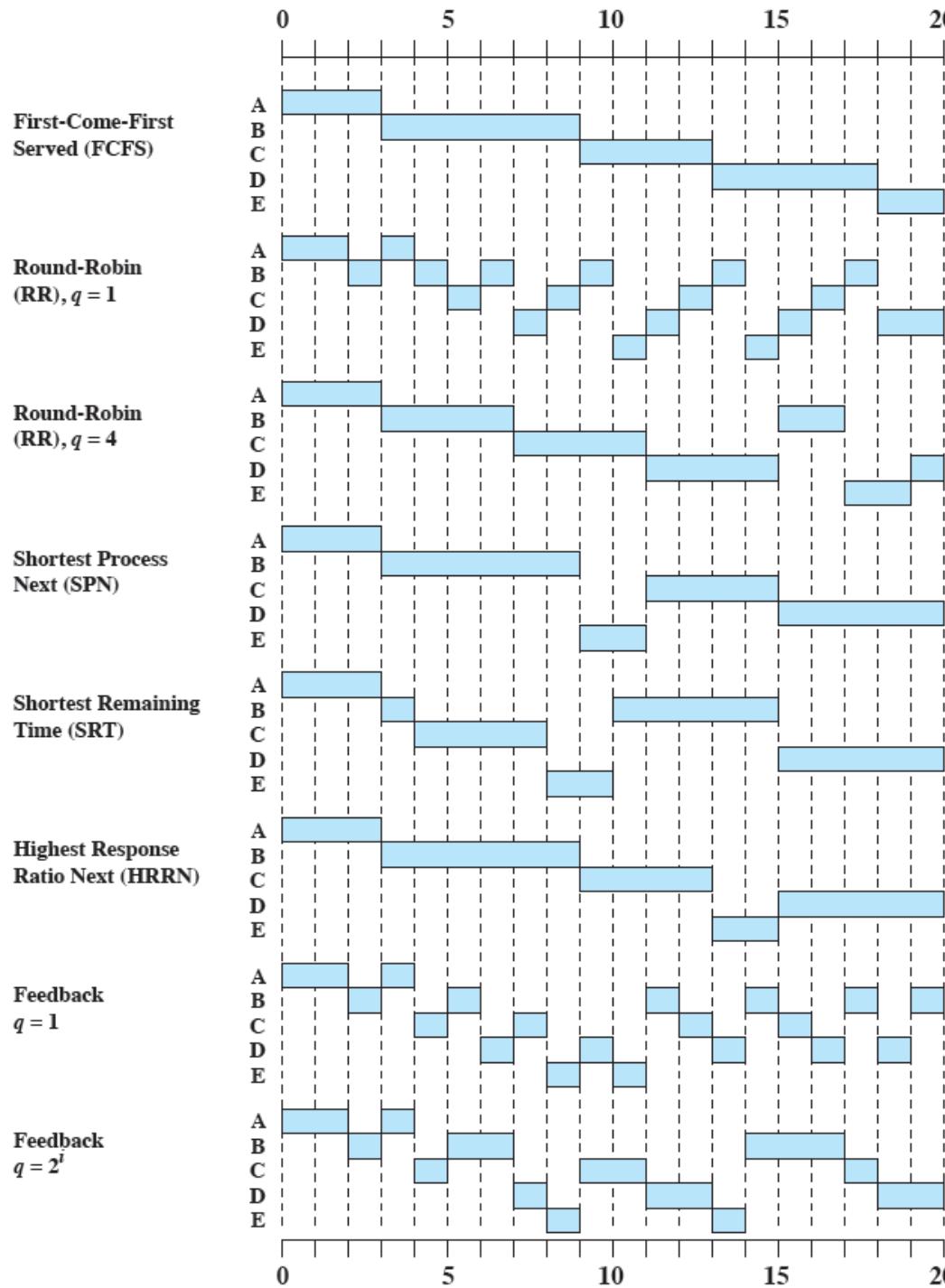
# Operation of Paging and TLB



# Nesting of Scheduling Functions



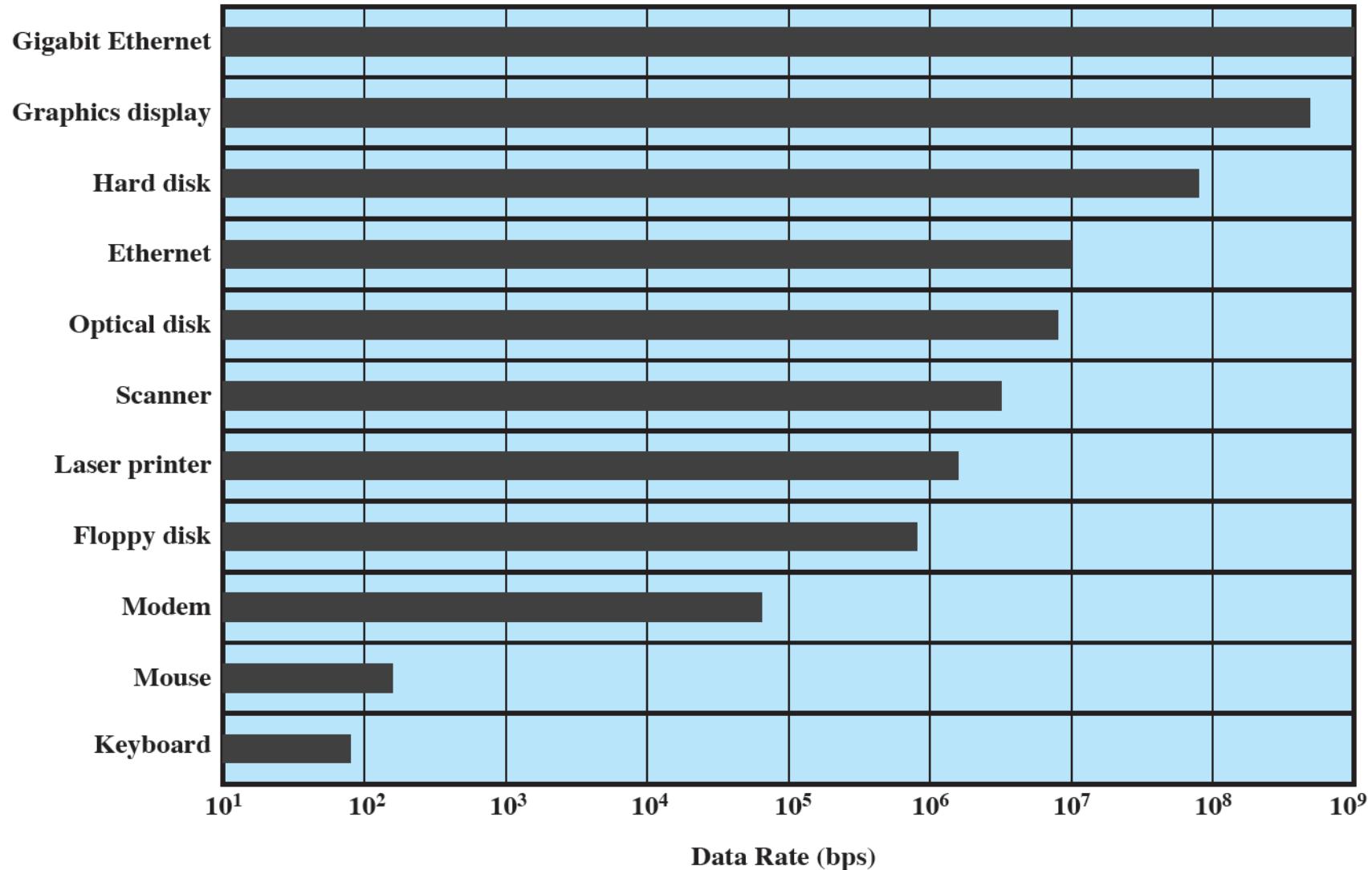
# A Comparison of Scheduling Policies



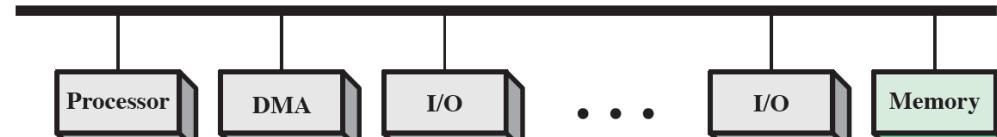
## A Comparison of Scheduling Policies

Process	A	B	C	D	E	
Arrival Time	0	2	4	6	8	
Service Time (Ts)	3	6	4	5	2	Mean
<b>FCFS</b>						
Finish Time	3	9	13	18	20	
Turnaround Time (Tr)	3	7	9	12	12	8.60
Tr/Ts	1.00	1.17	2.25	2.40	6.00	2.56
<b>RR q = 1</b>						
Finish Time	4	18	17	20	15	
Turnaround Time (Tr)	4	16	13	14	7	10.80
Tr/Ts	1.33	2.67	3.25	2.80	3.50	2.71
<b>RR q = 4</b>						
Finish Time	3	17	11	20	19	
Turnaround Time (Tr)	3	15	7	14	11	10.00
Tr/Ts	1.00	2.5	1.75	2.80	5.50	2.71
<b>SPN</b>						
Finish Time	3	9	15	20	11	
Turnaround Time (Tr)	3	7	11	14	3	7.60
Tr/Ts	1.00	1.17	2.75	2.80	1.50	1.84
<b>SRT</b>						
Finish Time	3	15	8	20	10	
Turnaround Time (Tr)	3	13	4	14	2	7.20
Tr/Ts	1.00	2.17	1.00	2.80	1.00	1.59
<b>HRRN</b>						
Finish Time	3	9	13	20	15	
Turnaround Time (Tr)	3	7	9	14	7	8.00
Tr/Ts	1.00	1.17	2.25	2.80	3.5	2.14
<b>FB q = 1</b>						
Finish Time	4	20	16	19	11	
Turnaround Time (Tr)	4	18	12	13	3	10.00
Tr/Ts	1.33	3.00	3.00	2.60	1.5	2.29

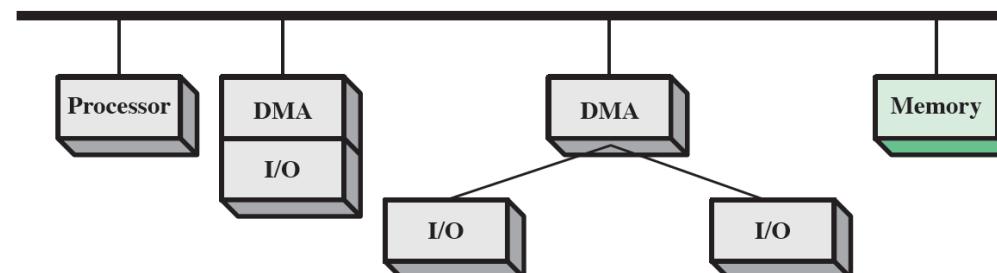
# Typical I/O Device Data Rates



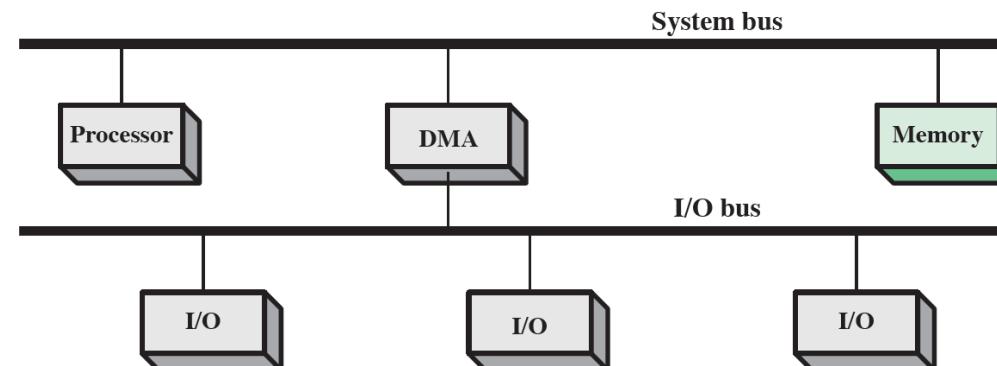
# Alternative DMA Configurations



(a) Single-bus, detached DMA

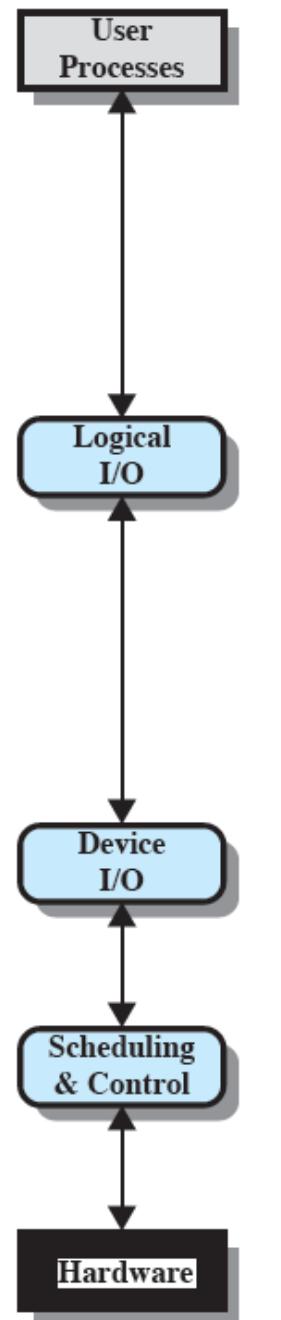


(b) Single-bus, Integrated DMA-I/O

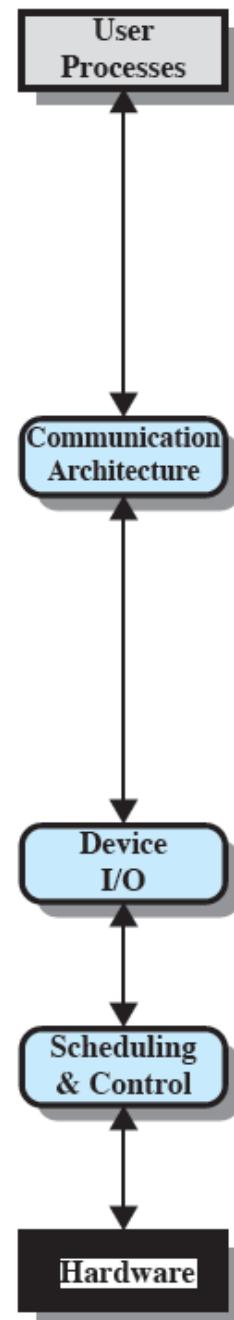


(c) I/O bus

# A Model of I/O Organization



(a) Local peripheral device

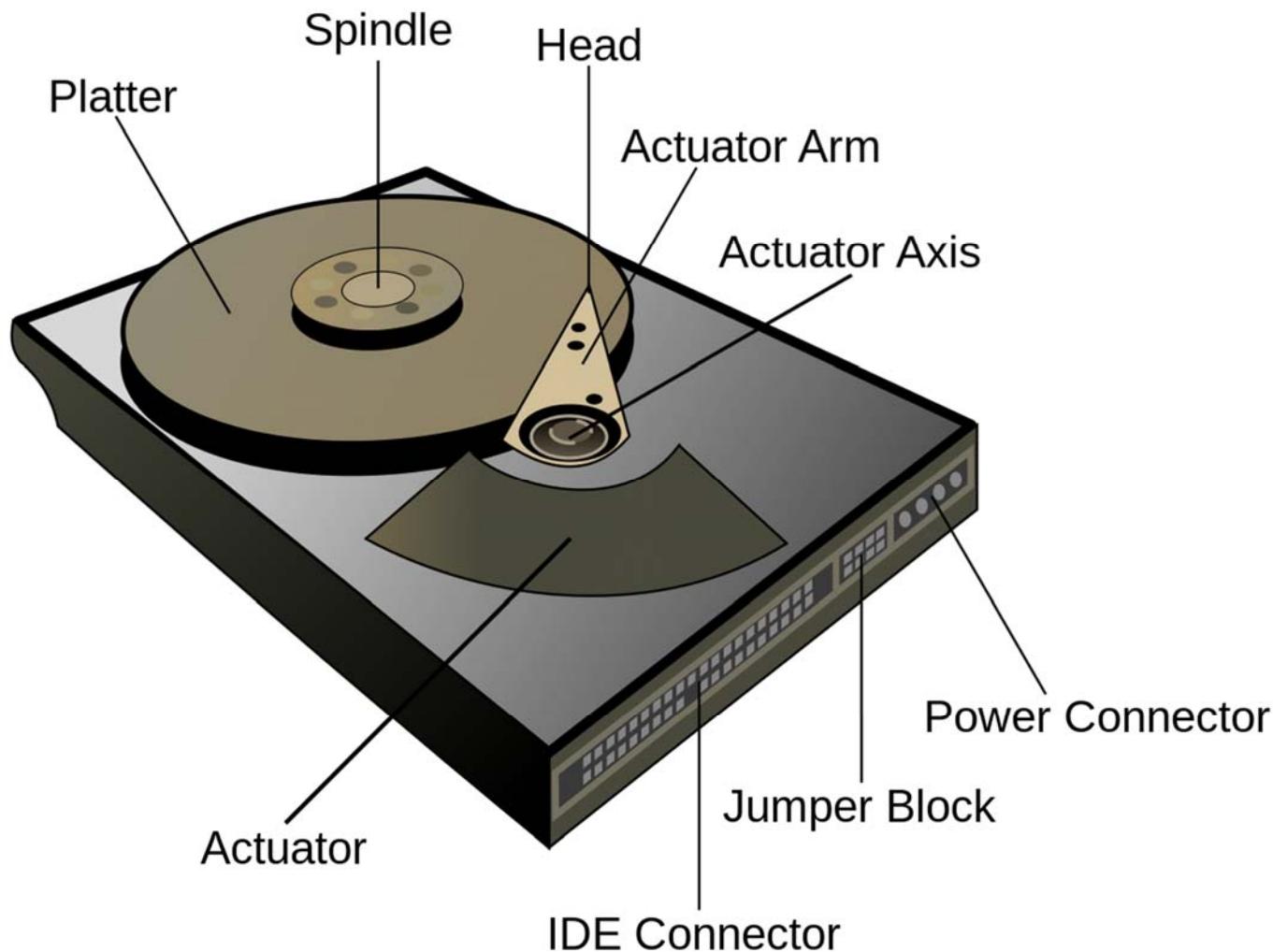


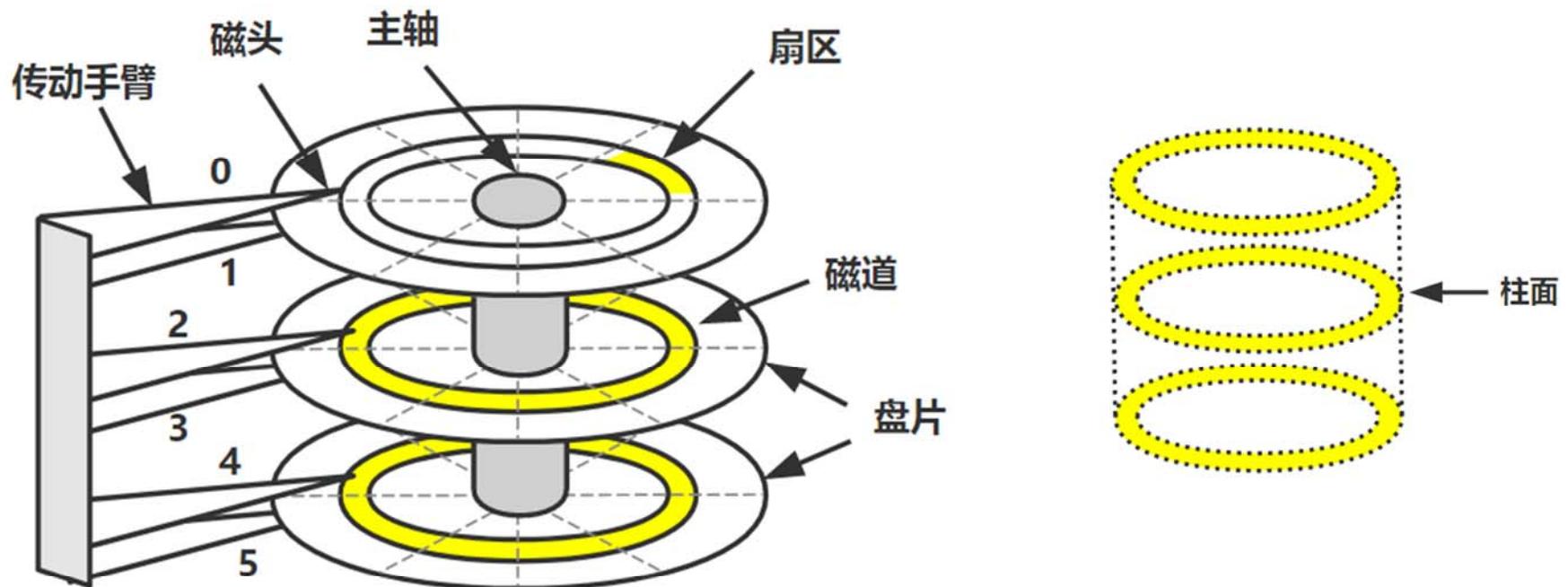
(b) Communications port



(c) File system

# Diagram labeling the major components of HDD





#### 9.3.3.4 扇区号与柱面号、当前磁道扇区号和当前磁头号的对应关系

假定硬盘的每磁道扇区数是 `track_secs`, 硬盘磁头总数是 `dev_heads`, 磁道总数是 `tracks`。对于指定的硬盘顺序扇区号是 `sector`, 对应的当前柱面号是 `cyl`, 在当前磁道上的扇区号是 `sec`, 当前磁头号是 `head`。那么若想从指定顺序扇区号 `sector` 换算成对应的当前柱面号、当前磁道上扇区号以及当前磁头号, 则可以使用以下步骤:

- $\text{sector} / \text{track\_secs} =$  整数是 `tracks`, 余数是 `sec`;
- $\text{tracks} / \text{dev\_heads} =$  整数是 `cyl`, 余数是 `head`;
- 在当前磁道上扇区号从 1 算起, 于是需要把 `sec` 增 1。

若想从指定的当前 `cyl`、`sec` 和 `head` 换算成从硬盘开始算起的顺序扇区号, 则过程正好与上述相反。换算公式和上面给出的完全一样, 即:

---

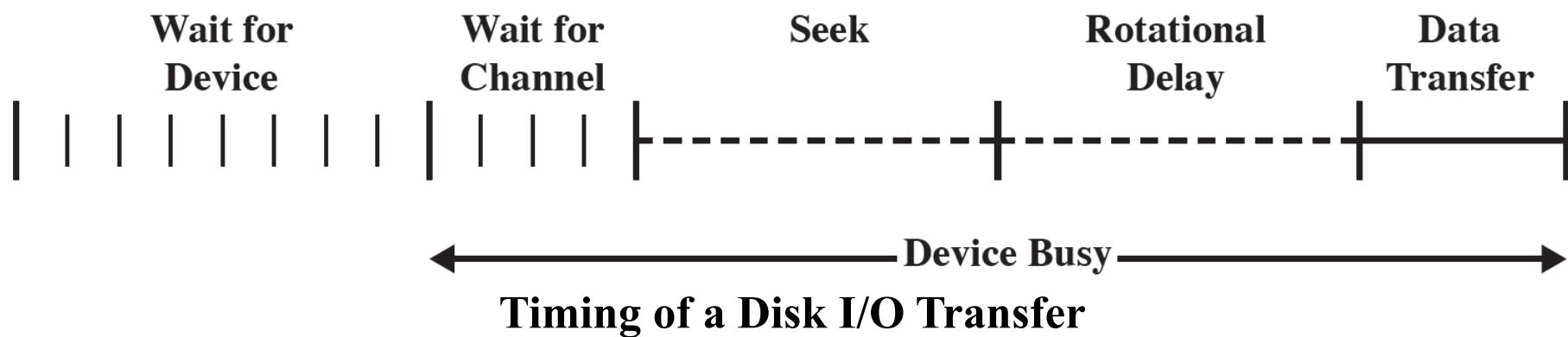

$$\text{sector} = (\text{cyl} * \text{dev\_heads} + \text{head}) * \text{track\_secs} + \text{sec} - 1$$


---

# Disk Performance Parameters

- The actual details of disk I/O operation depend on

- computer system
- operating system
- nature of the I/O channel and disk controller hardware



$T$  = transfer time,

$b$  = number of bytes to be transferred,

$N$  = number of bytes on a track, and

$r$  = rotation speed, in revolutions per second.

$$T_a = T_s + \frac{1}{2r} + \frac{b}{rN}$$

# 举例：一个典型的磁盘

- 平均寻道时间 4 ms
- 转速 7500 rpm
- 每个磁道有 500 个扇区、每个扇区有 512 个字节
- 假设读取一个包含 2500 个扇区，大小 1.28M 的文件
- 估计传送需要的总时间

$$T_a = T_s + \frac{1}{2r} + \frac{b}{rN}$$

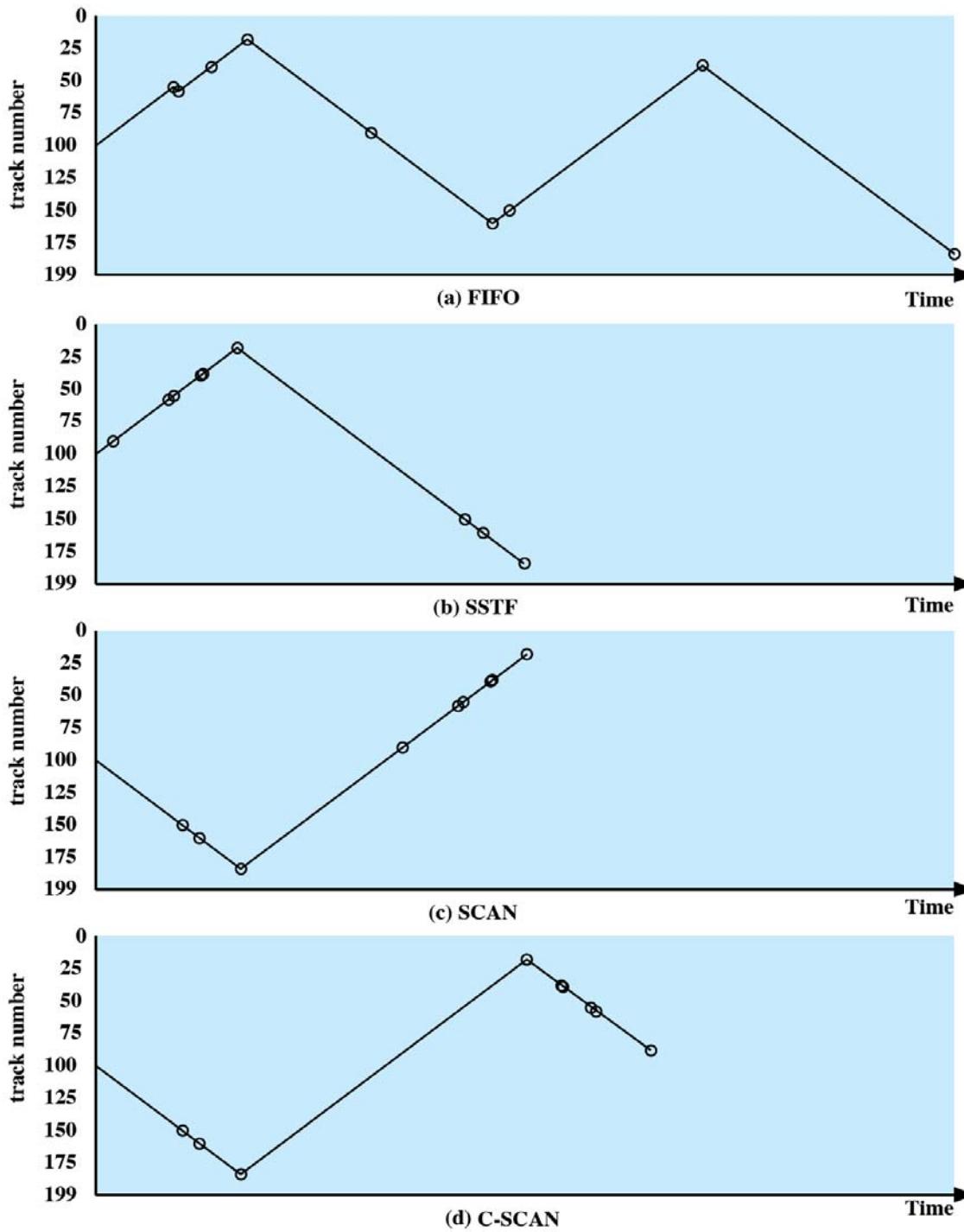
Average seek	4 ms
Rotational delay	4 ms
Read 500 sectors	<u>8 ms</u>
	16 ms

Average seek	4 ms
Rotational delay	4 ms
Read 1 sector	<u>0.016 ms</u>
	8.016 ms

Total time =  $16 + (4 \times 12) = 64$  ms = 0.064 seconds

Total time =  $2,500 \times 8.016 = 20,040$  ms = 20.04 seconds

# Comparison of Disk Scheduling Algorithms



# Comparison of Disk Scheduling Algorithms

(a) FIFO (starting at track 100)		(b) SSTF (starting at track 100)		(c) SCAN (starting at track 100, in the direction of increasing track number)		(d) C-SCAN (starting at track 100, in the direction of increasing track number)	
Next track accessed	Number of tracks traversed	Next track accessed	Number of tracks traversed	Next track accessed	Number of tracks traversed	Next track accessed	Number of tracks traversed
55	45	90	10	150	50	150	50
58	3	58	32	160	10	160	10
39	19	55	3	184	24	184	24
18	21	39	16	90	94	18	166
90	72	38	1	58	32	38	20
160	70	18	20	55	3	39	1
150	10	150	132	39	16	55	16
38	112	160	10	38	1	58	3
184	146	184	24	18	20	90	32
Average seek length	55.3	Average seek length	27.5	Average seek length	27.8	Average seek length	35.8

# 硬盘分区表

表 9-10 硬盘主引导扇区 MBR 的结构

偏移位置	名称	长度(字节)	说明
0x000	MBR 代码	446	引导程序代码和数据。
0x1BE	分区表项 1	16	第 1 个分区表项，共 16 字节。
0x1CE	分区表项 2	16	第 2 个分区表项，共 16 字节。
0x1DE	分区表项 3	16	第 3 个分区表项，共 16 字节。
0x1EE	分区表项 4	16	第 4 个分区表项，共 16 字节。
0x1FE	引导标志	2	有效引导扇区的标志，值分别是 0x55, 0xAA。



主引导扇区

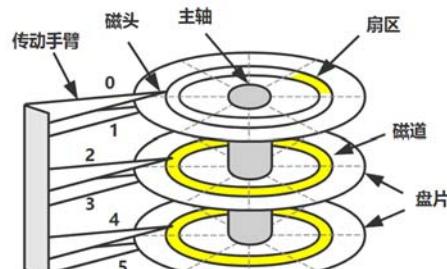
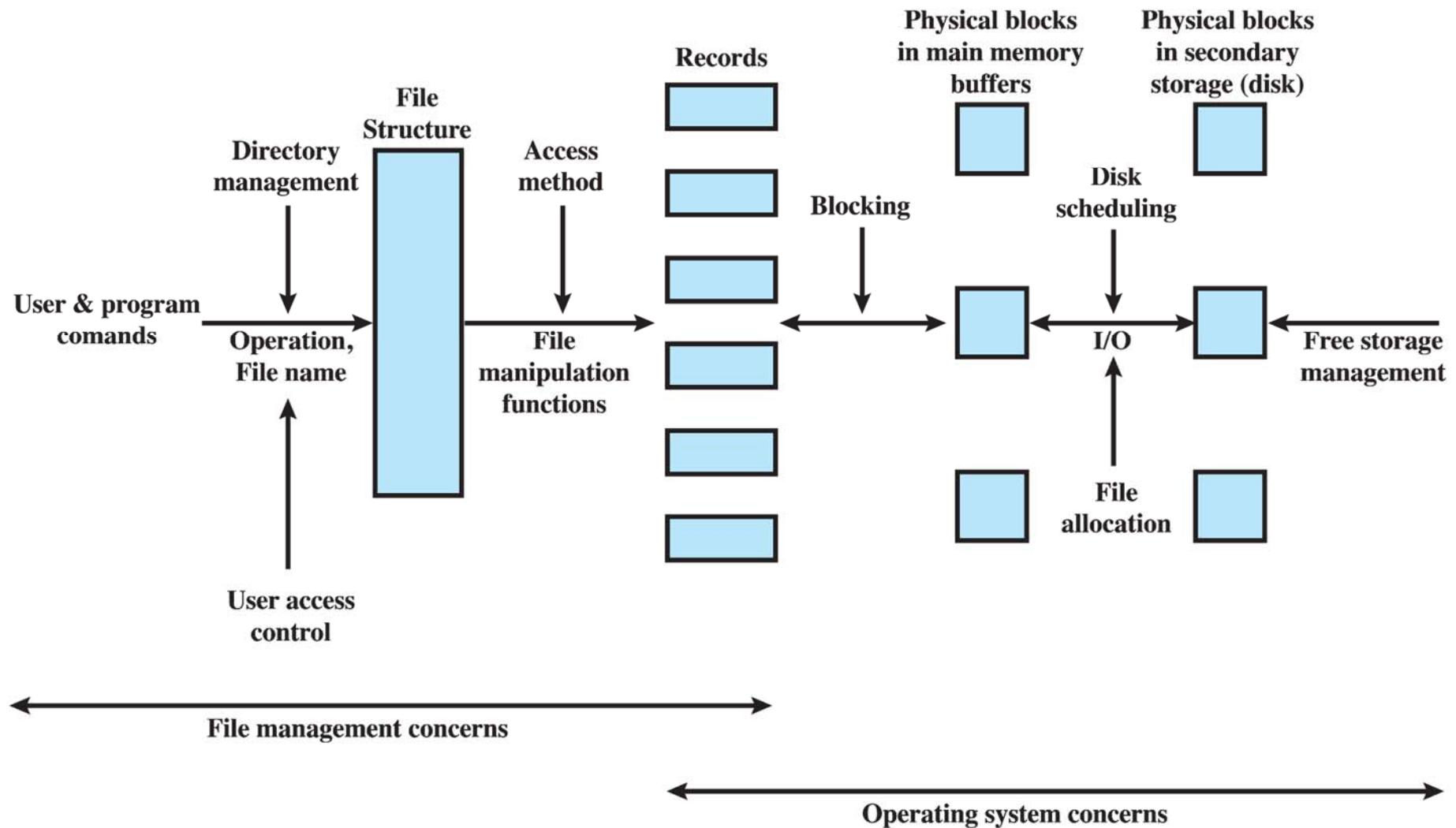


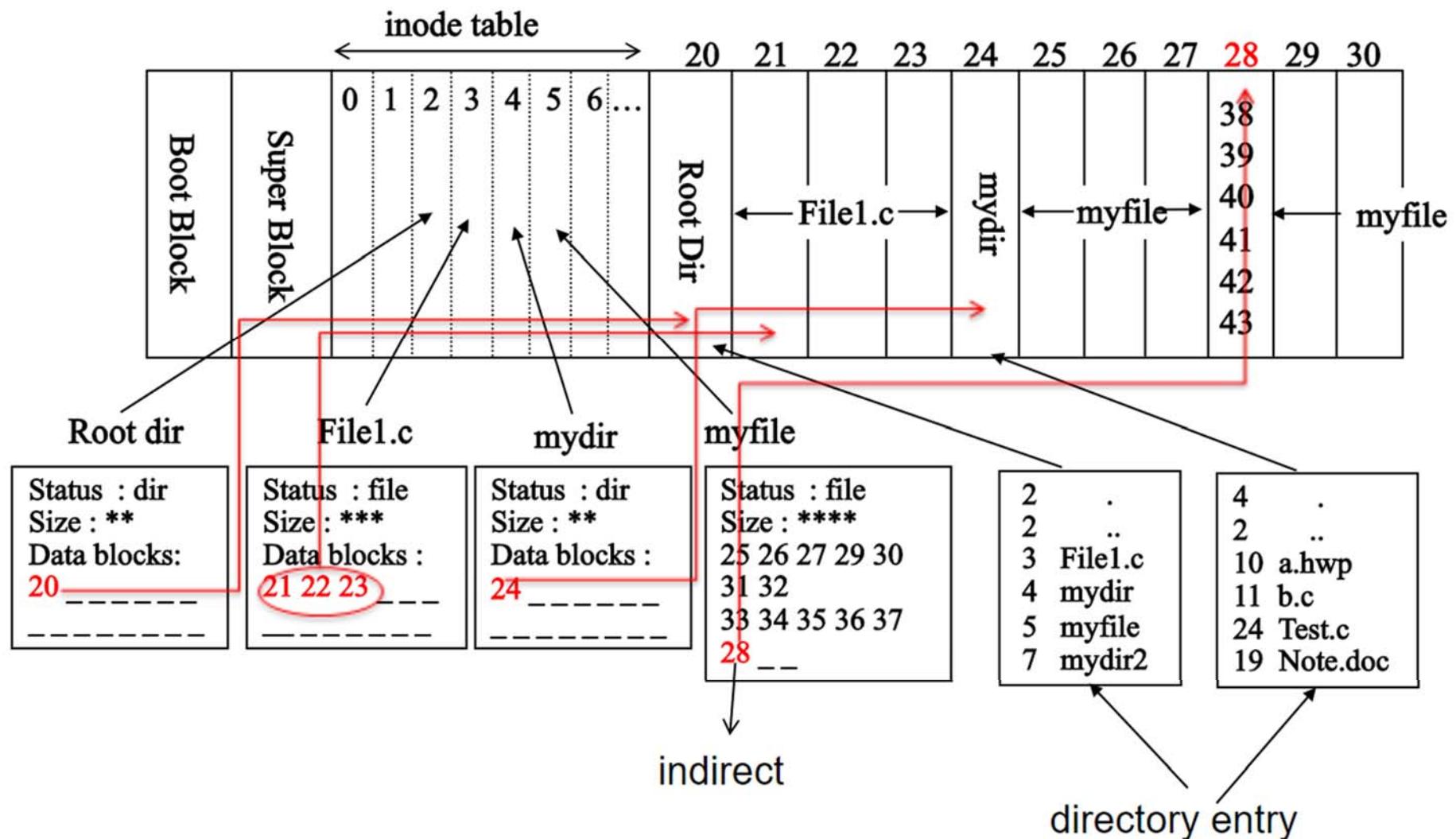
表 9-11 硬盘分区表项结构

位置	名称	大小	说明
0x00	boot_ind	字节	引导标志。4 个分区中同时只能有一个分区是可引导的。 0x00-不从该分区引导操作系统；0x80-从该分区引导操作系统。
0x01	head	字节	分区起始处的磁头号。磁头号范围为 0--255。
0x02	sector	字节	分区起始的当前柱面中扇区号(位 0-5) 和柱面号高 2 位(位 6-7)。
0x03	cyl	字节	分区起始处的柱面号低 8 位。柱面号范围为 0--1023。
0x04	sys_ind	字节	分区类型字节。0xb-DOS; 0x80-Old Minix; 0x83-Linux ...
0x05	end_head	字节	分区结束处的磁头号。磁头号范围为 0--255。
0x06	end_sector	字节	分区结束的当前柱面中扇区号(位 0-5) 和柱面号高 2 位(位 6-7)。
0x07	end_cyl	字节	分区结束的柱面号低 8 位。柱面号范围为 0--1023。
0x08-0x0b	start_sect	长字	分区起始的物理扇区号。它以整个硬盘的扇区号顺序从 0 计起。
0x0c-0x0f	nr_sects	长字	分区占用的扇区数。

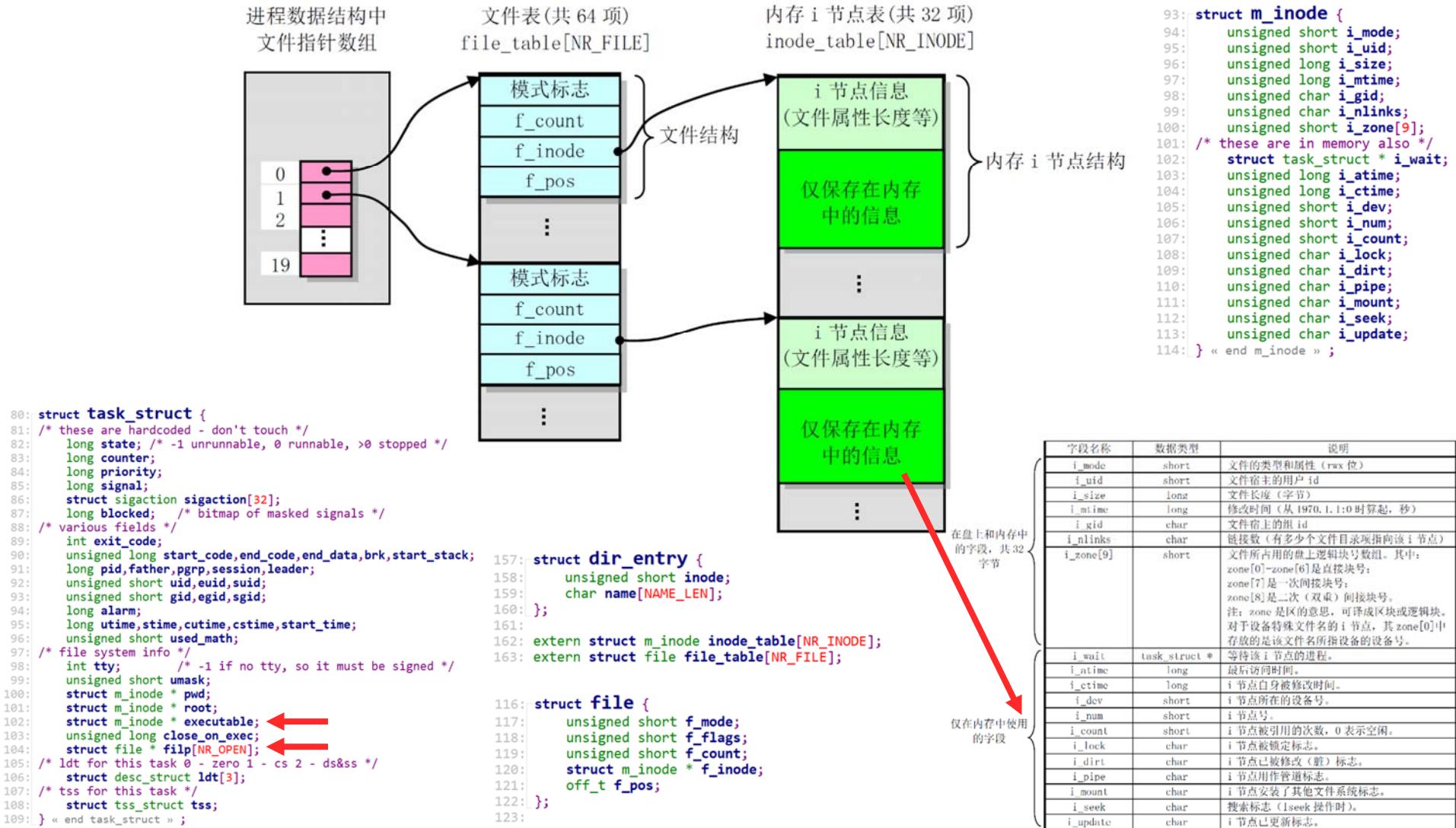
# Elements of File Management

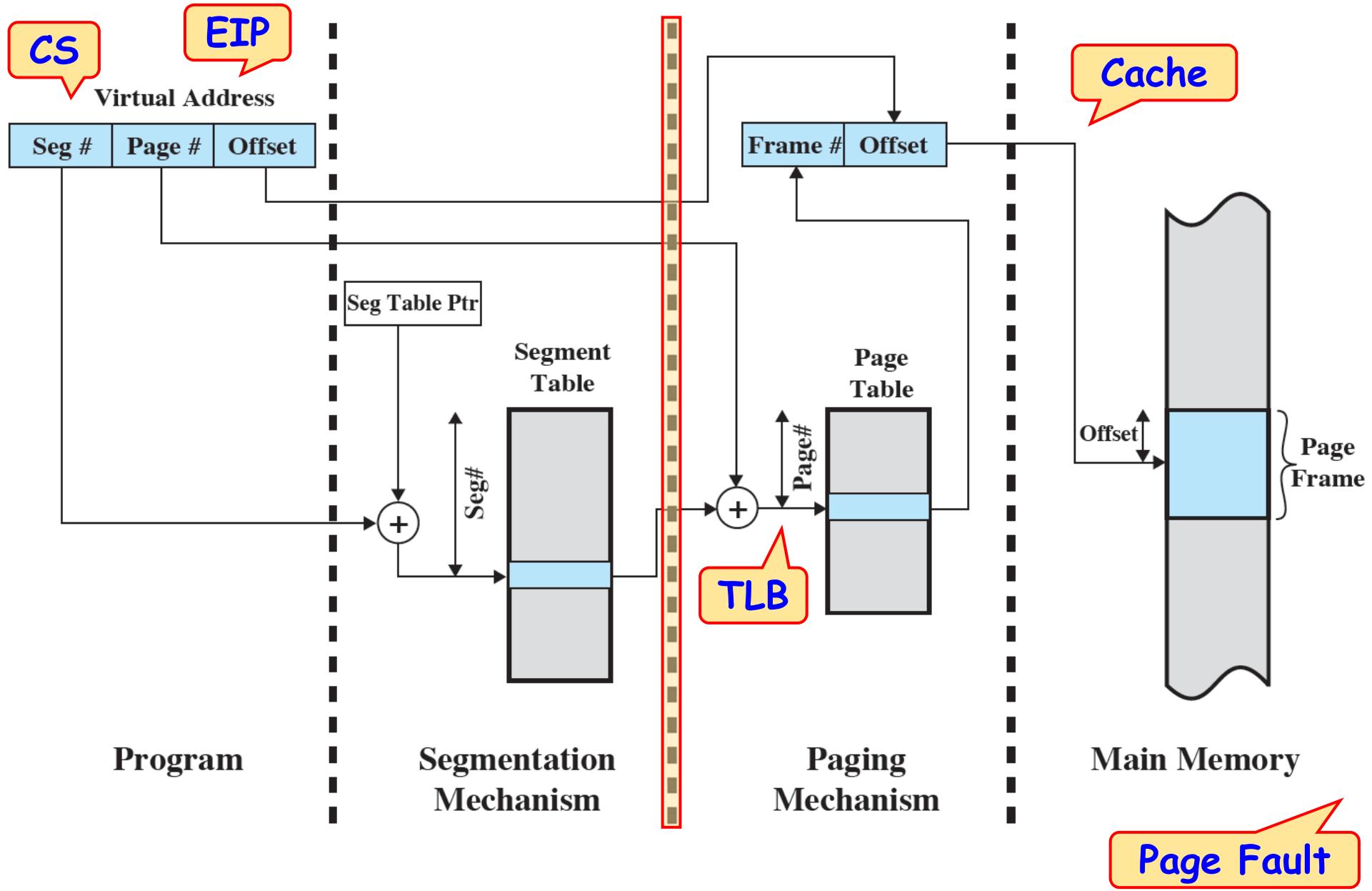


# Demo



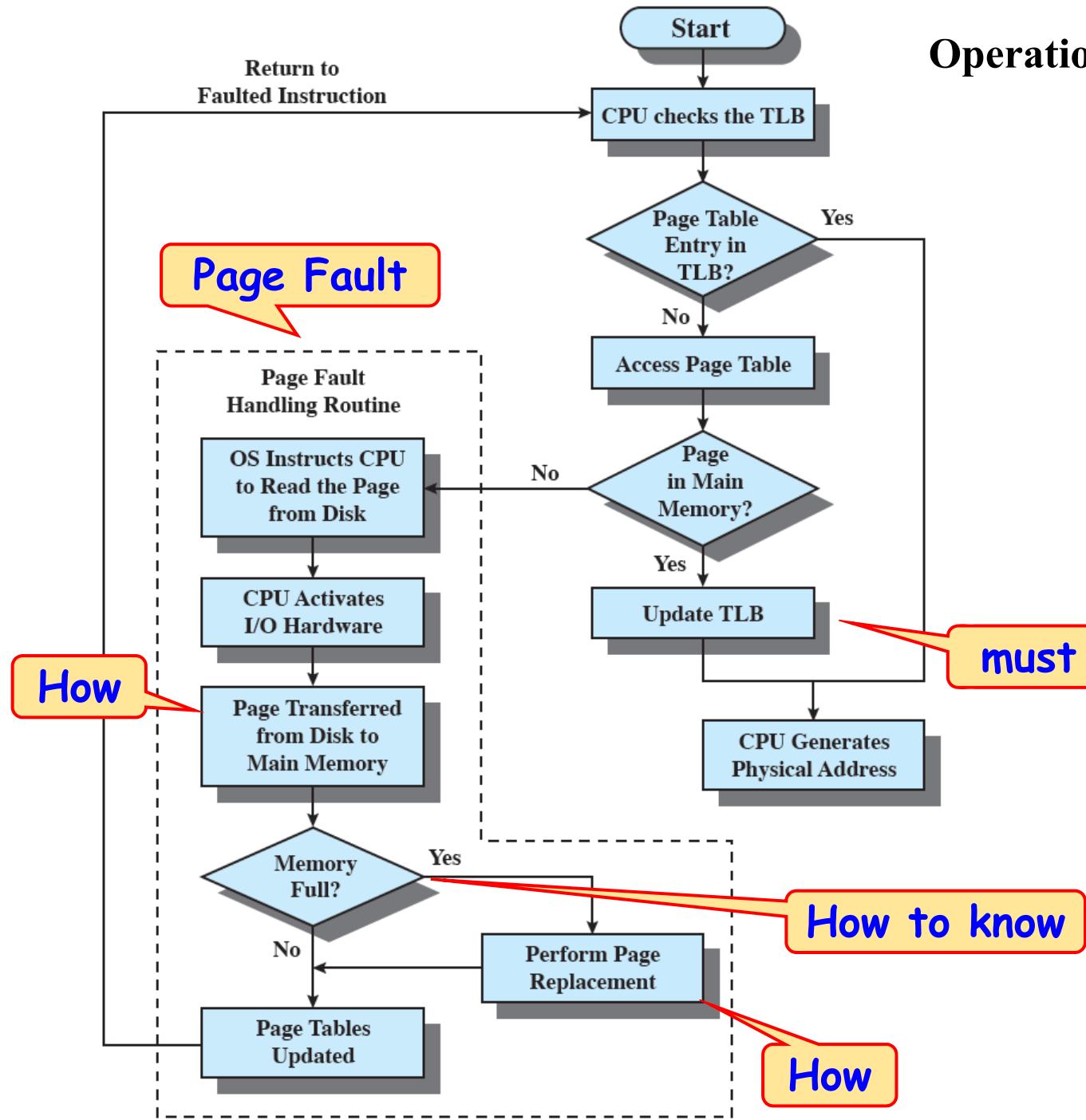
# 进程打开文件使用的内核数据结构



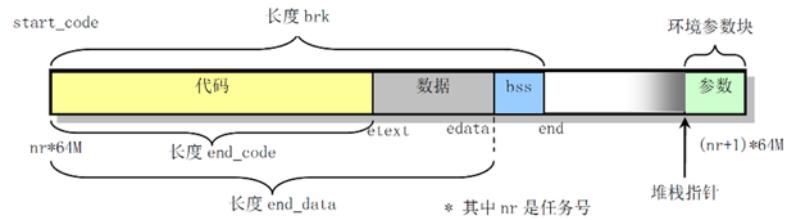


## Address Translation in a Segmentation/Paging System

# Operation of Paging and TLB



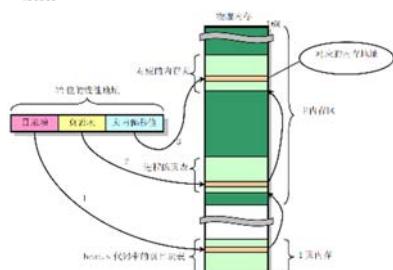
# Page fault



```

366: void do_no_page(unsigned long error_code,unsigned long address)
367: {
368:     int nr[4];
369:     unsigned long tmp;
370:     unsigned long page;
371:     int block,i;
372:
373:     address &= 0xfffff000;
374:     tmp = address - current->start_code;
375:     if (!current->executable || tmp >= current->end_data) {
376:         get_empty_page(address);
377:         return;
378:     }
379:     if (share_page(tmp))
380:         return;
381:     if (!(page = get_free_page()))
382:         oom();
383:     /* remember that 1 block is used for header */
384:     block = 1 + tmp/BLOCK_SIZE;
385:     for (i=0 ; i<4 ; block++,i++)
386:         nr[i] = bmap(current->executable,block);
387:     bread_page(page,current->executable->i_dev,nr);
388:     i = tmp + 4096 - current->end_data;
389:     tmp = page + 4096;
390:     while (i-- > 0) {
391:         tmp--;
392:         *(char *)tmp = 0;
393:     }
394:     if (put_page(page,address)
395:         return;
396:     free_page(page);
397:     oom();
398: } « end do_no_page »

```



```

page_fault:
    xchgl %eax,(%esp)
    pushl %ecx
    pushl %edx
    pushl %ds
    pushl %es
    pushl %fs
    movl $0x10,%edx
    movl %dx,%ds
    movl %dx,%es
    movl %dx,%fs
    movl %cr2,%edx
    pushl %edx
    pushl %eax
    testl $1,%eax
    jne 1f
    call do_no_page
    jmp 2f
    call do_wp_page
    addl $8,%esp
    popl %fs
    popl %es
    popl %ds
    popl %edx
    popl %ecx
    popl %eax
    iret
1:
2:

```

```

140: int bmap(struct m_inode * inode,int block)
141: {
142:     return _bmap(inode,block,0);
143: }

```

```

72: static int _bmap(struct m_inode * inode,int block,int create)
73: {
74:     struct buffer_head * bh;
75:     int i;
76:
77:     if (block<0)
78:         panic("_bmap: block<0");
79:     if (block >= 7+512+512)
80:         panic("_bmap: block>big");
81:     if (block<7) {
82:         if (create && !inode->i_zone[block])
83:             if ((inode->i_zone[block]=new_block(inode->i_dev)) {
84:                 inode->i_ctime=CURRENT_TIME;
85:                 inode->i_dirt=1;
86:             }
87:         return inode->i_zone[block];
88:     }
89:     block -= 7;

```

```

286: #define COPYBLK(from,to) \
287: __asm__("cld\n\t" \
288: "rep\n\t" \
289: "movsl\n\t" \
290: ::"c" (BLOCK_SIZE/4),"S" (from),"D" (to) \
291: )

```

```

293: /*
294:  * bread_page reads four buffers into memory at the desired address. It's
295:  * a function of its own, as there is some speed to be got by reading them
296:  * all at the same time, not waiting for one to be read, and then another
297:  * etc.
298: */

```

```

299: void bread_page(unsigned long address,int dev,int b[4])
300: {

```

```

301:     struct buffer_head * bh[4];
302:     int i;
303:
304:     for (i=0 ; i<4 ; i++)
305:         if (b[i]) {
306:             if ((bh[i] = getblk(dev,b[i])))
307:                 if (!bh[i]->b_uptodate)
308:                     ll_rw_block(READ,bh[i]);
309:             } else
310:                 bh[i] = NULL;
311:             for (i=0 ; i<4 ; i++,address += BLOCK_SIZE)
312:                 if (bh[i]) {
313:                     wait_on_buffer(bh[i]);
314:                     if (bh[i]->b_uptodate)
315:                         COPYBLK((unsigned long) bh[i]->b_data,address);
316:                     brelse(bh[i]);
317:                 }
318: } « end bread_page »

```

```

319: static inline void wait_on_buffer(struct buffer_head * bh)
320: {
321:     cli();
322:     while (bh->lock)
323:         sleep_on(&bh->b_wait);
324:     sti();

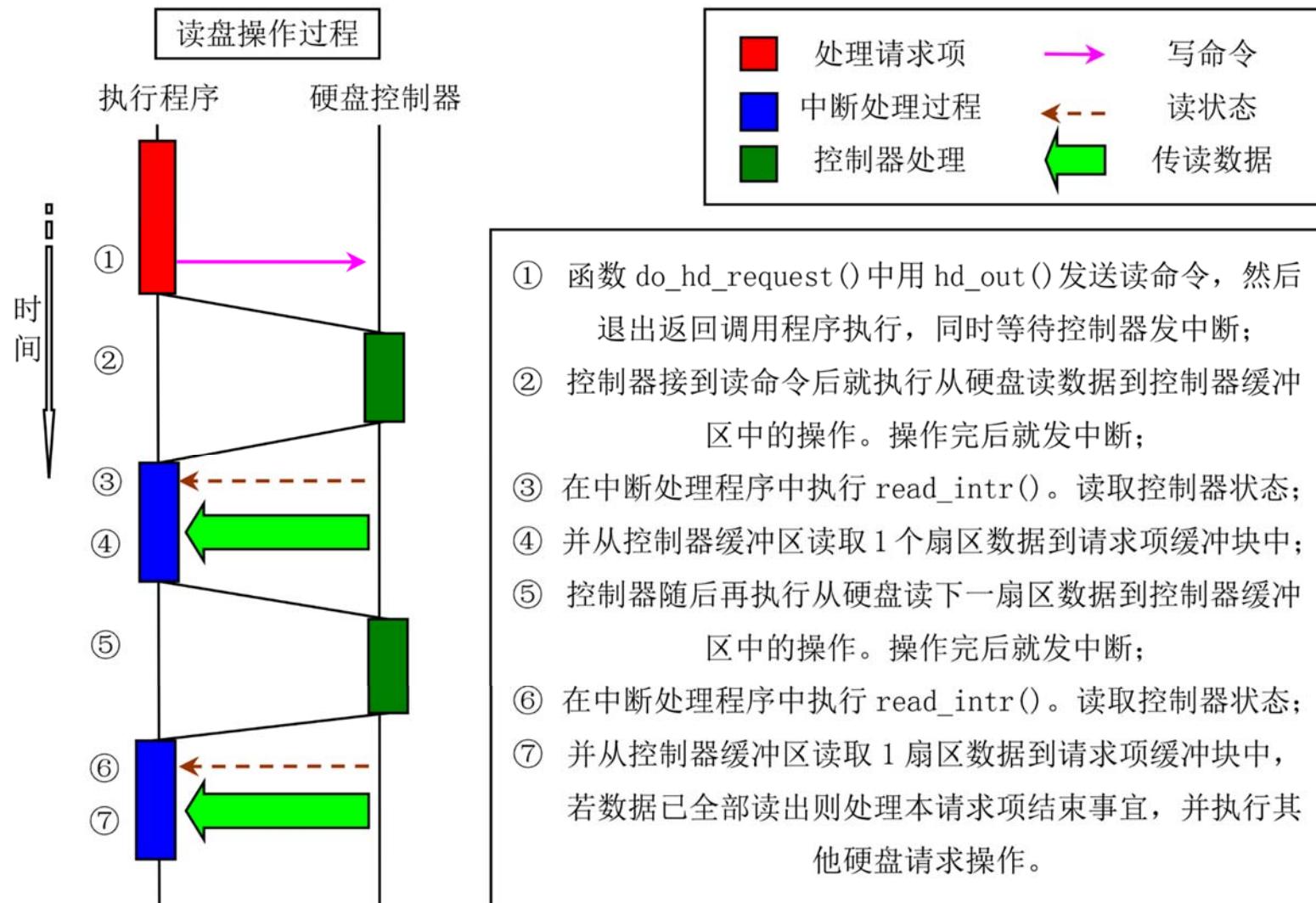
```

```

325: void sleep_on(struct task_struct **p)
326: {
327:     struct task_struct *tmp;
328:
329:     if (!p)
330:         return;
331:     if (current == &(init_task.task))
332:         panic("Task[0] trying to sleep");
333:     tmp = *p;
334:     *p = current;
335:     current->state = TASK_UNINTERRUPTIBLE;
336:     schedule();
337:     if (tmp)
338:         tmp->state=0;
339: }

```

# 读硬盘数据操作的时序关系



```

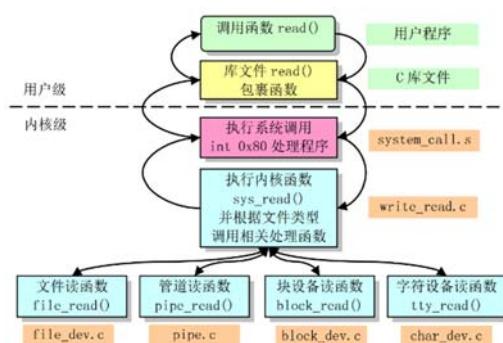
74: fn_ptr sys_call_table[] = { sys_setup, sys_exit, sys_fork, sys_read,
75: sys_write, sys_open, sys_close, sys_waitpid, sys_creat, sys_link,
76: sys_unlink, sys_execve, sys_chdir, sys_time, sys_mknod, sys_chmod,
77: sys_chown, sys_break, sys_stat, sys_lseek, sys_getpid, sys_mount,
78: sys_umount, sys_setuid, sys_getuid, sys_stime, sys_ptrace, sys_alarm,
79: sys_fstat, sys_pause, sys_utime, sys_stty, sys_gtty, sys_access,
80: sys_nice, sys_ftime, sys_sync, sys_kill, sys_rename, sys_mkdir,
81: sys_rmdir, sys_dup, sys_pipe, sys_times, sys_prof, sys_brk, sys_setgid,
82: sys_getgid, sys_signal, sys_geteuid, sys_getegid, sys_acct, sys_phys,
83: sys_lock, sys_ioctl, sys_fcntl, sys_mpx, sys_setpgid, sys_ulimit,
84: sys_uname, sys_umask, sys_chroot, sys_ustat, sys_dup2, sys_getppid,
85: sys_getpgrp, sys_setsid, sys_sigaction, sys_sgetmask, sys_ssetmask,
86: sys_setreuid,sys_setregid };

```

```

55: int sys_read(unsigned int fd,char * buf,int count)
56: {
57:     struct file * file;
58:     struct m_inode * inode;
59:
60:     if (fd>=NR_OPEN || count<0 || !(file=current->filp[fd]))
61:         return -EINVAL;
62:     if (!count)
63:         return 0;
64:     verify_area(buf,count);
65:     inode = file->f_inode;
66:     if (inode->i_pipe)
67:         return (file->f_mode&1)?read_pipe(inode,buf,count):-EIO;
68:     if (S_ISCHR(inode->i_mode))
69:         return rw_char(READ,inode->i_zone[0],buf,count,&file->f_pos);
70:     if (S_ISBLK(inode->i_mode))
71:         return block_read(inode->i_zone[0],&file->f_pos,buf,count);
72:     if (S_ISDIR(inode->i_mode) || S_ISREG(inode->i_mode)) {
73:         if (count+file->f_pos > inode->i_size)
74:             count = inode->i_size - file->f_pos;
75:         if (count<0)
76:             return 0;
77:         return file_read(inode,file,buf,count);
78:     }
79:     printk("(Read)inode->i_mode=%06o\n\r",inode->i_mode);
80:     return -EINVAL;
81: } « end sys_read »

```



```

17: int file_read(struct m_inode * inode, struct file * filp, char * buf, int count)
18: {
19:     int left,chars,nr;
20:     struct buffer_head * bh;
21:
22:     if ((left=count)<=0)
23:         return 0;
24:     while (left) {
25:         if (((nr = bmap(inode,(filp->f_pos)/BLOCK_SIZE))) {
26:             if (!(bh=bread(inode->i_dev,nr)))
27:                 break;
28:         } else
29:             bh = NULL;
30:         nr = filp->f_pos % BLOCK_SIZE;
31:         chars = MIN( BLOCK_SIZE-nr , left );
32:         filp->f_pos += chars;
33:         left -= chars;
34:         if (bh) {
35:             char * p = nr + bh->b_data;
36:             while (chars-->0)
37:                 put_fs_byte(*p++,buf++);
38:             brelse(bh);
39:         } else {
40:             while (chars-->0)
41:                 put_fs_byte(0,buf++);
42:         }
43:     } « end while left »
44:     inode->i_atime = CURRENT_TIME;
45:     return (count-left)?(count-left):-ERROR;
46: } « end file_read »
47: struct buffer_head * bread(int dev,int block)
48: {
49:     struct buffer_head * bh;
50:
51:     if (!(bh=getblk(dev,block)))
52:         panic("bread: getblk returned NULL\n");
53:     if (bh->b_uptodate)
54:         return bh;
55:     ll_rw_block(READ,bh);
56:     wait_on_buffer(bh);
57:     if (bh->b_uptodate)
58:         return bh;
59:     brelse(bh);
60:     return NULL;
61: }
62: void ll_rw_block(int rw, struct buffer_head * bh)
63: {
64:     unsigned int major;
65:
66:     if ((major=MAJOR(bh->b_dev)) >= NR_BLK_DEV ||
67:         !(blk_dev[major].request_fn)) {
68:         printk("Trying to read nonexistent block-device\n\r");
69:         return;
70:     }
71:     make_request(major,rw,bh);
72: }

```

**Diagram of the Block I/O Path:**

The diagram shows the data flow from user-space to a block device:

- A blue box labeled "该可操作" (User Operation) has an arrow pointing to a yellow box labeled "bread".
- An arrow points from "bread" to a green box labeled "高缓存区" (High Cache Area).
- An arrow points from "高缓存区" to a green box labeled "ll\_rw\_block".
- An arrow points from "ll\_rw\_block" to a green box labeled "wait\_on\_buffer(bh)".
- An arrow points from "wait\_on\_buffer(bh)" to a green box labeled "brelse(bh)".
- An arrow points from "brelse(bh)" to a green box labeled "设备控制器" (Device Controller).
- An arrow points from "设备控制器" to a gray box labeled "块设备" (Block Device).

```

88: static void make_request(int major,int rw, struct buffer_head * bh)
89: {
90:     struct request * req;
91:     int rw_ahead;
92:
93: /* WRITEA/READA is special case - it is not really needed, so if the */
94: /* buffer is locked, we just forget about it, else it's a normal read */
95:     if ((rw_ahead = (rw == READA || rw == WRITEA))) {
96:         if (bh->b_lock)
97:             return;
98:         if (rw == READA)
99:             rw = READ;
100:        else
101:            rw = WRITE;
102:    }
103:    if (rw!=READ & rw!=WRITE)
104:        panic("Bad block dev command, must be R/W/RA/WA");
105:    lock_buffer(bh);
106:    if ((rw == WRITE & !bh->b_dirt) || (rw == READ & bh->b_uptodate)) {
107:        unlock_buffer(bh);
108:        return;
109:    }
110: repeat:
111: /* we don't allow the write-requests to fill up the queue completely:
112: * we want some room for reads: they take precedence. The last third
113: * of the requests are only for reads.
114: */
115:    if (rw == READ)
116:        req = request+NR_REQUEST;
117:    else
118:        req = request+((NR_REQUEST*2)/3);
119: /* find an empty request */
120:    while (--req >= request)
121:        if (req->dev<0)
122:            break;
123: /* if none found, sleep on new requests: check for rw_ahead */
124:    if (req < request) {
125:        if (rw_ahead) {
126:            unlock_buffer(bh);
127:            return;
128:        }
129:        sleep_on(&wait_for_request);
130:        goto repeat;
131:    }
132: /* fill up the request-info, and add it to the queue */
133:    req->dev = bh->b_dev;
134:    req->cmd = rw;
135:    req->errors=0;
136:    req->sector = bh->b_blocknr<<1;
137:    req->nr_sectors = 2;
138:    req->buffer = bh->b_data;
139:    req->waiting = NULL;
140:    req->bh = bh;
141:    req->next = NULL;
142:    add_request(major+blk_dev,req);
143: } end make_request

```

```

270: struct buffer_head * bread(int dev,int block)
271: {
272:     struct buffer_head * bh;
273:
274:     if (!(bh=getblk(dev,block)))
275:         panic("bread: getblk returned NULL\n");
276:     if (bh->b_uptodate)
277:         return bh;
278:     if (bh->b_uptodate)
279:         return bh;
280:     brelse(bh);
281:     return NULL;
282: }
283: 39: static inline void wait_on_buffer(struct buffer_head * bh)
284: 40: {
41:     cli();
42:     while (bh->b_lock)
43:         sleep_on(&bh->b_wait);
44:     sti();
45: }
46: }
sleep_on () -> Schedule()

145: void ll_rw_block(int rw, struct buffer_head * bh)
146: {
147:     unsigned int major;
148:
149:     if ((major=MAJOR(bh->b_dev)) >= NR_BLK_DEV ||
150:         !(blk_dev[major].request_fn)) {
151:         printk("Trying to read nonexistent block-device\n\r");
152:         return;
153:     }
154:     make_request(major,rw,bh);
155: }


```

```

54: static void add_request(struct blk_dev_struct * dev, struct request * req)
55: {
56:     struct request * tmp;
57:
58:     req->next = NULL;
59:     cli();
60:     if (req->bh)
61:         req->bh->b_dirt = 0;
62:     if (!tmp = dev->current_request) {
63:         dev->current_request = req;
64:         sti();
65:         (dev->request_fn)();
66:         return;
67:     }
68:     for ( ; tmp->next ; tmp=tmp->next)
69:         if ((IN_ORDER(tmp,req) ||
70:              !IN_ORDER(tmp,tmp->next)) &&
71:             IN_ORDER(req,tmp->next))
72:             break;
73:     req->next=tmp->next;
74:     tmp->next=req;
75:     sti();
76: } end add_request

```

**电梯算法**

一路返回到  
wait\_on\_buffer

```

145: void ll_rw_block(int rw, struct buffer_head * bh)
146: {
147:     unsigned int major;
148:
149:     if ((major=MAJOR(bh->b_dev)) >= NR_BLK_DEV || 
150:         !(blk_dev[major].request_fn)) {
151:         printk("Trying to read nonexistent block-device\n\r");
152:         return;
153:     }
154:     make_request(major,rw,bh);
155: }

294: void do_hd_request(void)
295: {
296:     int i,r = 0;
297:     unsigned int block,dev;
298:     unsigned int sec,head,cyl;
299:     unsigned int nsect;
300:
301:     INIT_REQUEST;
302:     dev = MINOR(CURRENT->dev);
303:     block = CURRENT->sector;
304:     if (dev >= 5*NR_HD || block+2 > hd[dev].nr_sects) {
305:         end_request(0);
306:         goto repeat;
307:     }
308:     block += hd[dev].start_sect;
309:     dev /= 5;
310:     __asm__ ("divl %4::=a" (block), "=d" (sec):"0" (block), "1" (0),
311:             "r" (hd_info[dev].sect));
312:     __asm__ ("divl %4::=a" (cyl), "=d" (head):"0" (block), "1" (0),
313:             "r" (hd_info[dev].head));
314:     sec++;
315:     nsect = CURRENT->nr_sectors;
316:     if (reset) {
317:         reset = 0;
318:         recalibrate = 1;
319:         reset_hd(CURRENT_DEV);
320:         return;
321:     }
322:     if (recalibrate) {
323:         recalibrate = 0;
324:         hd_out(dev,hd_info[CURRENT_DEV].sect,0,0,0,
325:                WIN_RESTORE,&recal_intr);
326:         return;
327:     }
328:     if (CURRENT->cmd == WRITE) {
329:         hd_out(dev,nsect,sec,head,cyl,WIN_WRITE,&write_intr);
330:         for(i=0 ; i<3000 && !(r=inb_p(HD_STATUS)&DRO_STAT) ; i++)
331:             /* nothing */;
332:         if (!r) {
333:             bad_rw_intr();
334:             goto repeat;
335:         }
336:         port_write(HD_DATA,CURRENT->buffer,256);
337:     } else if (CURRENT->cmd == READ) {
338:         hd_out(dev,nsect,sec,head,cyl,WIN_READ,&read_intr);
339:     } else
340:         panic("unknown hd-command");
341: } « end do_hd_request »

```

do\_hd\_request()

```

1: #define outb(value,port) \
2: __asm__ ("outb %%al,%%dx"::"a" (value),"d" (port))
3:
4:
5: #define inb(port) ({ \
6:     unsigned char _v; \
7:     __asm__ volatile ("inb %%dx,%%al":="a" (_v):"d" (port)); \
8:     _v; \
9: })
10:
11: #define outb_p(value,port) \
12: __asm__ ("outb %%al,%%dx\n" \
13: "\tjmpf 1f\n" \
14: "1:\tjmpf 1f\n" \
15: "1:::a" (value), "d" (port))

```

```

17: #define inb_p(port) ({ \
18:     unsigned char _v; \
19:     __asm__ volatile ("inb %%dx,%%al\n" \
20: "\tjmpf 1f\n" \
21: "1:\tjmpf 1f\n" \
22: "1:::a" (_v):"d" (port)); \
23:     _v; \
24: })

```

```

45: struct hd_i_struct {
46:     int head,sect,cyl,wpcom,lzone,ctl;
47: };
48: #ifdef HD_TYPE
49: struct hd_i_struct hd_info[] = { HD_TYPE };
50: #define NR_HD ((sizeof (hd_info))/(sizeof (struct hd_i_struct)))
51: #else
52: struct hd_i_struct hd_info[] = { {0,0,0,0,0,0}, {0,0,0,0,0,0} };
53: static int NR_HD = 0;
54: #endif
55:
56: static struct hd_struct {
57:     long start_sect;
58:     long nr_sects;
59: } hd[5*MAX_HD]={ {0,0} };

```

```

1: #define port_read(port,buf,nr) \
2: __asm__ ("cld;rep;insw"::"d" (port),"D" (buf),"c" (nr))

```

read\_intr()

硬盘信息相关，  
端口读写

```

108: static inline void end_request(int uptodate)
109: {
110:     DEVICE_OFF(CURRENT->dev);
111:     if (CURRENT->bh->b_uptodate) {
112:         CURRENT->bh->b_uptodate = uptodate;
113:         unlock_buffer(CURRENT->bh);
114:     }
115:     if (!uptodate) {
116:         printk(DEVICE_NAME " I/O error\n\r");
117:         printk("dev %04x, block %d\n\r",CURRENT->dev,
118:               CURRENT->bh->b_blocknr);
119:     }
120:     wake_up(&CURRENT->waiting);
121:     wake_up(&wait_for_request);
122:     CURRENT->dev = -1;
123:     CURRENT = CURRENT->next;
124: }

```

wake\_up

硬盘信息相关，  
端口读写

```

180: static void hd_out(unsigned int drive,unsigned int nsect,unsigned int sect,
181:                     unsigned int head,unsigned int cyl,unsigned int cmd,
182:                     void (*intr_addr)(void))
183: {
184:     register int port asm("dx");
185:
186:     if (drive>1 || head>15)
187:         panic("Trying to write bad sector");
188:     if (!controller_ready())
189:         panic("HD controller not ready");
190:     do_hd = intr_addr;
191:     outb_p(hd_info[drive].ctl,HD_CMD);
192:     port=HD_DATA;
193:     outb_p(nsect++,port);
194:     outb_p(sec,++port);
195:     outb_p(cyl,++port);
196:     outb_p((drive<<4)|head,++port);
197:     outb_p(cmd,++port);
198:     outb_p(0xa0,(drive<<4)|head,++port);
199:     outb(cmd,++port);
200: } « end hd_out »

```

```

250: static void read_intr(void)
251: {
252:     if (win_result()) {
253:         bad_rw_intr();
254:         do_hd_request();
255:         return;
256:     }

```

```

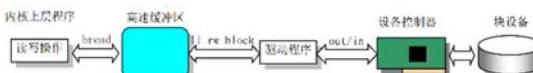
257:     port_read(HD_DATA,CURRENT->buffer,256);
258:     CURRENT->errors = 0;
259:     CURRENT->buffer += 512;
260:     CURRENT->sector++;
261:     if (--CURRENT->n_sectors) {
262:         do_hd = &read_intr;
263:         return;
264:     }

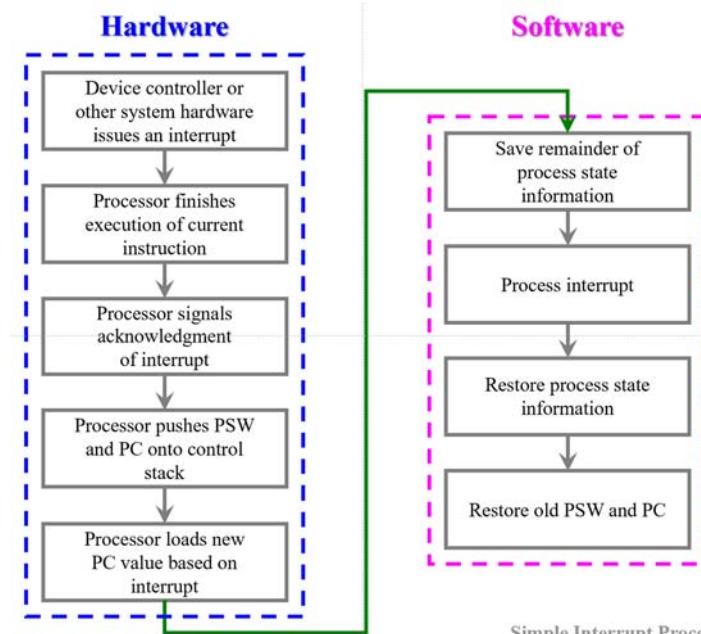
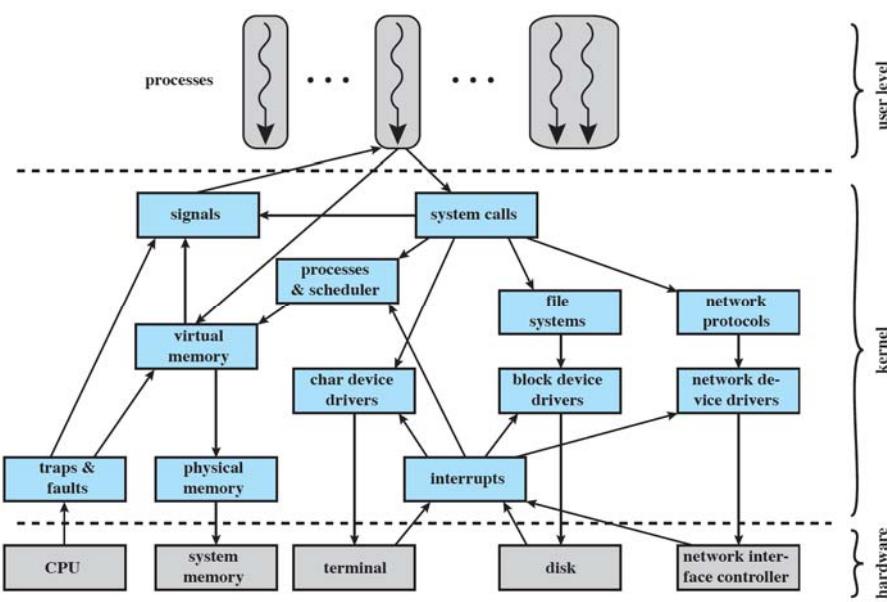
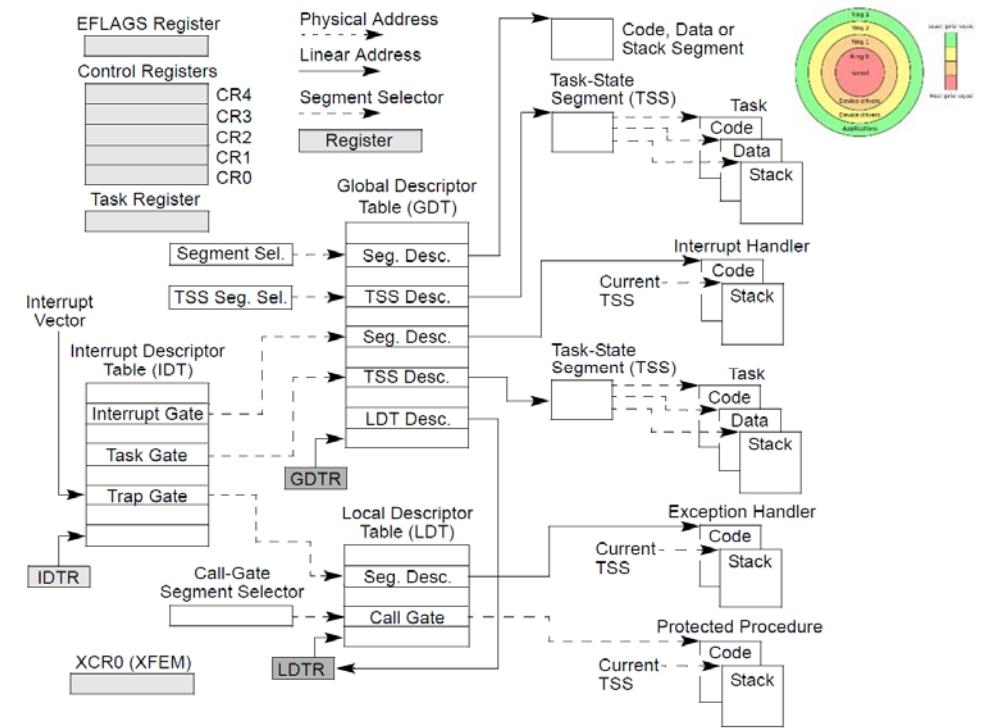
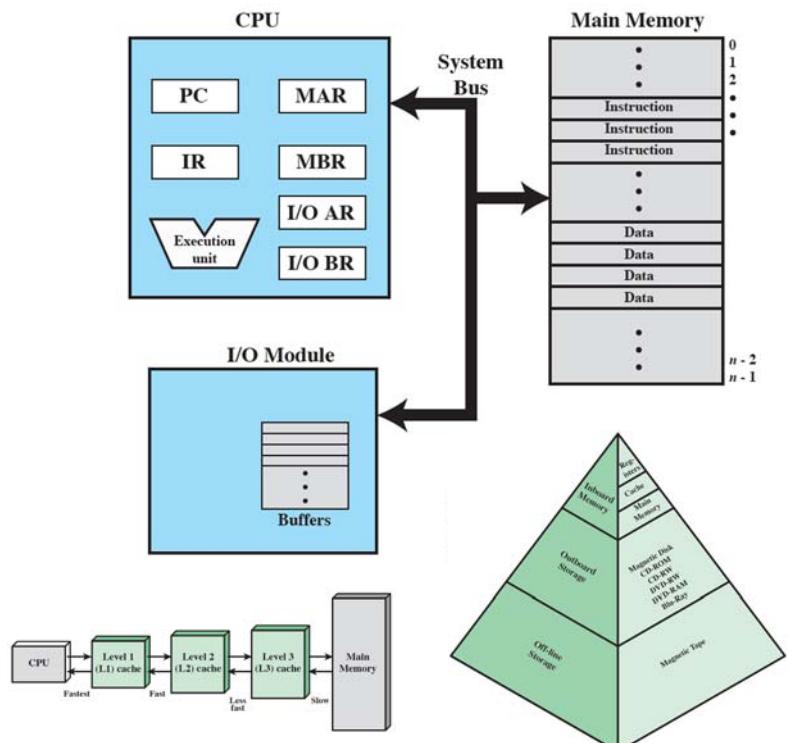
```

```

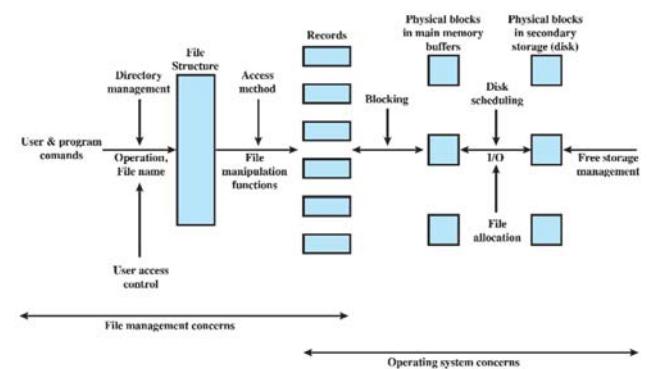
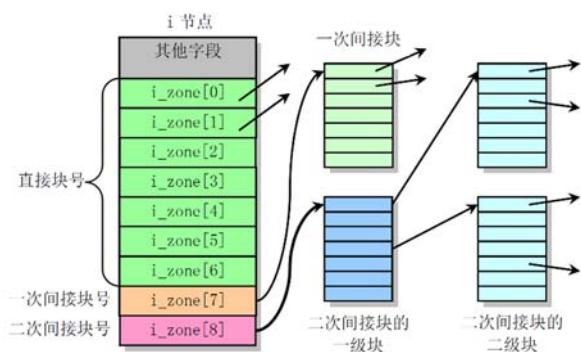
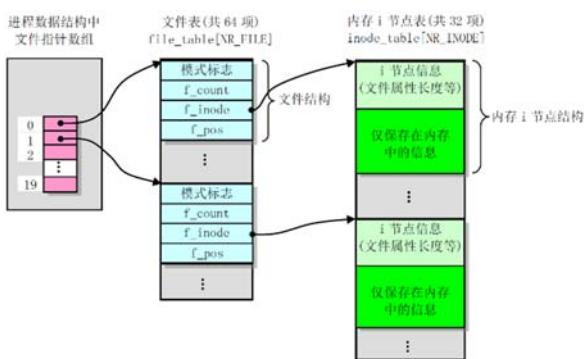
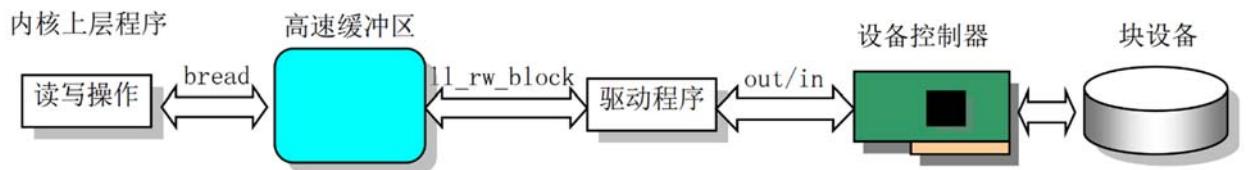
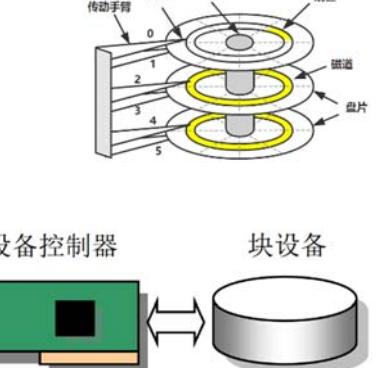
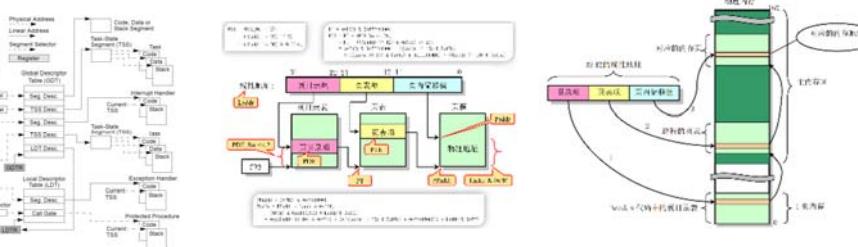
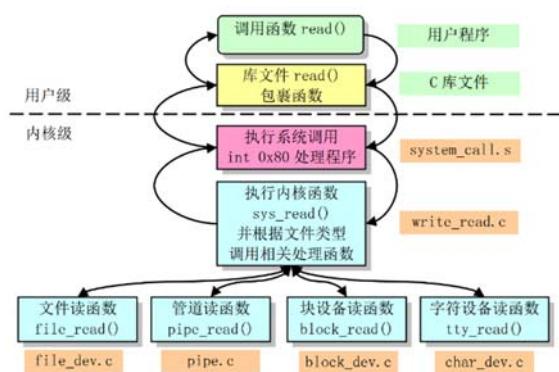
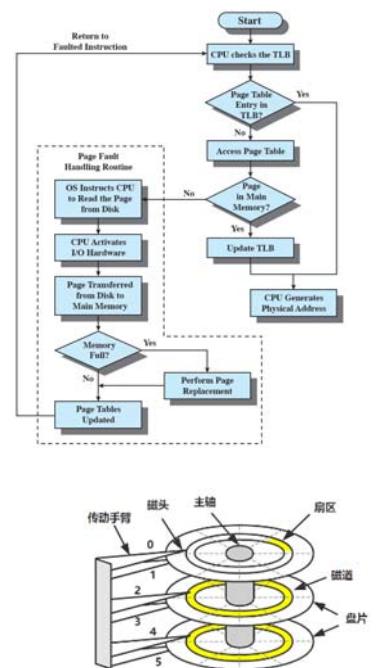
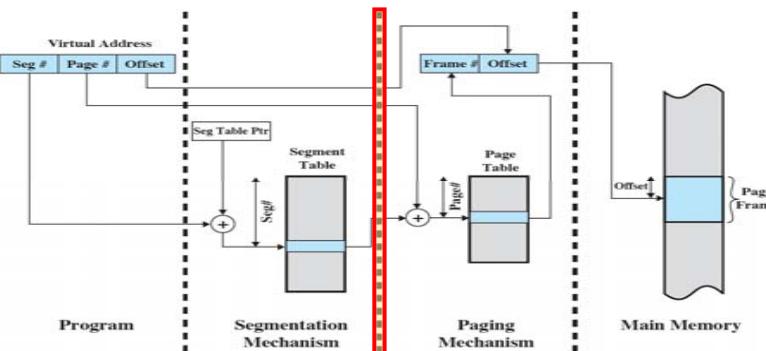
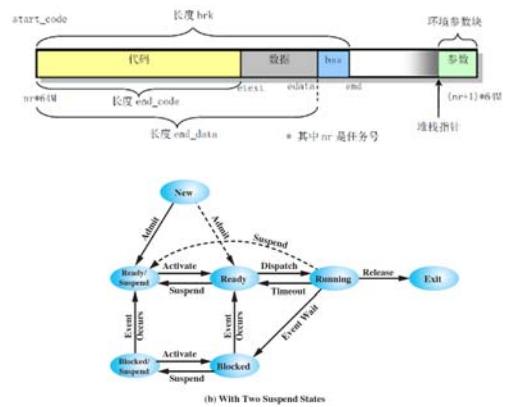
265:     end_request(1);
266:     do_hd_request();
267: }

```





Simple Interrupt Processing



# 漫谈 OS

---

- 虚拟
  - 虚拟内存、虚拟文件系统
- 并行
  - 流水、进程、线程
- 硬件加速
  - Cache, TLB, 缓冲, 磁盘缓冲
- 抽象
  - 层次模型
- 主奴机制 -- **operating**
  - 软硬结合
  - 模式
- 并发、共享
  - 加锁, 信号量
- 中断
- 栈
  - 各种巧妙用法
- 内存映射
  - Map、变换
- 调度
  - Trade-off
- 视角
  - Hardware、OS、User