一、x86 位数小的在小内存

1. 储存方式: little end first | OX78 小内存
OX12345678→ OX12 大内存 | OX56
| OX34
| OX12

immediate number: 井号
$0x — hex | $0 — octal
$53decimal

2. 操作: INC, DEC, XOR(异或取反位), SHL, SARI(算术右移,补符号位),
SHRI(逻辑右移,补0), ROL(循环左移), ROR(循环右移)

ORL %ECX %EBX → destination and first source
↑data type ↑second source →%EBX←EBX OR ECX
operation L:long(32bit) W:word(16bit) B:byte(8bit) default: 可以只用1reg:
偏移值 除JESP 2,4,8 (%ecx, 4)
②memory operand: displacement (SRI, SR2, scale) → SR2*scale +SRI +displacement?

MOV src,dst: 将 src 系列之值到 dst "dst"是reg/memory
LEA src,dst "src"是memory, "dst"是reg 把src的个数输入'dst'中
④CMPL %EAX, %EBX #flags←EBX-EAX '2'-'1' 又改其flag,reg不变!
TESTL %EAX, %EBX #flags← (EBX AND EAX) '2' AND '1'
jne jb jbe je jae ja-unsigned
≠ < ≤ = ≥ >
jne jl jle je jge jg-signed

④data size conversion: MOVS/Z from to
movsBL %AH, %ECX # ECX←sign extend to 32bit (AH)
movzBL 4(%EBP), %EAX #EAX←zero extend to 32bit (M[EBP+4])
⑤input, output independent →call data to/from AL, AX, EAX
IN port, dest.reg  OUT src.reg, port  port. 8bit imm或DX
INB $0X40, %AL #AL←M[0X40]  INW (%DX), %AX #AX←M[DX], AH←M[DX+1]
OUTB %AL, (%DX)#M[DX]←AL OUTW %AX, $8 #M[8]←AL, M[6]←AH

3. calling convention
①参数(argument): 从右往左入栈 | local 变量 | ←EBP, ESP
int fun (int key, int* array, int size) | old EBP |
why: I.参数内存地址fixed. lst arg | ret. add | argument可能会reuse:
is 8(%ebp) II.允许 variable number 的参数 | int key | 可能会变!
②caller saved reg(主程序保护), EAX, ECX, EDX, EFLAGS | int* array |
callee saved reg(子程序保存,出去还原), ESP, EBP, EBX, EDI, ESI | int size |
[caller sequence]:
I.save caller-saved reg(push) II.argument入栈 (push) III. call子程序
IV. argument 出栈(add → %esp) V. restore caller-saved reg(pop)
[callee sequence]: save old base pointer, get a new one
I. push %ebp; movl %esp, %ebp II.save callee-saved reg(push)
III.make space for local变量 (subl $4, %esp) IV. 正常执行
V. tear down local变量 (addl $4, %esp) VI. restore callee-saved reg(pop)
VII. load old base pointer, return (leave; ret)

主程序 初始化reg 子程序:
XORL %ESI, %ESI 赋予ESI 到 0 push %ebp 固定流程
push ... argument入栈 movl %esp, %ebp
call func 调用子程序 subl $4, %esp local变量
addl $8, %esp movl %eax, -4(%ebp) 在local变量
movl -4(%ebp), %eax 返回值
ESP←EBP+4; EBP←M[EBP] leave why we need: ①protect shared resources
EIP←M[ESP], ESP←ESP+4 ret ②prevent interrupts from corrupting data
二、Synchronization ← ③race condition ④critical section

1. system software's role: ①virtualization (the illusion of multiple/practically
unlimited resource) ②protection ③abstraction 隐藏底层细节, 提供simpler接口
2. I/O: ①independent I/O: using distinct instructions to separates I/O
ports from memory address IN, OUT← keyboard (device)
②memory-mapped I/O: no new instructions, has a region of memory
address that is set aside for I/O VGA
3. shared data / structure interrupt不能改变reg中内容
①interrupt handler: preserve the contents of regs, avoid overwriting memory
location used by interrupted program interrupt: preserve all regs used
I.volatile: this value may change at any time so compiler won't optimize it
volatile int busy = 0 不要什么变量都volatile, always reload 变量⇒more memory
ops → slower program. critical section should be short: avoid delaying
②critical section: 用IF实现; device; long delays crash system

①CLI: clear IF
(code): atomically: the entire critical section either executed or none executed.
②STI: set IF (no moving memory ops in/out critical section)
all processors have equal latency to all memory banks
4. multiprocessors, SMP (symmetric multiprocessor)

①spin-lock, lock是"1", unlock是"0"
program 要等即解锁而被执行其他 other processor
(用于interrupt handler(短!)) spin-lock-irqsave 先CLI再lock (避免deadlock)
locking must be atomic with ispin_unlock_irqstore,然unlock再STI
respect to other processors! (不然出现多个processor实用lock)
[deadlock]: a thread always hold the lock is waiting for an interrupt to finish.
But the interrupt is tring to obtain the lock and fail, thus the program can't
work. Lock always belong to the thread.单个 processor上发生!
[livelock]: multiple threads all need to acquire some same locks, but none of
them can get all locks. They will release the locks they acquired and try for
new, but never get all. The locks never belong to one thread.多个 processor!
解决办法: 按照同一顺序acquire locks ← lock ordering
②semaphore: 允许fixed number 的 processor流访问一个critical section(FIFO)
P(down): decrement 空闲获取 V(up): increment 用完释放 | MUTEX: thread = 1
当一个程序等待semaphore时, 允许其他程序在processor上运行 用于system call, 尤其
↳ 不能用于interrupt handler, 要等其完毕执行; 可能deadlock 是long critical section.
↳先acquire semaphore(mutex)再acquire spin lock: tring acquire semaphore
may put calling process to sleep, shouldn't happen while holding spinlock(deadlock)
②reader/writer problem: 允许无数个reader, 但writer只有1个, 且在无reader下进行
I. reader/writer spin-lock: 用于short critical sections / interrupt handler
↳ admit writer starvation (一直有reader请求, 不能写)
II. reader/writer semaphore: 用于system call
↳不admit writer starvation (writer请求时, reader排writer后面)
三、task和linkage
1. process / task: unit of scheduling ⇒独立提供resources来run program (有virtual
address space, executable code. process descriptor and at least one thread.
thread: 在process中, scheduled for execution. 一个process可以有很多thread, 它们共享
virtual addr space, system resources, 但有独立优先级和scheduled时要用的数据, 如
kernel stack和user stack.
①user-level view: process 有 process descriptor (PCB), pid(1~32767), tgid (thread
group id)(在multithread process中表示主pid)
② kernel view: 同时处理多个process, 在single per-process area(8kB)中存:
thread_info → pointer to task_struct或process descriptor(PCB)
↑ → dynamically allocated
kernel stack 不要在kernel中用recursion: overflow the kernel stack
②task structure: cyclic, double-linked list 还有 (soft irqd: soft interrupt daemon.
↳ kernel thread: boot时创造, 一直到
init→next next init_task, shutdown/重启
task→prev prev pid=1 → init_mm fs/files/signal...
[for systemcall]. map kernel pages
I.持pid*: small structure refering task II. find pid用hash table来�t large,
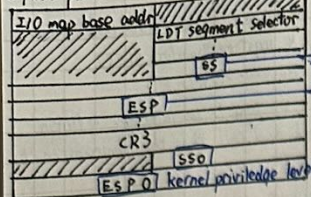sparse space map到small, clense space上 双向链表: back link 用于删除
2. 创造 process / task 到P5
→ create a copy of current program
I. user-level: 先fork再exec ← load a new program and start it
fork 创造一个子进程, 子进程和父进程同时进行, 子进程复制父进程address_space
↳ 在早期系统很浪费, 而且在exec过程中most data discarded; 还有个clone: 没见过
父进程 中返回子进程pid, 子进程 中返回0.
vfork: 创造一个子进程, 父进程 blocks, 子进程用父进程的pid和address space, 直接
共享数据, 子进程exit后, 父进程重夺control of address space.
"lazy approach" → copy-on-write: 子进程复制 page tobbs, 将write权限, 只有当子
进程write(fork中)时才会private copy 以作用: avoid work that won't be useful.
II. kernel view: fork, vfork, clone 三个system call都会调用 do_fork (在kernel中)
int do_fork (clone-flags, 如CLONE_PARENT 创造sibling task而不是child
返回new pid stack.start, new task's ESP (user level; kernel 为0)
或error负数 regs, regs value for the new task, stack-size, parent/child-tidptr);
x86又用2个参数
do_fork 先call copy_process来设置 process descriptor和child 用到的 kernel data struture
[kernel thread]: 没有对应的address space, 从last user stack kernel stack, thread-info
to excute 中继承 address space (page mapping)
3. process / task switch (单个process) (kernel stack)
创造和切时; TSS里有ESP(user stack), ESP0
p2 ___ 切换完time slice时: 换esp和esp0, 下次唤醒该process时,
p1 ___ 继续执行而非restart 可以 restoration (会存到 process descriptor中)
time

# 4. x86 task state segment TSS

## ① first part:

TSS descriptor中 segment limit定义了 I/O bitmap长度



- I/O map base addr
- LDT segment selector
- SS
- ESP
- CR3
- SS0
- ESP0 — kernel priviledge level

→ save value for user priviledge level.
→ I/O bitmap priviledge checking:
 IOPL (I/O优先级 level field,在EFLAG中) ≥ 执行I/O instruction, 否则
 check I/O port在I/O permission bitmap
 对应的bit, cleared则执行I/O instruction,
 否则 general protection exception

## ② second two parts:
interrupt redirection, I/O permission bitmap (I/O map base addr指向的start)

囧 →scheduling →goal: fair, efficient, responsive

1. jobs分类
① interactive: 响应, respond quickly to editors, GUI →用heuristics处理, 保证jobs很"互动"
② can't use lots of CPU
   实现: 超时间公平没time slice, each task run for a slice, 一直执行到no runnable, task has time left
② batch: 主要care time to completion到compilation →always叫non-real-time job优先级高
③ real-time: 将我们完成时间, periodic deadline

2. task state(存在task_struct /processor descriptor中) 只有这些状态会被执行
① TASK_RUNNING: 正在数据时/等待被数据, task被放在run queue中, 只有这些状态会被执行
② TASK_INTERRUPTIBLE: sleeping on semaphore/condition/signal,在wait queue中被唤醒 device 等待外部事件
 signal受唤醒 & TASK_UNINTERRUPTIBLE:不可中断的休眠状态(task is busy with sth. that can't stopped)
 主要针对I/O, network data, 不能被signal唤醒 (must be woken by signal) (规定义为停止)
④ TASK_STOPPED: not in a queue, must be woken up
⑤ TASK_ZOMBIE: task已经 结束, 不能删除, (但data structure还在)
 处理内存交换 (unactive用来或者 page到列(disk))

3. scheduling data structure — swapper process is always runnable (idle task)

① run queue, num of runnable task    ② priority array structure



- lock
- # runnable — number of switches
- # switches — for interactive jobs
- timestamp — task
- current task — task ←if nothing run, it will run.
- idle task
- active array
- expired array
- priority array #1
- priority array #2
- load bal for SMPs

- # active — active task数组
- bitmap — 指示哪个list非空
- real time task: 0~99/standard task:100~139
- →□□□□→□ (double linked list(存 task))

active process: runnable process但未超时
expired process: runnable process但已超时

active跑完后, change pointer better performance, 减少图像更新时的 (active指向#2) 实现 double-buffering 闪烁

## 4. scheduler policies
① SCHED_FIFO (first in, first out real-time process): 自执行直到它relinquish the CPU voluntarily (自愿放弃CPU) no pre-empted! 进程执行过程中被系统中断, 切换到另一进程 同列 run queue末端
② SCHED_RR (Round Robin real-time process): 给process分配time slice, 到时后 切换到 process 到列 run queue末端
③ SCHED_NORMAL (conventional, time-shared process): 用predefined formula重新计算 task priority, 不同的 priority 用 round robin policy.

[rescheduling and yielding (放弃)]: (主动放弃CPU)
① task change (context switch) I. current task yield by calling scheduler II. run out of time ② other place task may yield: semaphore/wake_up process
③ 在every time tick(由IRQ0上发 interrupt(IF=0)): reduce current task's time

## 五. memory allocation
1. 小函数
① kmalloc (size, flag)
用exponentially-sized slob caches (8B到4MB), allocation在physical memory上连续用于a few small items

[flag]: GFP_ATOMIC: does not sleep(不用于大内存request)
 may sleep GFP_KERNEL: kernel, driver···
 GFP_NOFS (no file sys call (不会push pages到disk))
 to wait GFP_NOIO (no I/O ops →等待内存影响可用 (page在 for page 内存中出现)
 GFP_USER: user(low priority)
 GFP_DMA: DMA accessible — direct memory access, 直接访问physical memory
 GFP_HIGHMEM: high memory(PAE) acceptable
 无法在32位系统上超过4MB物理 memory

② slab cache: 用于a lot of items, repeatedly →reduce internal fragmentation
I. slab最小是4kB(page大小), slab是连续physical memory, 被slab allocator切成individual cache 的, 会被整理成切成32, 64, 256··· 大小, 去分配

例如: file system 想 allocate 一个inode structure: 先调用 kmem_cache_create (), 计算 最优分配并返回一个 kmem_cache_t pointer (描述new inode cache 因为file system 需要 叫明, 调用 kmem_cache_alloc (以此 kmem_cache_t pointer), slab allocator会在slab 找free inode object, 若找到 (或 slab 内存不足): fetch a new slab from page and 返回 an inode object 当 file system 用完此inode后, 调用 kmem_cache_free () 去 release the inode 并将此inode在slab中的inode标为free. IV. 如果 当slab 中所有object 都是 free, 在memory不足时会将 slab 返回free page.

---

③ free pages: 用于a big, physically contiguous region
只能allocate multiples of page size    order要match!
get_free_page (flags) ←→ free_page (long) ← order是log(2为底数) of num
get_free_page (flags, order) ←→ free_page (long, order) of pages requested

④ vmalloc: 用于allocate virtual memory
[与 kmalloc对比]: kmalloc 更快, 但可能会因 memory too fragmented和 virtual mem
vmalloc 更慢. 要 modify page table entry to map physical memory and virtual mem

[memory fragmentation]:
I. external fragmentation: small blocks of free pages are scattered inside blocks of allocated pages thus can't used together
II. internal fragmentation: mismatch between the size of the memory request and the size of the memory allocated to satisfy the request

2. buddy system        order↓   [allocation]: 1 page   partially busy pair bitmaps



① why buddy system prefered?
I. contiguous page frames
II. page contiguous → page table unchanged →less TLB flush [deallocation]: 1 page
III. access by kernel through
 4MB pages →less TLB miss [eg]:
② partially busy bit: 1 if one buddy in use, "0" if both/neither in use [two blocks are buddy 叫: I. same size=2^order, 1 page
II. allocated in contiguous 物理内存
III. physical address of first block. 是2·b·(page size)的倍数

比较与现在block的大小
allocate 时优先小的 块块, 再拆到大的block
D8 (用bitmap D15 =1系判断)

| 25 | | | 30 | | D8 |
| 25 | 10 | | 30 | | D15 |
| 25 | 10 | 15 | | | A30 |
| 25 | 10 | 15 | | | A8 |
| | | 15 | | | D10 |
| | 10 | 15 | | | D20 |
| 32·25 | 10 | 15 | 20 | | A25 |
| | 10 | 15 | 20 | | A10 |
| | | 15 | 32·20 | | A15 |
| | | | 32·20 | | A20 |
| 128 | | | | | |

## 六. memory map

1. allocating memory to user mode process
① process requests for dynamic memory 被kernel 认为 non-urgent. I. 在executable file load后, 不可能马上用all code pages II. 每个 program在调用 malloc()时, process不会 马上access all memory obtained
② user mode process 类 获得dynamic memory时, 意味着address space加入了 new range of virtual addr, 但kernel会 defer allocating dynamic memory to user mode process.

2. memory map
0x00500000 — dynamically loaded libraries (not necessary in the same place in virtual memory (mapping can change)
code 0x00800000 — process addr space 被存在 memory descriptor 中, data structure mm_struct, is referenced
data — 映射mapping is the same by mm field of process descriptor
heap — for all program at kernel level, 虚拟physical memory is same (同一个kernel)
stack — virtual, physical TLB no need to flush
kernel 0x00C00000 pages — number of pages in map +地址 用r-b-tree of regions 记录

[eg]. cat和tac process (memory map) permission: read/write/execute
0x048000-0x04c000 [r-x] ④ 2286660  /bin/cat the file used to provide data
mapping region: 开始和 rw- 33 inode of  /bin/cat
结束的 virtual memory  r-x 312  0 [heap]
RW data region  rw- 2  /lib/libc-2.7.so
 r-- 2  /lib/libc-2.7.so  shared library
 r-x 26  /lib/libc-2.7.so
 rw- 2  /lib/libc-2.7.so
 rw- 21  0 [stack]
0x848000-0x04c000 r-x 4 121274894  /usr/bin/tac
 rw- 1  /usr/bin/tac
 rw- 33  0 [heap]
"shared library"
bf0f000-bf1b4000 rw- 21 0 [stack]

init ↔ mm ↔ mm_struct → virtual memory 中一块连续区域
用r-b-tree of regions 记录 page directory pointer
对flags, range, ops, file info

用file路径判断shared

[Q]: how much memory would be used by cat and tac process together?
cat process: 402 pages  tac process: 402 pages  shared library: 343 pages
[A]: solution while excluding RW data regions with RW permissions:
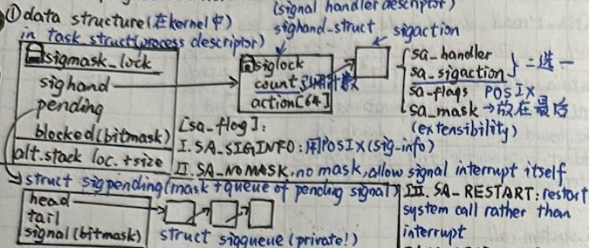$2 \cdot 402 - 343 + 4_{RW data} = 465$
solution while including all regions that contain shared data:
$2 \cdot 402 - 343 = 461$

七. signal 相当于 user-level 的 interrupt

1. 与 interrupt 异同

① 同: asynchronous (异步); not queued (send twice, handler 可能只执行1次).
can be ignored, blocked (masked) 或 caught (都有 handler func); 有 default
action (kernel 决定), can be changed by user program; 只传递 signal 不传递 data
给 handler; signal blocked while handler execute

② 异: software 生成 (no device associated). only software with permission can send

③ signal default action: ignore →signal, terminate program, terminate program
and dump core file to disk, stop program's execution, let program
continue execution. 唯 I (default action can't change): SIGKILL, SIGSTOP

2. control of signal behavior

① data structure (在 kernel 中)
in task struct (process descriptor)



② 函数: sigaction (signum, newact, oldact) ←read,
replace signal behavior [how]: I.SIG_BLOCK (add) II.SIG_UNBLOCK (rm)
sigprocmask (how, set, oldset): set bit vector of blocked signal III.SIG_SETMASK

3. signal generation: 而用 sys_kill system call 实现

① permission: sysadmin (kernel/priviledged call), process with same user id can always
send, process with same login session can send SIGCONT + continue.

② generated signal not ignored, 加列 pending 中, woken up sleeping recipient (可能
→唤醒进程)

③ 通过 send_sig() 和 force_send_sig() 函数 send signal, force 意味着 forcibly
reset an ignored signal to default behavior, unblock the signal. forced signal
always privileged (在 kernel 中生成) →一般用于 exception ←如果 exception
被 block 的话, program can't execute next instruction, kernel can't deliver
signal →deadlock.

④ info 版本: 属于一个 info structure pointer, (void*) 0 表示 traditional signals from
users, (void*) 1 表示 traditional signals from kernel, (void*) 2 表示 forced signal.

4. signal delivery ←从 interrupt, exception 和 system call 返回时触发
唯 only deliver to currently executing process.



5. 三个函数

① signal (sig, handler): 给对应 signal 赋上 handler
② exit (0): 终止程序执行并退出,
③ alarm (t): 在定时 "t" 结束后, 发送 "SIGALRM"

八. device driver

1. kernel role

① process management: creating/destroying/scheduling process
② memory management: 支持 virtual memory
③ filesystem: treat almost everything as a file
④ device control: 支持 system interactions with hardware devices, device control
operations are defined/implemented by device driver.

2. device driver

① kernel interacts with I/O devices by device drivers ←在 kernel level, 包括
data structure 和 functions. 好处是: I. encapsulated in a module II. h/w device
有 well-defined interface III. hide the detail IV. dynamic load/unload device
driver

---

② block device →用于实现 file systems     : character device
I. data accessible only in blocks of fixed size    I. contiguous space of bytes
II. addressed randomly (随机访问)         II. 有的 random access (magnetic tape driver)
III. transfer to/from device are buffered (to) 有的 sequentially (临界演奏) (sound card)
或 cached (from) 到 CD ROM, DVD, disk     III. not buffered when transfer

③ kernel abstractions for device             eg. keyboard, terminal, printers
I. device file store as real file (inode 中记录表示是 block/ character device 的地方)
II. device 由 major/minor numbers 区分: major number 是 device type, minor number 是
instance number (如果 driver 允许 more than one of a given model attached to computer)
④ registration: statically compiled in the kernel ⇒ register during kernel init phase
compiled as a kernel module ⇒ register when module loaded
[eg] insmod ./hello.ko  rmmod ./hello
⑤ 使用的 data structure (kernel): file operations, file object, inode

九. I-mail (example of device driver)

1. security: 用 I-mail administrator (user with priviledge) 成立是 sysadmin 可以 hand
off I-mail rights; without giving away other right.

2. user data / operation / lock

① data: authentication data (认证数据), message list, writing message (正在写的),
association with a program (one file a time)

| operation | user list R/W | user data | comment |
|---|---|---|---|
| read | X | ✓ | |
| write | X | ✓ | rait for message→用 * 表示需要 delivery |
| poll | X | ✓ | [message delivery]: 需要 user list R/W |
| fsync * | read (X) | ✓ | 和收 message 的 recipient's user data |
| release * | X | ✓ | semaphore. 不能同时持有两个 user data |
| authenticate | read | ✓ | semaphore! ⇒会陷入 deadlock (考虑同时 |
| set password | write | X | 多个 delivery 发生) |
| start writting * | read | ✓ | |
| delete message | X (ioct↓) | ✓ | |
| add new user | write | X → | hold writer lock, 不让后引的程序在写作 的 user. |
| delete user | write | ✓ → | sychronization (其他的需要引用时) |

② lock:   user list R/W semaphore



lock order:
I. user list r/w semaphore
(因 user list 中 read 务必 write 少)
(支型是 concurrent authorization)
II. user data semaphore
(lots of read, write)
(by one user process)

3. wait queue: 在一个 task 不能 make progress immediately (等待 semaphore, page 时).
把 task sleep: 放到 wait queue 中等待唤醒, 让 scheduler 跑别的程序
唯 race condition: task checking sleep condition, sleep 和 other task search the
queue, wake it up 之间  wait queue 是 double-linked list

① ① if (!(condition)) { while (1) { 开始 critical section (wq→lock);
    add to wait queue; set task state to TASK_INTERRUPTIBLE; 结束 critical section
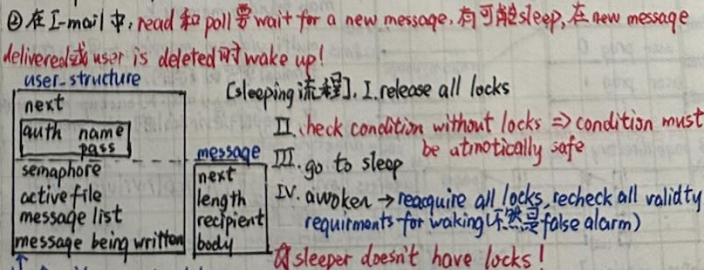    if (condition) break; ← 满足条件, wake up.
    if (!signal_pending (current)) { schedule() }; continue; } ← 继续 sleep.
    break; }} ← deliver signal, maybe return.

[函数] wake_up() →唤醒一个 task wake_up_all() →唤醒 wait_queue 中全部 task
wake_up interruptible() →唤醒一个 interruptible task

② 在 I-mail 中, read 和 poll 要 wait for a new message, 有可能 sleep, 在 new message
delivered 或 user is deleted 时 wake up!

user-structure
| next |
| auth name pass |
| semaphore |
| active file |
| message list |
| message being written |
[是一个 wait_queue_head_t]

message
| next |
| length |
| recipient |
| body |

[sleeping 流程]: I. release all locks
II. check condition without locks ⇒ condition must be atomically safe
III. go to sleep
IV. awoken → reacquire all locks, recheck all validity
requirements for waking (不然是 false alarm)
唯 sleeper doesn't have locks!

[check condition 例子]: I. 一定可行的: read an integer (并和某个常数); add integers
together; read a pointer 和与 NULL 比较 II. 不一定可行的: dereferencing a pointer;
遍历 pointer-based data structure

[reconstructure] 在 critical section 中 recalculate condition 并用 integer 存, woken 后访问这个
integer.

③I-mail read:
Ⅰ. 先 check 是否认证: unauthenticated state → throw it out
Ⅱ. while(1){ acquire user data lock:
　　check the state of file ⇒ NULL: return -EPERM; ← 在 loop 中 check, 因为可能 sleep,
　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　要么 delete 在 loop 中执行
　　同时 check 3 user 是否 deleted: 放在 lock 中执行, 避免 race condition。
　　find the message: leave the loop with the lock;
　　otherwise, release the lock and sleep; dereference 取 user data 只在 release 被 free,
　　　　　　　　　　　　　　　　　　　　　　　　← 可 release 不可能再 call,
　　if(wait_event_interruptible(···, (NULL == udata→active-file||NULL== udata→msg_list?)
　　　　　　user not deleted
　　　　　　user has a message
Ⅲ. read the data　hold the user data semaphore
　　release the semaphore, return.
4. dynamic allocation: 在 critical section外发送 避免 deadlock, 用 GFP-KERNEL
[Q]user is using I-mail, 但 user data is deleted?
[A]. delete 只 remove from user list, 不 free the structure. (active-file 的值突定是 delete?)
　[file]→NULL  unauthenticated (NOT AUTH)　| why user doesn't authorize after check?
　[file]→[user]→NULL  deleted (DELETED)　| auth need read lock, delete ops hold
　[file]⇄[user]  authenticated (Good)　| write locks on user list.
　(NOT_AUTH) authenticate (good) delete user (DELETED)
5. sub-operations 保证余不会与 delete 有影响问
①find by name: acquire the user list lock before calling the function.
②extract message being written: acquire the user data lock before calling function
③deliver mail: 允许 caller optionally hold user list lock(not held: acquire read user
　list lock), acquire recipient's user data lock, release user list lock if acquired
　within the function.
十. MP3: 一带了 TSS, LD T desc.
1.GDT: gdt-desc(48B): .word 的limit, .long的base →在 boot.s 中 lgdt gdt_desc 加入
2. idt: exception 用TRAP gate ( reserved 3 = 1), dpl=0, interrupt 用TRAP gate (~=0), dpl=0,
　system call 用 TRAP gate, dpl=3 因 INTR gate 会清空 IF flag (interrupt mask), TRAP gate 稳

| | | |
|---|---|---|
| 0X00~0X1F(32个) | exception | IF: interrupt enable flag (set: 有interrupt; clear: mask~) |
| 0X20~0X27 primary PIC | IRQ0~ IRQ7 | 因为 idt 用 linkage: c function 结束后默认用 ret, interrupt 需要 |
| 0X8~0X2F secondary PIC | IRQ8~ IRQ15 | 从kernel返回 user, 要用 iret ⇒ 用 asm linkage实现 没有 IRET 的话: kernel crashed after interrupt handler 执行完。 keyboard: 连到 irq1, handler 要 cli(), sti(). |
| 0X80 | system call INT 0X80 | rtc: 连到 irq8 |

先 master, 再 slave
3. PIC: slave 连到 master 的 irq2. 在 init 中, 先 mask interrupt: outb(0xff,0x21), 再:
ICW1: 开始初始化, edge-triggered input, cascade mode, 4ICWS outb(0x11,0x20)
ICW2: high bit of vector # outb(0x20+0, 0x21), outb(0x20+8, 0xA1)
ICW3: primary 是 bit vector of secondary outb(0x04,0x21) | enable_irq和disable_irq:
　　secondary 是 input pin of primary outb(0x02,0xA1) | 把对应 bit 的 0 或 1 来屏蔽
ICW4: ISA = x86, normal /auto EOI outb(0x01,0x21) | 对应 irq (初始化, 0xff)
因 send_eoi 在 interrupt 执行完后发送, 要向 master 和 slave 两个 PIC 发 (如果 irq 在 slave)
↳ 不然, 当前 irq 还会被认为在工作, 往后这个 irq # 发出的 interrupt 会被 blocked.

4. paging: physical mem　　virtual mem

| | |
|---|---|
| | 0 0~8MB, virtual 与 physical 相同 |
| video mem ← | video mem 0xB8000~0xB8FFF (4kB) |
| back 1, back 2, back 3 ← | back 1, back 2, back 3 4kB user-prog-paging: 将program image map 到对应的 user program 上 |
| kernel ← | kernel 4MB vidmap-paging: 将program的video page map 到video memory 上 |
| user prog 0 ← | program image 128MB multi-vidmem-mapping: 如果传入的 terminal是正在展示的,map到video 中, 否则 map 到 backup memory # (map的是 virtual 的 video memory) |
| user prog 1 ← | video map 4kB |

那是 physical memory
cr3 page directory　page table　　page directory entry
[虚拟内存矩阵]

| | 12 11 | 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| page (4kB) | address | AVL | PS | AVL | A | PCD | PWT | U/S | R/W | P |

page size (4MB)　user R/W
cr3 是 pointer to page directory table
[virtual memory]　PTE中bit6"G": set(1): TLB not flush when cr3 change
[directory # | page # | offset] 从 cr3 (PDBR) 出发, 用 10 MSB 在 page directory 中找对应 PDE,
too slow(involving too many memory ops) 跳到对应 page table, 用 next 10 bit 在 page table 中对应
　　PTE, 跳到对应 page, 用 last 12 bit 作为 page中 offset.
TLB keep translation of first 20 bits and reuse them　又又要 mapping 变了就
TLB flushed when cr3 is reloaded {movl %,%cr3, %,%eax;　要 flush TLB!
　　　　　　　　　　　　　　　　movl %,%eax, %,%cr3; }

5. file system: inode　→data block 每个block 4kB

| | | | |
|---|---|---|---|
| 0 | 0 | N-1 | 0 D-1 |

(boot block
#dir. entries 4B　总共 N 个 dentry
#inodes (N) 4B　总共 N 个 inode
4B #data blocks (D)　总共 D 个 datablock
52B reserved
4B dir. entries →

| | |
|---|---|
| length in B 4B | |
| 0th data block | 4B, data block 的编号 |
| | 最多 1023 个 data block |

| filename 32字 | ← 最多32个 character |
| file type 4B | ← 0: RTC; 1: directory; 2: regular file |
| inode # 4B | inode 编号 |
| 24B reserved | |

6个 dentry (第一个 refer to directory itself "."), 所以 max num of file: 62 (把读到的 块放入 dentry中 → 就是 抽出块!)

int32_t read-dentry-by-name (fname, dentry) 就是 抽块!
int32_t read-dentry-by-index (index, dentry)
　　　　　　　　　　　　　　dentry index
int32_t read-data (inode, offset, buf, length) 读的长度
each task 可以打开 8 个 file (但在 PCB 即 process descriptor 中)　← dynamic allocated
stdin | stdout | | | | | | |　→ open, read, write, close
stdin: read-only file (keyboard input)

| file operations table pointer 4B | |
|---|---|
| inode 4B | |
| file position 4B | |
| flags 4B | |

stdout: write-only file (terminal output)

6. system call.
①system call 要 protection boundary (user to kernel mode) ⇒ 用 asm linkage:
INT $0X80 调用, eax 存 想要的 system call #, 结束后在 eax 存 return value,
(error 时是 -1~-40 PS, 可以 定位在 error table 中找对应 error)
[assemble linkage]: Ⅰ. save regs to stack
Ⅱ. check for valid system call #
Ⅲ. call *sys-call-table (0, %eax, 4)
Ⅳ. 执行完 %eax 存 return value
Ⅴ. restore all regs, IRET

| | | |
|---|---|---|
| | ret_addr | ← ESP |
| EBX | argument 0 | |
| ECX | argument 1 | |
| EDX | argument 2 | |
| ESI | | |

② system call wrapper 以 open 为例:
open: pushl %ebx callee save regs
　　movl 0x08(esp), %ebx 装 argument
　　movl $0x05 %eax 装 system call number #
　　int $0x80 system call　　← 看返回值是不是
　　cmpl $0XFFFFFF001, %eax; jb done　对
error 况: xorl %edx, %edx　subl %eax, %edx　pushl %edx　←存正的 error number
　　call -errno-location  eax 是 pointer to errno
　　pop %ecx  movl %ecx (%eax) 把取正的 error number 装 errno 中
　　eax 装 -1　done: popl %ebx  ret

[IRET context]

| | | |
|---|---|---|
| USER-DS | (128MB+4MB-4) | |
| ESP | →user stack | |
| EFFLAG | | |
| CS | USER-CS | |
| EIP | first instruction of user program (file中byte 24-27) | |

因在 relocatable library 中找到 static data 的地址, 核心是知道现在的位置和代码里 errno 相对位置不变.
-errno-location:　　call getIP fake call, 目的: eax存 raddr
raddr:　addl $errno-raddr, %eax
　　　ret　　偏移 set1 不变
getIP: movl (%esp), %eax
　　　ret

pcb数据结构: pid, parent-pid,
file-dscs[8], saved-ebp, saved-esp
(自己程序用的 kernel stack) parent-ebp,
parent-esp(父进程用的 kernel stack)

③execute-halt 流程:　start →

kernel space
execute("shell")
↓
execute:
parse args　file的高4位 data 会有 magic
check for executable　number (0X7f,
set up paging　0X45, 0X4c, 0X46)
load file into memory
create PCB
prepare for context switch
push IRET context to kernel stack
IRET
return
→ system call linkage:
save regs
call correct system call
restore regs
IRET
halt:
restore parent data
restore parent paging
close files (all fds)
jump to execute return

user space
shell ─ type "ls" ─ execute "ls" ── 2 ── ls: ── 5 ── "do some work" ── halt() ── 6

[PCB]
PCB 1 start here　←8MB-16kB
process 1 kernel stack ←8MB-8kB
PCB 0 starts here

process 0 kernel stack ←8MB
get-pcb(pid): return (pcb-t*)(8MB-(pid+1)*8kB)

[字 cur-pcb方法]: asm volatile("movl %%esp, %0"
                            :"=r"(esp));
current_pid = (8MB - esp) / 8kB; cur-pcb = get-pcb(pid);

[execute "context switch"之后]: ← kernel stack(system call, exception, interrupt)
tss.ss0 = kernel_ds, tss_esp0 = 8MB - pid*8kB - 4 要用)
/* save parent esp, ebp */
/* IRET */

[halt "context switch"之后]:
tss 操作一样(此时的 pid是 parent_pid)
eax 装 return value, ebp 装 parent_ebp, esp 装 parent_esp
leave, return.

④ scheduler: SCHED_TASKS [3]: 一开始是 -2, 里面装每个 terminal里跑的 pid.
为什么用 pit: I. pit用 irq 0, priority 高, 不会被 interrupt II. frequency 不会像 rtc被
user program 改变

pit_handler: 先 send_eoi(1), 然后 scheduler()

[scheduler()]:
I. 先存 ebp, esp 到 cur-pcb的saved_ebp/esp中 注意判断有无 process!
II. 如果是"-2"的话, multi_vidmem_mapping(cur_index), execute "shell"
                                                已经(~+1)%3了
III. get the pcb of cur_index, user_prog_mapping
IV. multi_vidmem_mapping(pid), 改变 mapping!
V. change TSS: ss0 = kernel_ds, esp0 = 8MB - pid*8kB - 4
VI. context switch. 用当前这个 pcb的 esp, ebp来跳转! leave, ret 到
   pit handler(asm)的 iret, 再 switch 到 next process的 user space.

[switch_terminal]: (从1切换到2为例)
I. video map到 current terminal: 让 virtual 的 0XB8000 描到 physical video mem
                                        (video mem)
II. 将 video memory 复制到 back 1 上, 将 back 2 复制到 video memory 上
III. video map当前正在 schedule的 process, multi_vidmem_mapping(cur_index)
   在此时 keyboard 和 terminal 的 write 有区别. keyboard write 到当前正在展
   示的 terminal 上, 而 terminal write 要写到正在 schedule的 process 中.

十一. 补充
                        page attribute    cache disable
                   31   12 11  8  7  6↑  5  4↓  3  2  1  0
page table entry: | address | AVL | G | PAT | D | A | PCD | PWT | U/S | R/W | P |
                                                                      (在kernel中用)
                        global (set1): TLB not flushed when cr3 change

D(dirty flag): set"1"表示 page内容被修改

A(accessed flag): set"1"表示 page内容被访问(R/w)