

日曜研究室

技術的な観点から日常を綴ります

[xv6 #25] Chapter 2 – Traps, interrupts, and drivers – Code: Assembly trap handlers

テキストの34～36ページ

本文

xv6は、プロセッサにトラップを引き起こさせるようなint命令に遭遇したとき、目的に叶った何らかの処理を実行するために、x86のハードウェアをセットアップしなければならない。

x86は256個の個別の割り込みを扱うことが出来る。

割り込み0-31は、除算エラーやおかしなメモリアドレスにアクセスを試みたときのような、ソフトウェア例外のために定義されている。

xv6は、32個のハードウェア割り込みを32-63の領域に割り当てていて、64をシステムコール用の割り込みに割り当てている。

tvinit関数は、main関数から呼ばれ、idtテーブルの中の256個のエントリをセットアップする。

割り込み*i*は、`vectors[i]`の中のアドレスにあるコードによって制御される。

エントリポイントは全て違う。

なぜならxv6は割り込みハンドラに対してトラップ番号を提供しないからである。

(ここ訳があやしい。because the x86 provides does not provide the trap number to the interrupt handler.)

256個の個別のハンドラを使うということは、256個の場合を見分ける唯一の方法である。

(実際に割り込みハンドラが実行されるときにはそのハンドラからはトラップ番号を得る方法がないということかな。)

trap.cのtvinit関数付近

```
01 // Interrupt descriptor table (shared by all CPUs).
02 struct gatedesc idt[256];
03 extern uint vectors[]; // in vectors.S: array of 256 entry pointers
04 struct spinlock tickslock;
05 uint ticks;
06
07 void
```

```

08 tvinit(void)
09 {
10     int i;
11
12     for(i = 0; i < 256; i++)
13         SETGATE(idt[i], 0, SEG_KCODE<<3, vectors[i], 0);
14     SETGATE(idt[T_SYSCALL], 1, SEG_KCODE<<3, vectors[T_SYSCALL],
15     DPL_USER);
16
17     initlock(&tickslock, "time");
18 }

```

main.cのmain関数

```

01 // Bootstrap processor starts running C code here.
02 // Allocate a real stack and switch to it, first
03 // doing some setup required for memory allocator to work.
04 int
05 main(void)
06 {
07     kvmalloc();           // kernel page table
08     mpinit();             // collect info about this machine
09     lapicinit(mpbcpu());
10     seginit();            // set up segments
11     cprintf("\ncpu%d: starting xv6\n\n", cpu->id);
12     picinit();            // interrupt controller
13     ioapicinit();         // another interrupt controller
14     consoleinit();        // I/O devices & their interrupts
15     uartinit();           // serial port
16     pinit();              // process table
17     tvinit();             // trap vectors
18     binit();              // buffer cache
19     fileinit();           // file table
20     iinit();              // inode cache
21     ideinit();            // disk
22     if(!ismp)
23         timerinit();      // uniprocessor timer
24     startothers();        // start other processors (must come before
25     kinit)
26     kinit();              // initialize memory allocator
27     userinit();           // first user process (must come after kinit)
28     // Finish setting up this processor in mpmain.
29     mpmain();
30 }

```

tvinitは、ユーザのシステムコールをトラップするためにだけにT_SYSCALLを制御する。

それは、2番目の引数に1という値を渡す事によって"trap"タイプのゲートを指定する。

トラップゲートはFLフラグをクリアせず、システムコールハンドラ中に他の割り込みを許可する。

カーネルもまた、ユーザプログラムが直接的なint命令によってトラップを生成できるよう、システムコールゲートの権限をDPL_USERにセットする。

xv6は、プロセスにint命令を使って他の割り込み（例えば、デバイスの割り込み）を引き起こすことを許可しない。

もしプロセスがそれを試みたら、一般保護違反に遭遇するだろう。

（そしたらvectors[13]に処理が移る）

ユーザモードからカーネルモードに保護レベルを変更したとき、ユーザプロセスのスタックを使わないようにすべきである。

なぜなら、その内容がおかしくなってしまうかもしれないからである。

ユーザプロセスは、そのプロセスのユーザメモリじゃない部分を指すようなアドレスを%espにセットするようなバグや悪意を持ってるかもしれない。

xv6は、ハードウェアがスタックセグメントセレクタと%espのための新しい値を読み込む事を通して、タスクセグメントディスクリプタをセットアップする事によって、トラップ中にスタック切り替えを実行させるよう、x86ハードウェアをプログラムする。

switchvm関数は、ユーザプロセスのカーネルスタックのトップのアドレスを、タスクセグメントディスクリプタに保存する。

vm.cのswitchvm関数

```

01 // Switch TSS and h/w page table to correspond to process p.
02 void
03 switchvm(struct proc *p)
04 {
05     pushcli();
06     cpu->gdt[SEG_TSS] = SEG16(STS_T32A, &cpu->ts, sizeof(cpu->ts)-1,
07 0);
08     cpu->gdt[SEG_TSS].s = 0;
09     cpu->ts.ss0 = SEG_KDATA << 3;
10     cpu->ts.esp0 = (uint)proc->kstack + KSTACKSIZE;
11     ltr(SEG_TSS << 3);
12     if(p->pgdir == 0)
13         panic("switchvm: no pgdir");
14     lcr3(v2p(p->pgdir)); // switch to new address space
15     popcli();
16 }
```

トラップが発生したとき、プロセッサのハードウェアは以下の事を行う。

もしプロセッサがユーザモードで実行中だったら、タスクセグメントディスクリプタから%espと%ssをロードし、ユーザの%ssと%espを新しいスタックにプッシュする。

もしプロセッサがカーネルモードで実行中だったら、上記のような事は何も行わない。

そしたらプロセッサは%eflagsと%csと%eipレジスタをプッシュする。

いくつかのトラップでは、プロセッサはエラー語 (error word) もプッシュする。

そしたら、プロセッサは適切なIDTエントリから%eipと%csを読み込む。

xv6は、IDTエントリポイントへのエントリポイントを生成するためにPerlスクリプトを使う。

それぞれのエントリは、プロセッサがやらなかったらエラーコードをプッシュし、割り込み番号をプッシュし、そしてalltrapsにジャンプする。

vectors.pl

```

01 #!/usr/bin/perl -w
02
03 # Generate vectors.S, the trap/interrupt entry points.
04 # There has to be one entry point per interrupt number
```

```

05 # since otherwise there's no way for trap() to discover
06 # the interrupt number.
07
08 print "# generated by vectors.pl - do not edit\n";
09 print "# handlers\n";
10 print ".globl alltraps\n";
11 for(my $i = 0; $i < 256; $i++){
12     print ".globl vector$i\n";
13     print "vector$i:\n";
14     if(!($i == 8 || ($i >= 10 && $i <= 14) || $i == 17)){
15         print "    pushl \\\$0\n";
16     }
17     print "    pushl \\\$i\n";
18     print "    jmp alltraps\n";
19 }
20
21 print "\n# vector table\n";
22 print ".data\n";
23 print ".globl vectors\n";
24 print "vectors:\n";
25 for(my $i = 0; $i < 256; $i++){
26     print "    .long vector$i\n";
27 }
28
29 # sample output:
30 #   # handlers
31 #   .globl alltraps
32 #   .globl vector0
33 #   vector0:
34 #       pushl $0
35 #       pushl $0
36 #       jmp alltraps
37 #   ...
38 #
39 #   # vector table
40 #   .data
41 #   .globl vectors
42 #   vectors:
43 #       .long vector0
44 #       .long vector1
45 #       .long vector2
46 #   ...

```

`alltraps`は、プロセッサのレジスタの保存を続ける。

`%ds`, `%es`, `%fs`, `%gs`、それと汎用的なレジスタをプッシュする。（`alltraps`の`pushl %ds`から`pushal`の部分）

この努力の結果、カーネルスタックに、トラップの瞬間のプロセッサのレジスタを含む`trapframe`構造体が含まれる。（図2-2を参照）

プロセッサは、`%ss`, `%esp`, `%eflags`, `%cs`, `%eip`をプッシュする。

プロセッサもしくはトラップベクタはエラー番号をプッシュし、`alltraps`は残りをプッシュする。

トラップフレームは、カーネルから現在のプロセスに戻ったときにユーザモードのプロセッサのレジスタを復元するのに十分な全ての情報を含む。

その結果、プロセッサはトラップが開始した時点とちょうど同じ状態で続行することができる。

第1章を思い出すと、`userinit`はこのゴールを達成するために、手動でトラップフレームを構築している。

（図1-3参照）

図2-2 カーネルスタック上のトラップフレーム

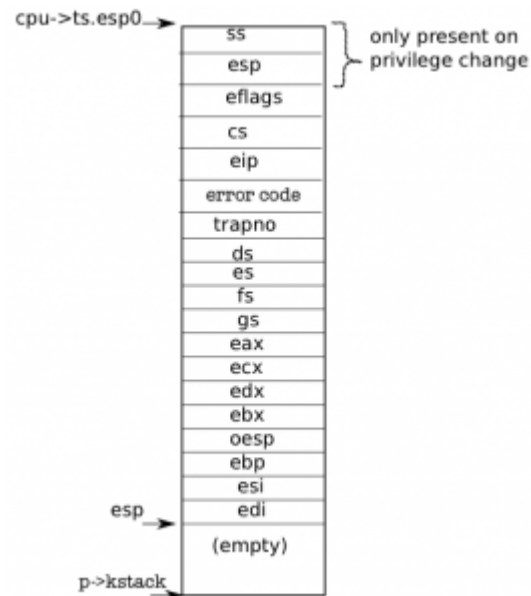
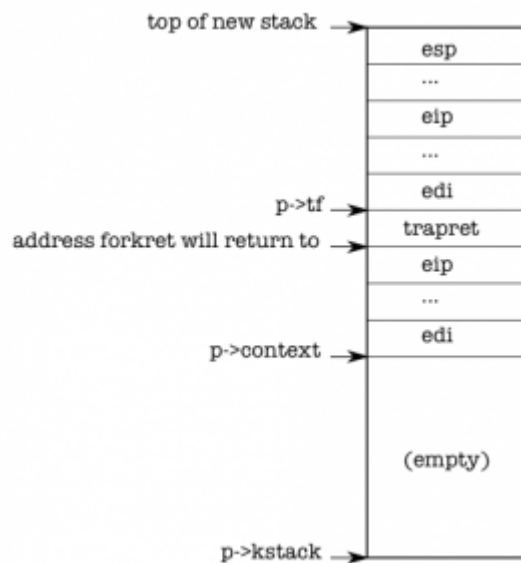


図1-3 新しいカーネルスタックのセットアップ (再掲)



trapasm.S

```

01 #include "mmu.h"
02
03 # vectors.S sends all traps here.
04 .globl alltraps
05 alltraps:
06 # Build trap frame.
07 pushl %ds
08 pushl %es
09 pushl %fs
10 pushl %gs
11 pushal
12

```

```

13  # Set up data and per-cpu segments.
14  movw $(SEG_KDATA<<3), %ax
15  movw %ax, %ds
16  movw %ax, %es
17  movw $(SEG_KCPU<<3), %ax
18  movw %ax, %fs
19  movw %ax, %gs
20
21  # Call trap(tf), where tf=%esp
22  pushl %esp
23  call trap
24  addl $4, %esp
25
26  # Return falls through to trapret...
27  .globl trapret
28  trapret:
29  popal
30  popl %gs
31  popl %fs
32  popl %es
33  popl %ds
34  addl $0x8, %esp # trapno and errcode
35  iret

```

x86.hのtrapframe構造体

```

01  // Layout of the trap frame built on the stack by the
02  // hardware and by trapasm.S, and passed to trap().
03  struct trapframe {
04      // registers as pushed by pusha
05      uint edi;
06      uint esi;
07      uint ebp;
08      uint oesp;      // useless & ignored
09      uint ebx;
10      uint edx;
11      uint ecx;
12      uint eax;
13
14      // rest of trap frame
15      ushort gs;
16      ushort padding1;
17      ushort fs;
18      ushort padding2;
19      ushort es;
20      ushort padding3;
21      ushort ds;
22      ushort padding4;
23      uint trapno;
24
25      // below here defined by x86 hardware
26      uint err;
27      uint eip;
28      ushort cs;
29      ushort padding5;
30      uint eflags;
31
32      // below here only when crossing rings, such as from user to
33      // kernel
34      uint esp;
35      ushort ss;
36      ushort padding6;
37  };

```

最初のシステムコールの場合、保存された`%eip`は、`int`命令の直後の命令のアドレスである。

`%cs`はユーザコードのセグメントセクタである。

`%eflags`は`int`命令を実行した瞬間の`eflags`レジスタの内容である。

保存された汎用レジスタの一部のように、`alltraps`は、カーネルが後で調べるためのシステムコール番号を含む`%eax`もまた保存する。

今、ユーザモードのプロセッサのレジスタは保存され、`alltraps`はプロセッサにカーネルのCコードを実行させるためのセットアップを完了することが出来る。

プロセッサは、ハンドラに入る前にセクタ`%cs`, `%ss`をセットする。

`alltraps`は`%ds`と`%es`をセットする。(`alltraps`の `movw $(SEG_KDATA<<3), %ax`から `movw %ax, %es`の部分)

それから、CPUごとのデータセグメントである`SEG_KCPU`を指すために`%fs`と`%gs`をセットする。

(`alltraps`の `movw $(SEG_KCPU<<3), %ax`から `movw %ax, %gs`の部分)

一度セグメントが適切にセットされたら、`alltraps`は、Cで書かれたトラップハンドラである`trap`関数(`trap.c`) を呼ぶことが出来る。

`alltraps`は、構築されたばかりのトラップフレームを指す`%esp`を、`trap`への引数としてスタックにプッシュする。(`alltraps`の `pushl %esp`の部分)

そして`trap`を呼ぶ。(`alltraps`の `call trap`の部分)

`trap`関数から返ってきた後、`alltraps`は、スタックポインタに加算することによってスタックの引数を取り除き、そして、`trapret`とラベルづけされたコードの実行を開始する。

我々は、第1章で、最初のユーザプロセスがユーザ空間に抜けるために`trapret`を実行したとき、このコードを追った。

同じ処理がここで起きた。

トラップフレームを通じてポップする事は、ユーザモードのレジスタを復元し、そして`iret`はユーザ空間に戻る。

trap.cのtrap関数

```

01 void
02 trap(struct trapframe *tf)
03 {
04     if(tf->trapno == T_SYSCALL) {
05         if(proc->killed)
06             exit();
07         proc->tf = tf;
08         syscall();
09         if(proc->killed)
10             exit();
11         return;
12     }
13
14     switch(tf->trapno) {
15     case T_IRQ0 + IRQ_TIMER:
16         if(cpu->id == 0) {
17             acquire(&tickslock);
18             ticks++;
19             wakeup(&ticks);

```

```

20     release(&tickslock);
21 }
22 lapiceoi();
23 break;
24 case T_IRQ0 + IRQ_IDE:
25     ideintr();
26     lapiceoi();
27     break;
28 case T_IRQ0 + IRQ_IDE+1:
29     // Bochs generates spurious IDE1 interrupts.
30     break;
31 case T_IRQ0 + IRQ_KBD:
32     kbdintr();
33     lapiceoi();
34     break;
35 case T_IRQ0 + IRQ_COM1:
36     uartintr();
37     lapiceoi();
38     break;
39 case T_IRQ0 + 7:
40 case T_IRQ0 + IRQ_SPURIOUS:
41     cprintf("cpu%d: spurious interrupt at %x:%x\n",
42             cpu->id, tf->cs, tf->eip);
43     lapiceoi();
44     break;
45
46 //PAGEBREAK: 13
47 default:
48     if(proc == 0 || (tf->cs&3) == 0){
49         // In kernel, it must be our mistake.
50         cprintf("unexpected trap %d from cpu %d eip %x (cr2=0x%x)\n",
51                 tf->trapno, cpu->id, tf->eip, rcr2());
52         panic("trap");
53     }
54     // In user space, assume process misbehaved.
55     cprintf("pid %d %s: trap %d err %d on cpu %d "
56             "eip 0x%x addr 0x%x-kill proc\n",
57             proc->pid, proc->name, tf->trapno, tf->err, cpu->id,
58             tf->eip,
59             rcr2());
60     proc->killed = 1;
61 }
62
63 // Force process exit if it has been killed and is in user space.
64 // (If it is still executing in the kernel, let it keep running
65 // until it gets to the regular system call return.)
66 if(proc && proc->killed && (tf->cs&3) == DPL_USER)
67     exit();
68
69 // Force process to give up CPU on clock tick.
70 // If interrupts were on while locks held, would need to check
71 nlock.
72 if(proc && proc->state == RUNNING && tf->trapno ==
73     T_IRQ0+IRQ_TIMER)
74     yield();
75
76 // Check if the process has been killed since we yielded
77 if(proc && proc->killed && (tf->cs&3) == DPL_USER)
78     exit();
79 }

```

ここまでの解説は、ユーザモードから発生したトラップについて話してきた。

しかしトラップはカーネルを実行してるときでも発生する。

そのような場合、ハードウェアはスタックを切り替えたり、スタックポインタやスタックセグメントセレクタを保存したりはしない。

しかしながら、ユーザモードからのトラップと同じ手順が発生し、そして同じxv6のトラップを制御するコードを実行する。

iretの後、カーネルモードの%csを復元し、プロセッサはカーネルモードでの実行を続ける。

感想

ユーザモードでシステムコールしたときに発生するトラップで、どうやってカーネルモードに移行してどうやって戻るかの話ですね。

本文にも書いてあるけど、`vectors.S`は`vectors.pl`から生成するようになってるので、一度`make`としないと存在しません。

タスクセグメントディスクリプタ周りの話がよくわかりません。

[プロテクトモード - Wikipedia](#)を読めば少しは分かるかな。

カテゴリー: 技術 | タグ: xv6 | 投稿日: 2012/3/1 木曜日 [<http://peta.okechan.net/blog/archives/1383>] |
