

# XV6进程线程源代码-部分阅读

---

## types.h

```
1 //定义了几个变量的别称
2 //所以叫做types.h
3 //uint ushort uchar pde_t
4
5 typedef unsigned int    uint;
6 typedef unsigned short ushort;
7 typedef unsigned char  uchar;
8 typedef uint pde_t;
9
10 //pde_t与uint是一样的, pde_t用在页表。
```

## param.h

```
1 //parameter.h 设置一些XV6操作系统中的参数
2
3
4 #define NPROC          64 // maximum number of processes 最大的进程数
5 #define KSTACKSIZE 4096 // size of per-process kernel stack 每个进程内核栈的最大空间是4096字节
6 #define NCPU           8 // maximum number of CPUs xv6操作系统支持的最大的CPU数是8个
7 #define NOFILE         16 // open files per process 每个进程打开的最大文件数为16个
8 #define NFILE          100 // open files per system 整个操作系统打开的最大文件数为100个
9 #define NBUF           10 // size of disk block cache 磁盘上的缓存大小
10 #define NINODE          50 // maximum number of active i-nodes 文件系统中最大的i-node数 (index_node 文件索引)
11 #define NDEV            10 // maximum major device number 最多10个驱动设备
12 #define ROOTDEV         1 // device number of file system root disk 文件系统的根目录放在1号驱动设备上
13 #define MAXARG          32 // max exec arguments 最大的执行参数的个数是32个, argv中元素的个数
14 #define LOGSIZE         10 // max data sectors in on-disk log 日志文件的=最大为10扇区 5120字节
15
```

## memlayout.h

```
1 // Memory layout
2 //内存布局的情况
3
4 #define EXTMEM 0x100000 // Start of extended memory
5 #define PHYSTOP 0xE000000 // Top physical memory
6 #define DEVSPACE 0xFE000000 // Other devices are at
   high addresses
7
8 // Key addresses for address space layout (see kmap in vm.c for
   layout)
9 //kmap表
10 #define KERNBASE 0x80000000 // First kernel virtual
   address
11 //从80000000地方开始，可以参考32位机器的elf格式
12 #define KERLINK (KERNBASE+EXTMEM) // Address where kernel is
   linked
13
14 //下面的都是虚拟地址和物理地址之间的转换
15 //v2p 虚拟地址转换为物理地址
16 //p2v 物理地址转换为虚拟地址
17 //为什么要这么转换？
18 //启动分页了，但是页表还没有设置好，所以这个时候虚拟地址和物理地址之间的转换
   是线性的转换.....
19 #ifndef __ASSEMBLER__
20
21 static inline uint v2p(void *a) { return ((uint) (a)) -
   KERNBASE; }
22 static inline void *p2v(uint a) { return (void *) ((a) +
   KERNBASE); }
23
24 #endif
25
26 #define V2P(a) (((uint) (a)) - KERNBASE)
27 #define P2V(a) (((void *) (a)) + KERNBASE)
28
29 #define V2P_WO(x) ((x) - KERNBASE) // same as V2P, but
   without casts
30 #define P2V_WO(x) ((x) + KERNBASE) // same as V2P, but
   without casts
31
```

## defs.h

定义了各种声明和全局函数，我这里只截取了进程相关的。

```
1 // proc.c
2 //常见的进程的定义函数
3 struct proc* copyproc(struct proc*);
4 void exit(void);
```

```

5  int          fork(void);
6  int          growproc(int); //给进程分配更多的内存资源
7  int          kill(int);
8  void         pinit(void);
9  /*
10 //初始化
11 void
12 pinit(void)
13 {
14     //主要是锁初始化
15     initlock(&ptable.lock, "ptable");
16 }
17 */
18 void         procdump(void); //将所有的进程打印出来（避免死锁的打
    印）
19 void         scheduler(void) __attribute__((noreturn)); //调
    度函数
20 void         sched(void);
21 void         sleep(void*, struct spinlock*);
22 void         userinit(void); //第一个进程设置好
23 int          wait(void);
24 void         wakeup(void*);
25 void         yield(void); //从running状态跳转到runnable状态
26

```

## x86.h

```

1  // Routines to let C code use special x86 instructions.
2  // 采用了内嵌汇编的方式 让C语言能够使用汇编
3
4  //介绍一些内嵌汇编的基本格式
5  //第一个：后面的是输出，第二个：后面的是输入
6  //第一个：前面的是指令
7
8  //涉及到asm中的in指令和out指令
9  //in指令读外设端口内容
10 //out指令向外设端口输出
11 //%1代表edx=port %0代表eax=data
12 static inline uchar
13 inb(ushort port)
14 {
15     uchar data;
16
17     asm volatile("in %1,%0" : "=a" (data) : "d" (port));
18     return data;
19 }
20
21 static inline void
22 insl(int port, void *addr, int cnt)
23 {
24     asm volatile("cld; rep insl" :

```

```

25         "=D" (addr), "=c" (cnt) :
26         "d" (port), "0" (addr), "1" (cnt) :
27         "memory", "cc");
28     }
29
30     static inline void
31     outb(ushort port, uchar data)
32     {
33         asm volatile("out %0,%1" : : "a" (data), "d" (port));
34     }
35
36     static inline void
37     outw(ushort port, ushort data)
38     {
39         asm volatile("out %0,%1" : : "a" (data), "d" (port));
40     }
41
42     static inline void
43     outsl(int port, const void *addr, int cnt)
44     {
45         asm volatile("cld; rep outsl" :
46             "=S" (addr), "=c" (cnt) :
47             "d" (port), "0" (addr), "1" (cnt) :
48             "cc");
49     }
50
51     //这个代码用来写入地址用的，就是修改addr的值
52     //https://docs.microsoft.com/zh-cn/cpp/intrinsics/stosb?
53     view=vs-2015
54     static inline void
55     stosb(void *addr, int data, int cnt)
56     {
57         asm volatile("cld; rep stosb" :
58             "=D" (addr), "=c" (cnt) :
59             "0" (addr), "1" (cnt), "a" (data) :
60             "memory", "cc");//cc c compiler
61     }
62
63     static inline void
64     stosl(void *addr, int data, int cnt)
65     {
66         asm volatile("cld; rep stosl" :
67             "=D" (addr), "=c" (cnt) :
68             "0" (addr), "1" (cnt), "a" (data) :
69             "memory", "cc");
70     }
71
72     struct segdesc;//一个数据结构，定义在mmu.h中
73     //各种段的描述符数据结构
74
75     //加载全局/中断描述符表格寄存器
76     //load

```

```

76 static inline void
77 lgdt(struct segdesc *p, int size)
78 {
79     volatile ushort pd[3];
80
81     pd[0] = size-1;
82     pd[1] = (uint)p;
83     pd[2] = (uint)p >> 16;
84
85     asm volatile("lgdt (%0)" : : "r" (pd));
86 }
87
88 struct gatedesc;
89
90 static inline void
91 lidt(struct gatedesc *p, int size)
92 {
93     volatile ushort pd[3];
94
95     pd[0] = size-1;
96     pd[1] = (uint)p;
97     pd[2] = (uint)p >> 16;
98
99     asm volatile("lidt (%0)" : : "r" (pd));
100 }
101
102 //装载任务状态段寄存器TR
103 //在任务内发生特权级变换时堆栈也随着自动切换，外层堆栈指针保存在内层堆栈
    中，而内层堆栈指针存放在当前任务的TSS中。
104 // 所以，在从外层向内层变换时，要访问TSS(从内层向外层转移时不需要访问
    TSS，而只需内层栈中保存的栈指针)。
105
106 static inline void
107 ltr(ushort sel)
108 {
109     asm volatile("ltr %0" : : "r" (sel));
110 }
111
112 static inline uint
113 readeflags(void)
114 {
115     uint eflags;
116     asm volatile("pushfl; popl %0" : "=r" (eflags));
117     return eflags;
118 }
119
120 //linux内核用于GS访问cpu特定的内存
121 //GS段寄存器
122
123 static inline void
124 loadgs(ushort v)
125 {

```

```

126     asm volatile("movw %0, %%gs" : : "r" (v));
127 }
128
129 //CLI汇编指令全称为Clear Interrupt，该指令的作用是禁止中断发生。
130 //
131 //STI汇编指令全称为Set Interrupt，该指令的作用是允许中断发生。
132
133 static inline void
134 cli(void)
135 {
136     asm volatile("cli");
137 }
138
139 static inline void
140 sti(void)
141 {
142     asm volatile("sti");
143 }
144
145 //交换两个寄存器的值
146
147 static inline uint
148 xchg(volatile uint *addr, uint newval)
149 {
150     uint result;
151
152     // The + in "+m" denotes a read-modify-write operand.
153     asm volatile("lock; xchgl %0, %1" :
154                 "+m" (*addr), "=a" (result) :
155                 "1" (newval) :
156                 "cc");
157     return result;
158 }
159
160 //读CR2寄存器
161 //CR2是页故障线性地址寄存器，保存最后一次出现页故障的全32位线性地址。
162
163 static inline uint
164 rcr2(void)
165 {
166     uint val;
167     asm volatile("movl %%cr2,%0" : "=r" (val));
168     return val;
169 }
170 //CR3是页目录基址寄存器，保存页目录表的物理地址，页目录表总是放在以4K字节
    为单位的存储器边界上，因此，它的地址的低12位总为0
171 /*
172  * //CR0中包含了6个预定义标志，0位是保护允许位PE(Protected
    Enable)，用于启动保护模式，
173  * 如果PE位置1，则保护模式启动，如果PE=0，则在实模式下运行。
174  * 1位是监控协处理位MP(Monitor coprocessor)，它与第3位一起决定：当
    TS=1时操作码WAIT是否产生一个“协处理器不能使用”的出错信号。

```

```

175  * 第3位是任务转换位(Task Switch), 当一个任务转换完成之后, 自动将它置
    1. 随着TS=1, 就不能使用协处理器。
176  * CR0的第2位是模拟协处理器位 EM (Emulate coprocessor), 如果EM=1,
    则不能使用协处理器, 如果EM=0, 则允许使用协处理器。
177  * 第4位是微处理器的扩展类型位ET(Processor Extension Type), 其内保存
    着处理器扩展类型的信息,
178  * 如果ET=0, 则标识系统使用的是287协处理器, 如果 ET=1, 则表示系统使用的
    是387浮点协处理器。
179  * //////////CR0的第31位是分页允许位(Paging Enable), 它表示芯片上的分
    页部件是否允许工作。
180  */
181 static inline void
182 lcr3(uint val)
183 {
184     asm volatile("movl %0,%%cr3" : : "r" (val));
185 }
186
187 //PAGEBREAK: 36
188 // Layout of the trap frame built on the stack by the
189 // hardware and by trapasm.S, and passed to trap().
190
191 //用户态陷入内核态的数据结构
192
193 struct trapframe {
194     // registers as pushed by pusha
195     uint edi;
196     uint esi;
197     uint ebp;
198     uint oesp;      // useless & ignored
199     uint ebx;
200     uint edx;
201     uint ecx;
202     uint eax;
203
204     // rest of trap frame
205     ushort gs;
206     ushort padding1;
207     ushort fs;
208     ushort padding2;
209     ushort es;
210     ushort padding3;
211     ushort ds;
212     ushort padding4;
213     uint trapno;
214
215     // below here defined by x86 hardware
216     uint err;
217     uint eip;
218     ushort cs;
219     ushort padding5;
220     uint eflags;
221

```

```

222 // below here only when crossing rings, such as from user
    to kernel
223 uint esp;
224 ushort ss;
225 ushort padding6;
226 };
227

```

## asm.h

这个基本上完全看不懂，只看懂了边界地址一定是4K的整数倍.....

```

1 //
2 // assembler macros to create x86 segments
3 //
4
5 //汇编程序宏来创建x86段
6
7 // 汇编的宏定义
8
9 #define SEG_NULLASM
10 \
11     .word 0, 0;
12 \
13     .byte 0, 0, 0, 0
14
15 // The 0xC0 means the limit is in 4096-byte units
16 //4096字节正好就是4k，无论是虚拟地址空间还是物理内存，都是按照4k划分的
17 // and (for executable segments) 32-bit mode.
18
19 //find in path 全局搜索
20 //生成elf中各种段
21 //看不懂.....
22 #define SEG_ASM(type,base,lim)
23 \
24     .word (((lim) >> 12) & 0xffff), ((base) & 0xffff);
25 \
26     .byte (((base) >> 16) & 0xff), (0x90 | (type)),
27 \
28     (0xC0 | (((lim) >> 28) & 0xf)), (((base) >> 24)
29 & 0xff)
30
31 #define STA_X      0x8      // Executable segment
32 #define STA_E      0x4      // Expand down (non-executable
33 segments)
34 #define STA_C      0x4      // Conforming code segment
35 (executable only)
36 #define STA_W      0x2      // Writeable (non-executable
37 segments)
38 #define STA_R      0x2      // Readable (executable segments)
39 #define STA_A      0x1      // Accessed

```



## mmu.h

定义了一些寄存器的值，一些数据结构，x86 memory management。

```

1 // This file contains definitions for the
2 // x86 memory management unit (MMU).
3 // 内存管理单元，进程地址空间详细数据结构
4
5 // Eflags register
6
7 //详细的定义了eflags这个寄存器
8 //32位的 下面正好也32位
9 #define FL_CF          0x00000001    // Carry Flag
10 #define FL_PF          0x00000004    // Parity Flag
11 #define FL_AF          0x00000010    // Auxiliary carry Flag
12 #define FL_ZF          0x00000040    // Zero Flag
13 #define FL_SF          0x00000080    // Sign Flag
14 #define FL_TF          0x00000100    // Trap Flag
15 #define FL_IF          0x00000200    // Interrupt Enable
16 #define FL_DF          0x00000400    // Direction Flag
17 #define FL_OF          0x00000800    // Overflow Flag
18 #define FL_IOPL_MASK    0x00003000    // I/O Privilege Level
19     bitmask
20 #define FL_IOPL_0        0x00000000    // IOPL == 0
21 #define FL_IOPL_1        0x00001000    // IOPL == 1
22 #define FL_IOPL_2        0x00002000    // IOPL == 2
23 #define FL_IOPL_3        0x00003000    // IOPL == 3
24 #define FL_NT          0x00004000    // Nested Task
25 #define FL_RF          0x00010000    // Resume Flag
26 #define FL_VM          0x00020000    // Virtual 8086 mode
27 #define FL_AC          0x00040000    // Alignment Check
28 #define FL_VIF          0x00080000    // Virtual Interrupt
29     Flag
30 #define FL_VIP          0x00100000    // Virtual Interrupt
31     Pending
32 #define FL_ID          0x00200000    // ID flag
33
34 // Control Register flags
35 //定义了32位的CR0寄存器
36 //CR0中含有控制处理器操作模式和状态的系统控制标志
37 #define CR0_PE          0x00000001    // Protection Enable
38 #define CR0_MP          0x00000002    // Monitor coProcessor
39 #define CR0_EM          0x00000004    // Emulation
40 #define CR0_TS          0x00000008    // Task Switched
41 #define CR0_ET          0x00000010    // Extension Type
42 #define CR0_NE          0x00000020    // Numeric Error
43 #define CR0_WP          0x00010000    // Write Protect
44 #define CR0_AM          0x00040000    // Alignment Mask
45 #define CR0_NW          0x20000000    // Not Writethrough

```

```

43 #define CR0_CD          0x40000000    // Cache Disable
44 #define CR0_PG          0x80000000    // Paging
45
46 #define CR4_PSE          0x00000010    // Page size extension
47
48
49 //定义了elf格式中的各种段
50 #define SEG_KCODE 1    // kernel code
51 #define SEG_KDATA 2    // kernel data+stack
52 #define SEG_KCPU 3    // kernel per-cpu data
53 #define SEG_UCODE 4    // user code
54 #define SEG_UDATA 5    // user data+stack
55 #define SEG_TSS 6    // this process's task state
56
57 //PAGEBREAK!
58 #ifndef __ASSEMBLER__ // if not define
59 // Segment Descriptor 段描述符
60 //看到这个数据结构，我终于对32位机器上的地址生成有了直观的感受.....
61 struct segdesc {
62     uint lim_15_0 : 16; // Low bits of segment limit
63     uint base_15_0 : 16; // Low bits of segment base address
64     uint base_23_16 : 8; // Middle bits of segment base address
65     uint type : 4; // Segment type (see STS_ constants)
66     uint s : 1; // 0 = system, 1 = application
67     uint dpl : 2; // Descriptor Privilege Level
68     uint p : 1; // Present
69     uint lim_19_16 : 4; // High bits of segment limit
70     uint avl : 1; // Unused (available for software use)
71     uint rsv1 : 1; // Reserved
72     uint db : 1; // 0 = 16-bit segment, 1 = 32-bit
73     segment
74     uint g : 1; // Granularity: limit scaled by 4K when
75     set
76     uint base_31_24 : 8; // High bits of segment base address
77 };
78
79

```

## elf.h

看名字就知道了，ELF的头文件。

```

1 // Format of an ELF executable file
2 // 定义了可执行文件的ELF格式
3
4 #define ELF_MAGIC 0x464C457FU // "\x7FELF" in little endian
5 //elf中叫做魔数的概念
6 //最开始的4个字节是所有ELF文件都必须相同的标识码，分别为0x7F、0x45、
7 //0x4c、0x46
8 //为什么要倒着写？小端模式.....
9 // File header

```

```

10 struct elfhdr {
11     uint magic; // must equal ELF_MAGIC
12     uchar elf[12];
13     ushort type;
14     ushort machine;
15     uint version;
16     uint entry;
17     uint phoff;
18     uint shoff;
19     uint flags;
20     ushort ehsize;
21     ushort phentsize;
22     ushort phnum;
23     ushort shentsize;
24     ushort shnum;
25     ushort shstrndx;
26 };
27
28 // Program section header
29
30 //定义的这两个数据结构是有差别的
31 //https://baike.baidu.com/pic/ELF/7120560/0/279759ee3d6d55fb68d
32 //f1ab16f224f4a20a4dd08?fr=lemma&ct=single
33 struct proghdr {
34     uint type;
35     uint off;
36     uint vaddr;
37     uint paddr;
38     uint filesz;
39     uint memsz;
40     uint flags;
41     uint align;
42 };
43
44 // values for Proghdr type
45 #define ELF_PROG_LOAD 1
46
47 // Flag bits for Proghdr flags
48 //定义了一些标志位
49 #define ELF_PROG_FLAG_EXEC 1
50 #define ELF_PROG_FLAG_WRITE 2
51 #define ELF_PROG_FLAG_READ 4

```

## proc.h

```

1 //声明了cpu、进程、进程上下文等数据结构
2 //声明了proc.c中要使用的数据结构
3 //proc.h中没有实现线程结构
4 // Segments in proc->gdt.
5 #define NSEGS 7
6

```

```

7 // Per-CPU state
8 //定义了一个CPU结构
9 struct cpu {
10     uchar id;                // Local APIC ID; index into
                                // cpus[] below//extern struct cpu cpus[NCPU];//索引号
11     struct context *scheduler; // switch() here to enter
                                // scheduler//调度的时候保存上下文环境
12     struct taskstate ts;      // Used by x86 to find stack for
                                // interrupt//任务状态 用在完成中断之后恢复自己的状态
13     struct segdesc gdt[NSEGS]; // x86 global descriptor table//
                                // 全局描述符表.....
14     volatile uint started;    // Has the CPU started?
15     //https://blog.csdn.net/qq_25426415/article/details/54631192
16     //关闭终端，记录关中断的次数，与锁有关.....
17     int ncli;                 // Depth of pushcli nesting.
18     int intena;               // Were interrupts enabled
                                // before pushcli?
19
20     // Cpu-local storage variables; see below
21     struct cpu *cpu; //定义了一个自己的指针
22     //当前跑在这个CPU上面的进程
23     struct proc *proc;        // The currently-running
                                // process.
24 };
25
26 //所以说，xv6是可用的，它支持多CPU
27 extern struct cpu cpus[NCPU]; //结构数组
28 extern int ncpu; //定义了CPU的个数
29
30 // Per-CPU variables, holding pointers to the
31 // current cpu and to the current process.
32 // The asm suffix tells gcc to use "%gs:0" to refer to cpu
33 // and "%gs:4" to refer to proc. seginit sets up the
34 // %gs segment register so that %gs refers to the memory
35 // holding those two variables in the local cpu's struct cpu.
36 // This is similar to how thread-local variables are
37 // implemented
38 // in thread libraries such as Linux pthreads.
39 //通过指针cpu和proc，能够准确引用当前所在CPU的cpu结构体
40 //https://blog.csdn.net/qq_25426415/article/details/54619488
41 extern struct cpu *cpu asm("%gs:0"); // &cpus[cpunum()]
42 extern struct proc *proc asm("%gs:4"); //
43 //cpus[cpunum()].proc
44
45 //PAGEBREAK: 17
46 // Saved registers for kernel context switches.
47 // Don't need to save all the segment registers (%cs, etc),
48 // because they are constant across kernel contexts.
49 // Don't need to save %eax, %ecx, %edx, because the
50 // x86 convention is that the caller has saved them.
51 // Contexts are stored at the bottom of the stack they
52 // describe; the stack pointer is the address of the context.

```

```

51 // The layout of the context matches the layout of the stack in
    swtch.S
52 // at the "Switch stacks" comment. Switch doesn't save eip
    explicitly,
53 // but it is on the stack and allocproc() manipulates it.
54 //保存上下文的寄存器 在上下文切换的时候用得着
55 //为内核上下文切换保存的寄存器。
56 ////不需要保存所有的段寄存器
57 struct context {
58     uint edi;
59     uint esi;
60     uint ebx;
61     uint ebp;
62     uint eip;
63 };
64
65 //进程的六种状态
66 //未使用, 胚胎(初期), 睡眠, 可运行, 正在运行, 僵尸
67 enum procstate { UNUSED, EMBRYO, SLEEPING, RUNNABLE, RUNNING,
    ZOMBIE };
68
69 // Per-process state
70 //定义了进程的结构
71 struct proc {
72     uint sz; // Size of process memory
    (bytes)//进程内存的大小(字节)
73     pde_t* pgdir; // Page table//页表的指针
74     char *kstack; // Bottom of kernel stack for
    this process//内核的最低位置在哪里
75     enum procstate state; // Process state //进程的状态 六种
    之一
76     volatile int pid; // Process ID //进程ID
77     struct proc *parent; // Parent process //父进程
78     struct trapframe *tf; // Trap frame for current
    syscall //当该进程启动系统调用的时候应该保存的信息
79     struct context *context; // swtch() here to run process
    //进程切换应该保存的寄存器信息
80     void *chan; // If non-zero, sleeping on chan
    //阻塞位的标志
81     int killed; // If non-zero, have been killed
    //被kill的标志位
82     struct file *ofile[NOFILE]; // Open files //打开的文件表
83     struct inode *cwd; // Current directory //当前目录
84     char name[16]; // Process name (debugging) //程
    序名字
85 };
86
87 //讲了一下进程地址空间
88 // Process memory is laid out contiguously, low addresses
    first:
89 // text
90 // original data and bss

```

```
91 //    fixed-size stack
92 //    expandable heap
93
```

## proc.c

```
1  //一些进程创建 退出 等待的具体函数
2  #include "types.h"
3  #include "defs.h"
4  #include "param.h"
5  #include "memlayout.h"
6  #include "mmu.h"
7  #include "x86.h"
8  #include "proc.h"
9  #include "spinlock.h"
10
11 struct {
12     struct spinlock lock; //互斥锁
13     struct proc proc[NPROC]; //NPROC=64
14 } ptable; //进程表
15 //进程索引表，64个进程以数组的形式记录
16
17 static struct proc *initproc; //初始进程，userinit调用中赋值
18
19 int nextpid = 1;
20 //forkret函数的函数体在后面
21 extern void forkret(void);
22 //陷入返回
23 extern void trapret(void);
24
25 static void wakeup1(void *chan);
26
27 //初始化
28 void
29 pinit(void)
30 {
31     //主要是锁初始化
32     initlock(&ptable.lock, "ptable");
33 }
34
35 //PAGEBREAK: 32
36 // Look in the process table for an UNUSED proc.
37 // If found, change state to EMBRYO and initialize
38 // state required to run in the kernel.
39 // Otherwise return 0.
40 static struct proc*
41 allocproc(void)
42 {
43     struct proc *p;
44     char *sp;
45
```

```

46 //锁住&ptable
47 acquire(&ptable.lock);
48 for(p = ptable.proc; p < &ptable.proc[NPROC]; p++)
49     if(p->state == UNUSED)//找到未被使用的进程
50         goto found;
51 release(&ptable.lock);//解锁
52 //如果没有找到，连一个unused的进程都找不到
53 return 0;//返回0
54
55 found:
56 p->state = EMBRYO;//将未被使用的进程设置为胚胎模式，初期模式.....反正
    就这么叫.....
57 p->pid = nextpid++;
58 release(&ptable.lock);
59
60 // Allocate kernel stack.分配内核栈
61 if((p->kstack = kalloc()) == 0){
62     /*kalloc()函数的作用是分配一个4096B大小的空间，而在param.h中
    规定了一个进程的内核栈大小为4096*/
63     p->state = UNUSED;//分配失败，状态变回UNUSED
64     return 0;
65 }
66 sp = p->kstack + KSTACKSIZE;//栈指针移动，因为已经分配好了内核栈
67
68 // Leave room for trap frame.
69 //不理解.....
70 //? ? ?
71 sp -= sizeof *p->tf;
72 p->tf = (struct trapframe*)sp;
73
74 // Set up new context to start executing at forkret,
75 // which returns to trapret.
76 sp -= 4;
77 *(uint*)sp = (uint)trapret;
78
79 sp -= sizeof *p->context;
80 p->context = (struct context*)sp;
81 memset(p->context, 0, sizeof *p->context);
82 p->context->eip = (uint)forkret;
83
84 return p;
85 }
86
87 //PAGEBREAK: 32
88 // Set up first user process.
89 //设置第一个用户进程
90 void
91 userinit(void)
92 {
93     struct proc *p;
94     extern char _binary_initcode_start[],
        _binary_initcode_size[];

```

```

95
96     p = allocproc(); //设置好了内核应该设置的东西
97     initproc = p; //将initproc作为第一个进程，所有进程的父亲
98     if((p->pgdir = setupkvm()) == 0)
99         panic("userinit: out of memory?");
100     inituvm(p->pgdir, _binary_initcode_start,
(int)_binary_initcode_size); //分配页表空间
101     p->sz = PGSIZE;
102     memset(p->tf, 0, sizeof(*p->tf));
103     p->tf->cs = (SEG_UCODE << 3) | DPL_USER; //做DPL的分配 CS段中有
DPL=3 代表用户态
104     p->tf->ds = (SEG_UDATA << 3) | DPL_USER;
105     p->tf->es = p->tf->ds;
106     p->tf->ss = p->tf->ds;
107     p->tf->eflags = FL_IF;
108     p->tf->esp = PGSIZE;
109     p->tf->eip = 0; // beginning of initcode.S
110
111     safestrcpy(p->name, "initcode", sizeof(p->name));
112     p->cwd = namei("/"); //分配目录 cwd Current directory 当前分
配的是根目录 /
113
114     p->state = RUNNABLE;
115 }
116
117 // Grow current process's memory by n bytes.
118 // Return 0 on success, -1 on failure.
119 //给进程空间增加内存
120 int
121 growproc(int n)
122 {
123     uint sz;
124
125     sz = proc->sz;
126     if(n > 0){
127         if((sz = allocuvm(proc->pgdir, sz, sz + n)) == 0)
128             return -1;
129     } else if(n < 0){
130         if((sz = deallocuvm(proc->pgdir, sz, sz + n)) == 0)
131             return -1;
132     }
133     proc->sz = sz;
134     switchuvm(proc);
135     return 0;
136 }
137
138 // Create a new process copying p as the parent.
139 // Sets up stack to return as if from system call.
140 // Caller must set state of returned proc to RUNNABLE.
141
142 //创建一个新进程，将p复制为父进程。
143 //设置堆栈以使其好像从系统调用中返回一样。

```



```

144 //调用者必须将返回的proc的状态设置为RUNNABLE。
145 int
146 fork(void)
147 {
148     int i, pid;
149     struct proc *np;
150
151     // Allocate process.
152     if((np = allocproc()) == 0)//allocproc
153         return -1;
154
155     // Copy process state from p.
156     if((np->pgdir = copyuvm(proc->pgdir, proc->sz)) == 0){
157         kfree(np->kstack);
158         np->kstack = 0;
159         np->state = UNUSED;
160         return -1;
161     }
162     np->sz = proc->sz;
163     np->parent = proc;//父子关系
164     *np->tf = *proc->tf;////理解为所有进程公用一个内核栈tf（我不太理
解.....）
165
166     // Clear %eax so that fork returns 0 in the child.
167     np->tf->eax = 0;
168
169     for(i = 0; i < NOFILE; i++)
170         if(proc->ofile[i])
171             np->ofile[i] = filedup(proc->ofile[i]);//复制打开的文件表
172     np->cwd = idup(proc->cwd);
173
174     pid = np->pid;
175     np->state = RUNNABLE;//fork出来的子进程 状态设置为可运行
176     safestrcpy(np->name, proc->name, sizeof(proc->name));
177     return pid;
178 }
179
180 // Exit the current process. Does not return.
181 // An exited process remains in the zombie state
182 // until its parent calls wait() to find out it exited.
183 void
184 exit(void)
185 {
186     struct proc *p;
187     int fd;
188
189     if(proc == initproc)//如果是initproc被exit painc
190         /*void
191         panic(char *s)
192         {
193             printf(2, "%s\n", s);
194             exit();

```

```

195     */
196     panic("init exiting");
197
198     // close all open files.
199     //关闭所有的文件
200     for(fd = 0; fd < NOFILE; fd++){
201         if(proc->ofile[fd]){
202             fclose(proc->ofile[fd]);
203             proc->ofile[fd] = 0;
204         }
205     }
206
207     iput(proc->cwd);
208     proc->cwd = 0;
209
210     acquire(&ptable.lock);
211
212     // Parent might be sleeping in wait().
213     //唤醒它的父进程
214     wakeup1(proc->parent);
215
216     // Pass abandoned children to init.
217     //将它的子进程给init，防止孤儿进程，僵尸进程的出现
218     for(p = ptable.proc; p < &ptable.proc[NPROC]; p++){
219         if(p->parent == proc){
220             p->parent = initproc;
221             if(p->state == ZOMBIE)
222                 wakeup1(initproc);
223         }
224     }
225
226     // Jump into the scheduler, never to return.
227     //将自己设置为僵尸进程
228     proc->state = ZOMBIE;
229     //exit之后就到了调度时机
230     sched();
231     panic("zombie exit");
232 }
233
234 // wait for a child process to exit and return its pid.
235 // Return -1 if this process has no children.
236 //父进程等待，等待子进程的pid
237 //扫描所有进程，如果进程没有子进程，则返回 - 1；如果有子进程且找到一个已
    进入
238 //zombie状态，则先对其资源进行回收，然后返回其pid；如果有子进程但都没有
    进入zombie状
239 //态，则进入sleeping状态。其中：
240 int
241 wait(void)
242 {
243     struct proc *p;
244     int havekids, pid;

```

```

245
246     acquire(&ptable.lock); //锁住
247     for(;;){
248         // Scan through table looking for zombie children.
249         havekids = 0;
250         for(p = ptable.proc; p < &ptable.proc[NPROC]; p++){
251             if(p->parent != proc)
252                 continue;
253             havekids = 1;
254             if(p->state == ZOMBIE){
255                 // Found one.
256                 pid = p->pid;
257                 kfree(p->kstack);
258                 p->kstack = 0;
259                 freevm(p->pgdir);
260                 p->state = UNUSED;
261                 p->pid = 0;
262                 p->parent = 0;
263                 p->name[0] = 0;
264                 p->killed = 0;
265                 release(&ptable.lock);
266                 return pid;
267             }
268         }
269
270         // No point waiting if we don't have any children.
271         if(!havekids || proc->killed){
272             release(&ptable.lock);
273             return -1;
274         }
275
276         // Wait for children to exit. (See wakeup1 call in
proc_exit.)
277         sleep(proc, &ptable.lock); //DOC: wait-sleep
278     }
279 }
280
281 //PAGEBREAK: 42
282 // Per-CPU process scheduler.
283 // Each CPU calls scheduler() after setting itself up.
284 // Scheduler never returns. It loops, doing:
285 //  - choose a process to run
286 //  - swtch to start running that process
287 //  - eventually that process transfers control
288 //    via swtch back to the scheduler.
289 // -选择要运行的进程
290 // -swtch开始运行该进程
291 // -最终该过程转移了控制权
292 //通过swtch返回调度程序。
293 void
294 scheduler(void)
295 {

```

```

296     struct proc *p;
297
298     for(;;){
299         // Enable interrupts on this processor.
300         sti();
301
302         // Loop over process table looking for process to run.
303         acquire(&ptable.lock);
304         for(p = ptable.proc; p < &ptable.proc[NPROC]; p++){
305             if(p->state != RUNNABLE)//简单的遍历 找到第一个runnable.....这
                么直白的吗.....
306                 continue;
307
308                 // Switch to chosen process. It is the process's job
309                 // to release ptable.lock and then reacquire it
310                 // before jumping back to us.
311                 proc = p;
312                 switchvm(p);
313                 p->state = RUNNING;
314                 swtch(&cpu->scheduler, proc->context);
315                 switchkvm();
316
317                 // Process is done running for now.
318                 // It should have changed its p->state before coming
319                 back.
320                 proc = 0;
321             }
322             release(&ptable.lock);
323         }
324     }
325
326     // Enter scheduler. Must hold only ptable.lock
327     // and have changed proc->state.
328     void
329     sched(void)
330     {
331         int inena;
332
333         if(!holding(&ptable.lock))
334             panic("sched ptable.lock");
335         if(cpu->ncli != 1)
336             panic("sched locks");
337         if(proc->state == RUNNING)
338             panic("sched running");
339         if(readeflags() & FL_IF)
340             panic("sched interruptible");
341         inena = cpu->intena;
342         swtch(&proc->context, cpu->scheduler);
343         cpu->intena = inena;
344     }
345

```

```

346 // Give up the CPU for one scheduling round.
347 //yield过程将running转变成runnable
348 void
349 yield(void)
350 {
351     acquire(&ptable.lock); //DOC: yieldlock
352     proc->state = RUNNABLE;
353     sched(); //CPU空出来了就要执行调度了
354     release(&ptable.lock);
355 }
356
357 // A fork child's very first scheduling by scheduler()
358 // will swtch here. "Return" to user space.
359 void
360 forkret(void)
361 {
362     static int first = 1;
363     // Still holding ptable.lock from scheduler.
364     release(&ptable.lock);
365
366     if (first) {
367         // Some initialization functions must be run in the
368         // context
369         // of a regular process (e.g., they call sleep), and thus
370         // cannot
371         // be run from main().
372         first = 0;
373         initlog();
374     }
375     // Return to "caller", actually trapret (see allocproc).
376 }
377
378 // Atomically release lock and sleep on chan.
379 // Reacquires lock when awakened.
380 void
381 sleep(void *chan, struct spinlock *lk)
382 {
383     if (proc == 0)
384         panic("sleep");
385
386     if (lk == 0)
387         panic("sleep without lk");
388
389     // Must acquire ptable.lock in order to
390     // change p->state and then call sched.
391     // Once we hold ptable.lock, we can be
392     // guaranteed that we won't miss any wakeup
393     // (wakeup runs with ptable.lock locked),
394     // so it's okay to release lk.
395     if (lk != &ptable.lock) { //DOC: sleeplock0
396         acquire(&ptable.lock); //DOC: sleeplock1

```

```

396     release(&lk);
397 }
398
399 // Go to sleep.
400 proc->chan = chan;
401 proc->state = SLEEPING;
402 sched();
403
404 // Tidy up.
405 proc->chan = 0;
406
407 // Reacquire original lock.
408 if(&lk != &ptable.lock){ //DOC: sleeplock2
409     release(&ptable.lock);
410     acquire(&lk);
411 }
412 }
413
414 //PAGEBREAK!
415 // Wake up all processes sleeping on chan.
416 // The ptable lock must be held.
417 static void
418 wakeup1(void *chan)
419 {
420     struct proc *p;
421
422     for(p = ptable.proc; p < &ptable.proc[NPROC]; p++){//很简单的
        循环
423         if(p->state == SLEEPING && p->chan == chan)
424             p->state = RUNNABLE;//改变状态到runnable
425     }
426
427 // Wake up all processes sleeping on chan.
428 //为什么要设置两个这样的函数，wakeup和wakeup1?
429 void
430 wakeup(void *chan)
431 {
432     acquire(&ptable.lock);
433     wakeup1(chan);//调用wakeup1
434     release(&ptable.lock);
435 }
436
437 // Kill the process with the given pid.
438 // Process won't exit until it returns
439 // to user space (see trap in trap.c).
440 int
441 kill(int pid)
442 {
443     struct proc *p;
444
445     acquire(&ptable.lock);//锁
446     for(p = ptable.proc; p < &ptable.proc[NPROC]; p++){

```

```

447     if(p->pid == pid){
448         p->killed = 1;//killed设置为1
449         // wake process from sleep if necessary.
450         //也可以在这里将阻塞进程变成可运行
451         //为什么要在这里设置这样一个功能?
452         //CPU空出来了?
453         if(p->state == SLEEPING)
454             p->state = RUNNABLE;
455         release(&ptable.lock);//解锁
456         return 0;
457     }
458 }
459 release(&ptable.lock);
460 return -1;
461 }
462
463 //PAGEBREAK: 36
464 // Print a process listing to console.  For debugging.
465 // Runs when user types ^P on console.
466 // No lock to avoid wedging a stuck machine further.
467 //将所有的进程打印在控制台上
468 //没有死锁
469 void
470 procdump(void)
471 {
472     static char *states[] = {
473         [UNUSED]    "unused",
474         [EMBRYO]     "embryo",
475         [SLEEPING]   "sleep ",
476         [RUNNABLE]   "runble",
477         [RUNNING]    "run   ",
478         [ZOMBIE]     "zombie"
479     };
480     int i;
481     struct proc *p;
482     char *state;
483     uint pc[10];
484
485     for(p = ptable.proc; p < &ptable.proc[NPROC]; p++){
486         if(p->state == UNUSED)
487             continue;
488         if(p->state >= 0 && p->state < NELEM(states) && states[p-
>state])
489             state = states[p->state];
490         else
491             state = "???";
492         cprintf("%d %s %s", p->pid, state, p->name);
493         if(p->state == SLEEPING){
494             getcallerpcs((uint*)p->context->ebp+2, pc);
495             for(i=0; i<10 && pc[i] != 0; i++)
496                 cprintf(" %p", pc[i]);
497         }

```

```
498     cprintf("\n");
499 }
500 }
501
502
503
```

## swtch.S

```
1  # Context switch
2  //上下文切换
3  #
4  # void swtch(struct context **old, struct context *new);
5  #
6  # Save current register context in old
7  # and then load register context from new.
8
9  .globl swtch
10 swtch:
11     movl 4(%esp), %eax//利用栈来传递传参数，32位机器
12     movl 8(%esp), %edx
13
14     # Save old callee-save registers
15     pushl %ebp
16     pushl %ebx
17     pushl %esi
18     pushl %edi
19
20     # Switch stacks
21     movl %esp, (%eax)//一个打了括号 一个不打括号 修改指针本身的内容 间接寻址
22     movl %edx, %esp
23
24     # Load new callee-save registers
25     popl %edi
26     popl %esi
27     popl %ebx
28     popl %ebp
29     ret
30
```



