% 实验代码

## "操作系统原理与实践"实验报告

基于内核栈切换的进程切换

## 基于内核栈切换的进程切换

预备知识

TSS切换

linux0.11 采用TS S和一条指令ljmp就能完成任务切换。 TS S切换 是指进程切换时依靠任务状态段(TS S)的切换来完

成。 在x86系统结构中,每个任务(进程和线程)都对应一个独立的TSS , TSS 是内存中的一个结构体,里面包含了 几乎所有的CPU寄存器的映像。有一个任务寄存器(Task Register,简称TR)指向当前进程对应的TSS结构体。所谓 的TSS切换就将CPU中几乎所有的寄存器都复制到TR指向的那个TSS结构体中保存起来,同时找到一个目标TSS,即 要切换到的下一个进程对应的TSS,将其中存放的寄存器映像"扣"CPU上,就完成了执行现场的切换。

实验目的

深入理解进程和进程切换的概念; 综合应用进程、CPU管理、PCB、LDT、内核栈、内核态等知识解决实际问题;

开始建立系统认识。

(1) 完成主体框架 (2) 完成PCB切换、内核栈切换、LDT切换等 (3) 修改fork(),进程基于内核栈切换,所以进程 需要创建出能完成内核栈切换的样子 (4) 完成PCB,即task\_struct结构,增加相应的内容域,同时处理由于 task\_struct所造成的影响 (5) 修改后的linux0.11仍然可以启动、可以正常使用 (6) 分析实验3的日志体会修改前后 系统运行的差别 具体步骤

内核创建一个新进程,接口调用如下 fork->int 0x80->sys call->sys fork->copy process (1) 用户程序调用fork 中通过中断指令INT 0x80 进入内核态。INT指令会把用户态的寄存器的值自动压入栈(SS, ESP, eflags, CS,IP) ,并根据中断向量表设置CS和IP。用户态的其他寄存器怎么处理? INT指令自动压栈的这个栈是用户的栈 还是内核的栈? (2) sys\_call—开始将用户态的寄存器(ds, es, fs, edx, ecx, ebx)压入栈,接着设置好内

实验内容

编写汇编的switch\_to:

核的数据段ds,fs。 此时eax存放是还是sys\_fork的系统调用号。 然后用call指令调用sys\_fork. call指令会自动把 返回地址入栈。

进程切换分析如下:

(3) sys\_fork首先查找新的进程号,并将进程号放入到eax寄存器。接着将用户态的一些寄存器入栈 (gs,esi,edi,ebp,eax),接着调用copy\_poress (4) copy\_process是一个C函数,汇编调用C函数会参数默认放在栈中。copy\_process的参数与栈中的寄存器的顺序是 一一对应的。 首先申请一块4K内存,接着这块内存开始地址为0存放PCB结构体,4K地址处初始化栈指针。内核栈需 要copy放到新进程的栈中,存放顺序与进程切换switch\_to处理对应。 有一点需要注意的是参数eax是pid号,但是子进 程应该用0 替换,因为子进程返回pid为0。Q: 为什么不一次性把所有的寄存器入栈。

(5) copy\_process执行完后eax中存放的是返回值,再加个返回值是pid,返回到sys\_fork中调用add \$20,%esp,然后

ret,这表示丢掉sys\_fork压入的(gs,esi,edi,ebp,eax),此时sp指向sys\_call调用sys\_fork的返回地址,ret指令将该返 回地址压给PC。 Q: 为什么要执行add \$20, %esp? (6) sys\_call调用sys\_fork返回后,执行pushl%eax,这个eax是sysfork的返回值此时sysfork栈帧是 (ds,es,fs,edx,ecx,ebx,eax),sys\_call 最后将栈中的值弹出,并执行iRet。 再次返回到用户态中。至此一个完成的系 统调用就完成了。 进程切换switch\_to 从时钟中断说进程切换,时钟中断的事件如下 硬件时钟中断->timer\_interrupt->do\_timer->schedule->switch\_to->新进程

(1) 硬件时钟中断,硬件自动将寄存器(SS, SP, eflags, cs, ip) 入栈, 跳转到中断处理函数timer\_interrupt处执

行。 (2) timer interrupt首先将寄存器(ds, es,fs,edx, ecx, ebx, eax)入栈,然后设置ds, es,fs,接着调用call指令 调用C函数do\_timer, 其参数也通过栈传输。do\_timer返回时, esp+4 相当于丢弃do\_timer压入栈的那个参数。 (3) do\_timer函数中检查当前进程的时间片是否为0,假如为0则调用schedule()进行任务切换。

(ss,sp,eflags,cs,eip, ds,es,fs,edx,ecx,ebx,eax+do\_timer+schedule) Q: C函数调用汇编代码如何传参数? (5) switch\_to首先处理c函数调用汇编将ebp入栈,并将当前esp保存到ebp,接着讲寄存器(ecx,ebx,eax)入 栈,此时栈帧的内容是(ebp, ecx, ebx, eax)。假如当前进程与切换的进程是同一个进程则不需要切换,弹出栈中 的4个寄存器,返回schedule->do\_timer->timer\_interrupt。 (6) 如果不是同一个进程则需要切换进程,切换步骤: 切换PCB->切换内核栈指针->切换LDT

a.切换PCB:将current和switch\_to传的pnext换内容,现在current指向的是pnext进程,ebx指向pnex进程的内存开始

b.TS S中内核栈指针的重写:中断发送时,需要把当前进程的(ss,esp,eflag,cs,eip)压入当前进程的内核栈中。如何

d.切换LDT:切换完ldt表,下一个进程执行用户态程序时使用的映射表就是自己的LDT表,地址空间实现了分离。

这些寄存器。接着ret, ret的地址也应该继续弹栈将(ds,es,fs,gs,esi,edi,edx)弹栈,最后执行iret将终端压入的

d.切换PC: switch\_to最后退出需要返回下一个子进程的用户态,current指向下一个进程的pcb,寄存器esp指向下一

个进程的内核栈。先是执行弹出switch\_to压入的寄存器(ebp,ecx,ebx,eax),所以copy\_process的栈也要对应的压入

(4) C函数shedule 里面,首先对所有进程的信号进行处理,如果进程不阻塞则把进程状态修改等待。 然后根据调度

算法找出下一个要运行的进程。 接着调用汇编s witch\_to切换到下一个进程。 此时内核栈的内容时

找到内核栈的指针呢? intel的中断机制是这么处理的,取TR寄存器存指向的TSS结构体中内核栈指针作为当前进程的 内核栈指针。所有的进程都共用这一个TSS,所以一个进程允许前应该设置好tss中的内核栈指针。 c.切换内核栈:将esp寄存器保存到当前进程的PCB中,同时将下一个进程的内核栈指针赋值给寄存器esp。注意:原 来的PCB中并没有内核栈指针,增加这个成员变量后,需要修改与这个结构体相关的硬编码。

(ss,esp,eflags,cs,eip) 弹出返回用户态。

index 21ab8ea..3482f65 100644 --- a/linux - 0.11/include/linux/sched.h +++ b/linux - 0.11/include/linux/sched.h

> long counter; long priority;

long signal;

+++ b/linux - 0.11/kernel/Makefile

\$(CC) \$(CFLAGS) \

signal.o mktime.o

index 2486b13..405c971 100644

--- a/linux - 0.11/kernel/fork.c +++ b/linux - 0.11/kernel/fork.c

#include <asm/segment.h> #include <asm/system.h>

@@ 17,6 +17,8 @@

- c - o \$\*.o \$<

@@ 82,6 +82,7 @@struct task\_struct {

struct sigaction sigaction[32];

@@ 114,6 +115,7 @@struct task\_struct {

位置。

问题回答 针对下面的代码片段: movltss,%ecx addl \$4096,%ebx movl %ebx,ESP 0(%ecx) 回答问题: (1) 为什么要加 4096;(2)为什么没有设置tss中的ss0。答:(1)ebx是进程pcb的存放的开始地址,而栈在一页4K内存地址 的最后。中断还需要 (2) 改用内核栈切换后调用init指令寄存器ss的值会自动压到内核栈中,调用iret返回时,恢 复ss寄存器的值。所以tss中的ss0不再需要。

针对代码片段: (--krnstack) = ebp; (--krnstack) = ecx; (--krnstack) = ebx; (--krnstack) = 0; 回答问题: (1) 子进

程第一次执行时, eax=? 为什么要等于这个数? 哪里的工作让eax等于这样一个数? (2)这段代码中的ebx和ecx

来自哪里,是什么含义,为什么要通过这些代码将其写到子进程的内核栈中? (3) 这段代码中的ebp来自哪里,

是什么含义,为什么要做这样的设置?可以不设置吗?为什么?答: (1) eax=0,子程序第一次执行,fork返回

0表示子程序返回。父进程的fork返回子程序的pid。 (2) ebx, ecx来自于switch\_to,表示的是switch\_to函数的参

为什么要在切换完LDT之后要重新设置fs=0x17? 而且为什么重设操作要出现在切换完LDT之后,出现在LDT之前

数,写入到子进程的内核栈是进程切换switch (3) C调用汇编要的处理要求。

diff -- git a/linux - 0.11/include/linux/sched.h b/linux - 0.11/include/linux/sched.h

long state; /\* - 1 unrunnable, 0 runnable, >0 stopped \*/

long blocked; /\* bitmap of masked signals \*/

@@ 24,7 +24,7 @@CPP =gcc - 3.4 - E - nostdinc - I./include

-OBJS = sched.o system\_call.o traps.o asm.o fork.o \

panic.o printk.o vsprintf.o sys.o exit.o \

diff -- git a/linux - 0.11/kernel/fork.c b/linux - 0.11/kernel/fork.c

又会怎么样? 答: 待定 代码补丁

> \*/ #define INIT TASK \ /\* state etc \*/ { 0,15,15, \ /\* signals \*/ 0,{{},},0, \ /\* ec,brk... \*/ 0,0,0,0,0,0, \ /\* pid etc.. \*/ 0, -1,0,0,0, diff -- git a/linux - 0.11/kernel/Makefile b/linux - 0.11/kernel/Makefile index 0afa1dc..c73e97c 100644 --- a/linux - 0.11/kernel/Makefile

extern void write\_verify(unsigned long address); long last\_pid=0; @@ 66,6 +68,7 @@int copy\_mem(int nr,struct task\_struct \* p) \* information (task[nr]) and sets up the necessary registers. It \* also copies the data segment in it's entirety. \*/ int copy\_process(int nr,long ebp,long edi,long esi,long gs,long none, ebx,long ecx,long edx, long long fs,long es,long ds, @@ 90,6 +93,7 @@int copy\_process(int nr,long ebp,long edi,long esi,long gs,long p->utime = p->stime = 0; p->cutime = p->cstime = 0; p- >start\_time = jiffies; p- >tss.back\_link = 0; p->tss.esp0 = PAGE\_SIZE + (long) p; p->tss.ss0 = 0x10;none, p- >tss.trace\_bitmap = 0x80000000; if (last\_task\_used\_math == current) \_\_asm\_\_("clts ; fnsave %0"::"m" (p - >tss.i387)); if (copy\_mem(nr,p)) { task[nr] = NULL; free\_page((long) p); @@ 132,6 +138,108 @@int copy\_process(int nr,long none, p->state = TASK\_RUNNING; /\* do this last, just in case \*/ return last\_pid;

none,

@@ 113,6 +117,8 @@int copy\_process(int nr,long ebp,long edi,long esi,long gs,long ebp,long edi,long esi,long gs,long } /\* NOTE! this doesn't copy the supervisor stack \*/

find\_empty\_process(void) 100644 \*last\_task\_used\_math = NULL; \* task[NR\_TASKS] = {&(init\_task.task), task\_struct };

{ diff -- git a/linux - 0.11/kernel/sched.c b/linux - 0.11/kernel/sched.c index 15d839b..e8c05c1 --- a/linux - 0.11/kernel/sched.c +++ b/linux - 0.11/kernel/sched.c @@ 64,6 +64,8 @@struct task\_struct struct long user\_stack [ PAGE\_SIZE>>2 ] ; struct { @@ 105,6 +107,7 @@void schedule(void) { int i,next,c; struct task\_struct wake up any interruptible tasks that have got a signal \*/ /\* check alarm, @@ 130,7 +133,7 @@void schedule(void) if (!\* -- p) continue; if ((\*p) ->state == TASK\_RUNNING && (\*p) ->counter > c) c = (\*p) ->counter, next = i; if (c) break; for(p = &LAST\_TASK; p > &FIRST\_TASK; - p) @@ 138,7 +141,8 @@void schedule(void) (\*p) ->counter = ((\*p) ->counter >> 1) + (\*p) - >priority; switch\_to(next); } int sys\_pause(void) diff -- git a/linux - 0.11/kernel/system\_call.s b/linux - 0.11/kernel/system\_call.s index 05891e1..e89a1ec 100644 --- a/linux - 0.11/kernel/system\_call.s +++ b/linux - 0.11/kernel/system\_call.s @@ 45,12 +45,16 @@EFLAGS = 0x24 OLDESP = 0x28OLDSS = 0x2C= 0 # these are offsets into the task - struct. state counter = 4 priority = 8 - signal = 12 = 16 # MUST be 16 (=len of sigaction) sigaction -blocked = (33\*16) # offsets within sigaction = 0sa handler @@ 67,6 +71,7 @@nr\_system\_calls system\_call,sys\_fork,timer\_interrupt,sys\_execve .globl .globl hd\_interrupt,floppy\_interrupt,parallel\_interrupt device not available, .globl coprocessor\_error .align 2 bad\_sys\_call: @@ 92,7 +97,7 @@system\_call: movl \$0x17,%edx # fs points to local mov %dx,%fs sys\_call\_table(,%eax,4) call pushl %eax movl current,%eax cmpl \$0,state(%eax) # state jne reschedule @@ 283,3 +288,55 @@parallel\_interrupt: outb %al,\$0x20 %eax popl iret %ebp # ebp基址寄存器 : 保存进入函数时sp栈顶位置; 这里是保存上一个函数的ebp 8(%ebp), %ebx # ebp+8是pnext的地址 addl \$4096, %bx # 栈底=ebx+4096 movl %ebx. ESPO(%ecx) # 设置当前任务的内核栈esp0 12( %ebp), %ecx#取出第二个参数, \_LDT( next )