

Operating Systems Concepts:

Chapter 6: The System Nucleus

William Knottenbelt (wik@doc.ic.ac.uk)
(after Jeff Magee, Kevin Twidle, Steve Vickers and Ariel Burton)

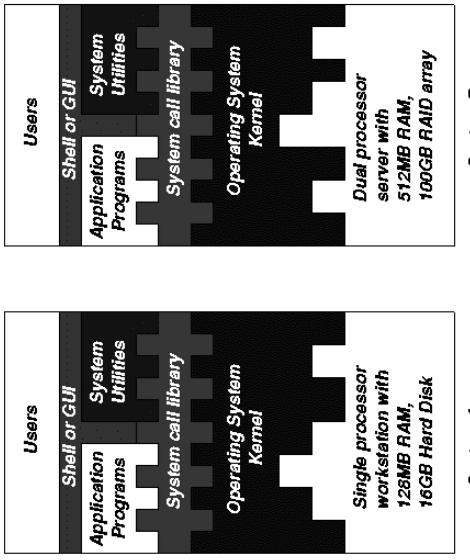
An operating system:

- manages a system's resources so that they are used efficiently and safely, e.g.,
 - CPU(s)
 - memory
 - devices (modems, disks, network interfaces, video interfaces, etc.)
- presents a virtual machine that provides convenient abstractions, e.g.,
 - files rather than disk locations (device independence)
 - inter-process synchronisation and communication.

Hardware + OS = Usable Virtual Machine

1
6 - Simple Kernel

The Operating System as Virtual Machine



2
6 - Simple Kernel

The System Nucleus (Kernel)

- Shares processor among multiple processes
 - scheduling
 - context switching
- Process management
 - creation, deletion, initiation, termination
 - interprocess synchronisation e.g. P & V, signal & wait
- Synchronisation with real-time
- Interrupt handling
- *Monolithic kernels* additionally include higher-level services, e.g. memory management, I/O, filing system, networking etc.
- *Microkernel*/systems delegate higher-level services to servers that run as applications in user mode

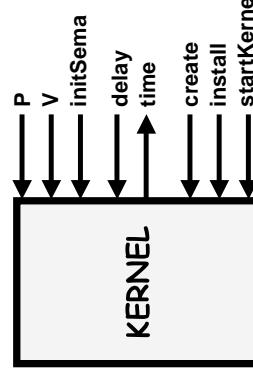
Resources (real source code of kernels):

*Simple Kernel: see <http://www.doc.ic.ac.uk/~wik/OperatingSystemsConcepts>

*Linux: <http://www.tanacom.com/tour/linux/or/usr/src/linux>.

3
6 - Simple Kernel

Simple Kernel Design



Simplifying assumptions:

- 1) Processes are assumed to run forever
- 2) The code for the kernel + user processes is loaded at start-up time

⇒ Suitable for embedded applications.

4
6 - Simple Kernel

Kernel Interface

Today, kernels are usually written in a high-level programming language, with some low-level assembly code.

module "kernel.hpp"

```
create(P, size, pr) creates a new process executing function P with workspace size and priority pr.
/* processes synchronization: */ opaque means "black box" – internal structure is hidden
opaque struct Semaphore; void P(Semaphore *s);
void V(Semaphore *s); initSemaphore(Semaphore *s, int value);

/* synchronization with real time: */
int time();
void delay(int ticks),

/* system set up: */
enum Priority {high, normal, low, idle};
typedef void Proc();
void create(Proc *P, int size, Priority pr);
void install(Proc *P, int intno);
void startKernel();
```

6 - Simple Kernel

time() returns clock ticks since kernel was started.
delay(t) suspends a process for t clock ticks.

install(P, intno) makes function P the handler for interrupt intno, so that P is invoked whenever that interrupt occurs.

startKernel() starts processes running under kernel control.

5

Module "asterisk.cpp"

```
#include "kernel.hpp"
// clock interrupts every 20 ms
const int ticksper10s = 500;
const int clockvector = 32;
int count; // incremented every tick
Semaphore sync;
```

*Two phases of execution –
1: Initialization (creates, installs)
2: StartKernel to run processes.*

```
void tick() {
    count = (count+1) % ticksper10s;
    if (count == 0) {
        V(&sync);
    }
}
```

```
void main() { /* main routine starts here */
    Count = 0;
    initSemaphore(sync, 0)
    install(tick, clockvector);
    create(print, 500, normal); // normal is a priority level
    startKernel();
}

void print() {
    while (1) {
        P(&sync);
        put("**");
    }
}
```

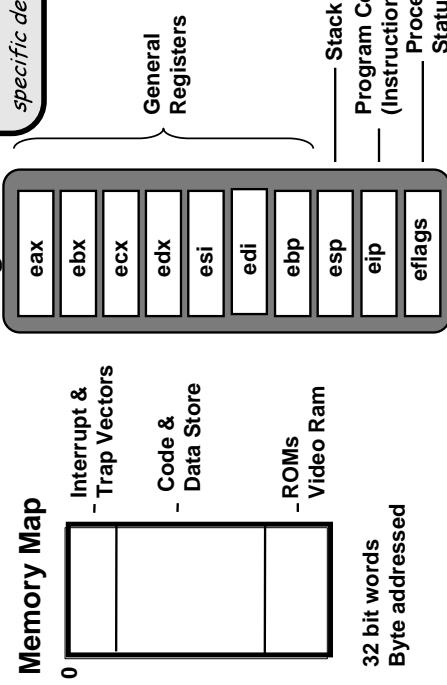
6 - Simple Kernel

6

Target Architecture: Intel 386

You need to understand the general features of this, but not the machine-specific details.

Registers



7

6 - Simple Kernel

Vectorized Interrupts on the 386.

Flags Register

11	10	9	8	7	6	5	4	3	2	1	0
OF	DF	IF	TF	SF	ZF	AF	PF	CF			

IF – Interrupt Flag. When set interrupts are enabled otherwise they are disabled.

```
cli      clear IF (disable interrupts)
sti      set IF (enable interrupts)
pushf   put flags on stack
popf   set flags from top of stack
iret
```

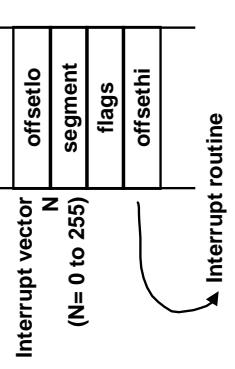
Interrupting device – signals interrupt, puts vector number N (8 bits) on CPU pins:
CPU hardware then – pushes flags and PC (eip) onto the stack;
– sets flag values from the vector;

– calculates new PC (eip) from the segment and offsets ...
... thus effecting a jump to interrupt routine.
At end of interrupt routine, iret instruction pops return PC and flags from stack.

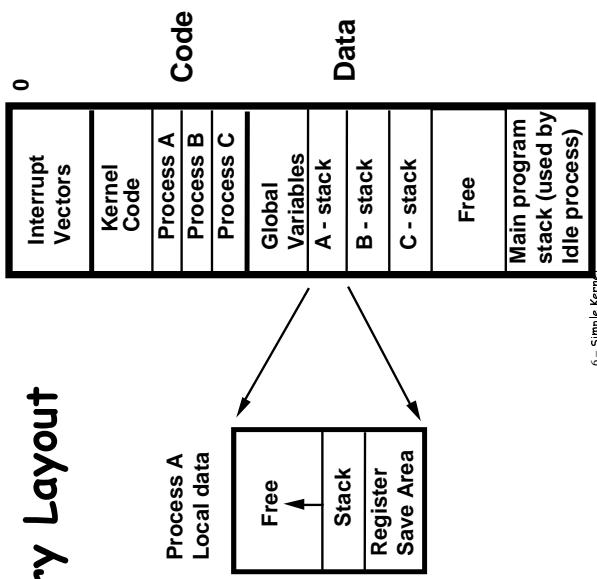
6 - Simple Kernel

8

Vector configuration

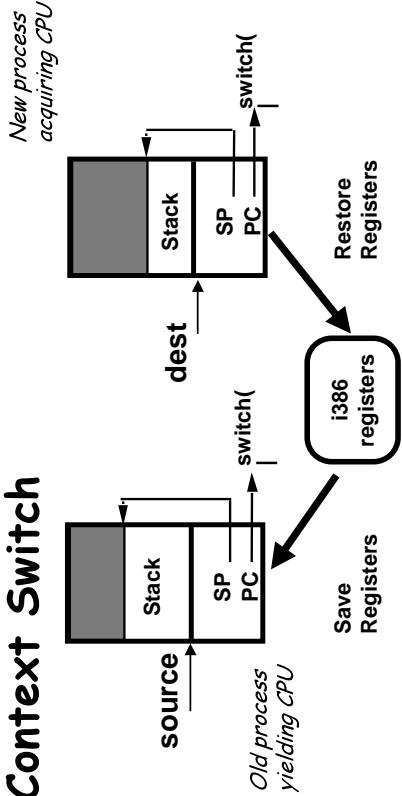


Memory Layout



9

Context Switch



```

New process
acquiring CPU

typedef struct Context {
    int eax, ebx, ecx, edx, esi, edi, ebp, eip, esp;
    /*   0   4   8   12  16  20  24  28  32  36 note these offsets*/
} Context;

void switch(Context *source, Context *dest);
  
```

6 - Simple Kernel

10

Implementation of switch

```

.globl proc_switch
proc_switch:
# (It helps to first see stack layout on next slide)
        # Save current eax
        movl %eax,-4(%esp)          # Get addr of 'source' ctxt
        # Save ebx in 'source' ctxt
        movl %ebx,4(%esp),%eax
        # Get orig. eax into ebx
        movl -4(%esp),%ebx
        # Save orig. eax in 'source' ctxt
        movl %ebx,(%esp)
        # Save other regs
        movl %ecx,8(%esp)
        # (see offsets on previous slide)
        movl %edx,12(%esp)
        movl %esi,16(%esp)
        movl %edi,20(%esp)
        movl %ebp,24(%esp)
        # Only way to get flags
        pushf
        popl %ebx
        movl %ebx,28(%esp)
        # Save them in 'source' ctxt
        # Get return address off stack
        popl %ebx
        movl %ebx,32(%esp)
        # Save as 'source' IP
        movl %esp,36(%esp)
        # Save an esp that has return
        # addr popped to simulate a ret
  
```

6 - Simple Kernel

Implementation of switch cont.)

```

        movl 4(%esp),%eax          # Get addr of 'dest' ctxt
        movl 36(%esp),%esp          # Switch to the 'dest' ctxt stack frame
        # Following 3 pushl's set up the stack
        # so we can use an iret to return to
        # the task (iret is equivalent to ret, pop)
        movl 28(%eax),%ebx
        pushl %ebx                 # Push eflags
        movl $KERNEL_CS,%ebx
        pushl %ebx                 # Push code segment of the 'dest' proc
        movl 32(%eax),%ebx
        pushl %ebx                 # Push execution addr of 'dest' proc
        movl 24(%eax),%ebp          # Now just load the rest of the regs
        movl 20(%eax),%edi          # from the 'dest' ctxt into the regs
        movl 16(%eax),%esi
        movl 12(%eax),%edx
        movl 8(%eax),%ecx
        movl 4(%eax),%ebx
        movl (%eax),%eax
  
```

6 - Simple Kernel

In Linux this routine is called `switch_to`, and is implemented in `arch/i386/kernel/process.c` and include/`asm-i386/system.h`.

6 - Simple Kernel

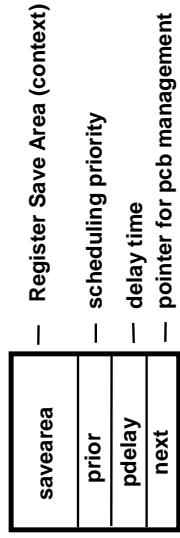
11

6 - Simple Kernel

12

Kernel Representation of Processes

Process Control Block



```
enum Priority {high, normal, low, idle};
typedef PCB *Process;
typedef struct PCB {
    Context savearea;
    Priority prior;
    int pdelay;
    Process next;
} PCB;
```

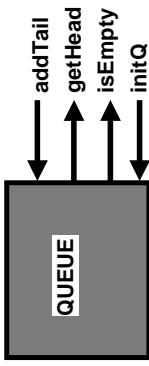
Processes are managed by forming "linked lists" or queues of process control blocks.

In linux, a process is represented by a *struct task_struct*, defined in include/linux/sched.h

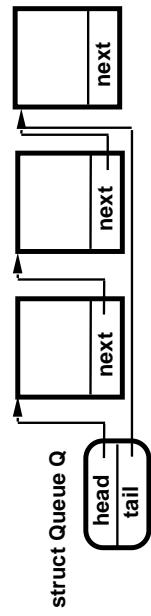
6 - Simple Kernel

13

Kernel FIFO Queue Operations



```
typedef struct Queue {
    Process head, tail;
} Queue;
```



Advantages of this data structure?

14

6 - Simple Kernel

15

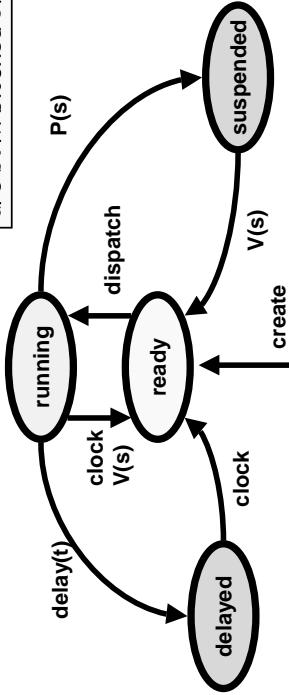
Queue Implementation

```
bool isEmpty(Queue *Q) {
    return (Q->head == NULL);
}

void addTail(Queue *Q, Process p) {
    if (Q->head == NULL)
        Q->head = p;
    else
        Q->tail->next = p;
    Q->tail = p;
    p->next = NULL;
}
```

Delayed and suspended
are both **blocked** states.

Process State Transitions



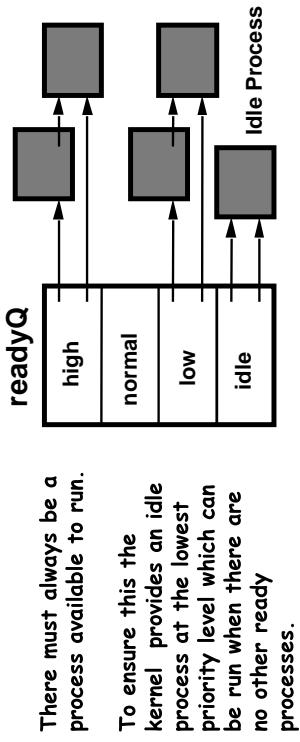
A process will be in exactly one of these states. The kernel maintains data structures to keep track of the processes in each state:
A **ready** process – is on one of the four ready queues, readyQ[priority] (see below)
A **delayed** process – is on delayQ
A **suspended** process – is on the queue associated with a semaphore S
The **running** process – is pointed to by the variable running
(Just one processor, so only one process at a time in the running state.)

6 - Simple Kernel

16

Ready Queue - Static multi-level priority

Actually *four* ready queues, one for each priority.



```
Process running; //points to currently running process
Queue readyQ[4]; //one queue for each priority level
enum Priority {high,normal,low,idle};
```

6 - Simple Kernel

17

Scheduling Operations

Selecting a process to run

```
void setReady(Process p) {
    addTail(readyQ[int(p->prior)], p);
}

void dispatch() {
    Process oldrunning = running;
    Priority pr = high;
    while (!isEmpty(readyQ[int(pr)])) {
        pr = Priority(int(pr) + 1);
        // gets next priority
        running = getHead(readyQ[int(pr)]);
        switch(oldrunning->savearea,
              running->savearea);
    }
}
```

18

6 - Simple Kernel

Kernel Exit & Entry (cont.)

To ensure mutual exclusion to the kernel data structures, kernel procedures must run with interrupts disabled. We will use the functions:

int int_mutexon()
void int_mutexoff()
Assembler:

```
.globl int_mutexon
int_mutexon:
    pushf
    cli
    popl %eax
    ret
```

```
void initSched() {
    for (int p=0; p<4; p++) {
        InitQ(readyQ[int(p)]);
    }
}

int psr = getHead(readyQ[int(idle)]);
create(null, 0, idle);
running = getHead(readyQ[int(idle)]);
}
```

```
ENTERKERNEL
int psw = int_mutexon();

EXITKERNEL
int_mutexoff(psw);

void ... (...)
    int psw = int_mutexon(); // disable interrupts
    :
    int_mutexoff(psw); // restore interrupts to previous state
end ...
```

Why can't int_mutexoff simply enable interrupts?

Hint: what happens if one kernel procedure calls another?

6 - Simple Kernel

19

20

Semaphore Implementation

```
typedef struct Semaphore {
    int count;
    Queue waiting;
} Semaphore;

void P(Semaphore *s) {
    ENTERKERNEL
    if (s->count > 0)
        s->count--;
    else {
        addTail(s->waiting, running);
        dispatch();
    }
    EXITKERNEL
}

void initSem(Semaphore *s, int value) {
    ENTERKERNEL
    s->count = value;
    initQ(s->waiting);
    EXITKERNEL
}
```

6 - Simple Kernel

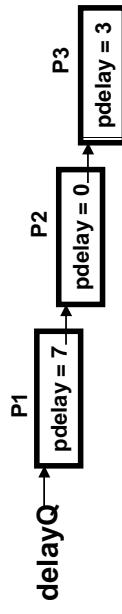
21

```
SEMAPHORE
P
V
initSem
Process delayQ;
int clock;
```

The delayQ is not a FIFO queue. We will use a "delta" queue where the delay time of a process in the queue is the sum of its pdelay field and the pdelay fields of those processes that precede it in the queue.

Example:

```
At time 100 P1 executes delay(7);
At time 100 P2 executes delay(7);
At time 100 P3 executes delay(10);
```



22

6 - Simple Kernel

Time Operations

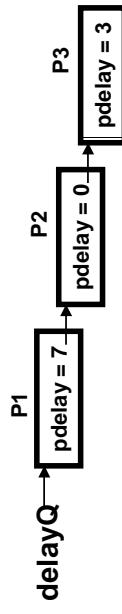
```
TIMEOPS
clock
delayQ
process
tick
initTime
```

```
Process delayQ;
int clock;
```

The delayQ is not a FIFO queue. We will use a "delta" queue where the delay time of a process in the queue is the sum of its pdelay field and the pdelay fields of those processes that precede it in the queue.

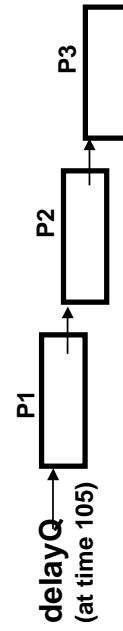
Example:

At time 100 P1 executes delay(7);
At time 100 P2 executes delay(7);
At time 100 P3 executes delay(10);



Time Implementation

```
At time 95 P1 executes delay(20);
At time 100 P2 executes delay(30);
At time 105 P3 executes delay(35);
```



```
void tick() {
    // Update current time
    clock++;
    // wake up any processes whose delays have expired
    if (delayQ != NULL) {
        delayQ->pdelay--;
        Process p = NULL;
        while (delayQ != NULL && delayQ->pdelay == 0) {
            // loop setting ready all processes with
            // pdelay = 0, or until delayQ is empty
            p = delayQ;
            delayQ = p->next;
            setReady(p);
        }
        // reschedule
        setReady(running);
        dispatch();
    }
}
```

Exercise
What is the advantage of this organisation?
Answer: efficiency---only the process at the head of the queue needs to be examined by the interrupt handler on each clock tick (see below).

23

6 - Simple Kernel

24

6 - Simple Kernel

Interrupts

Delay

```
void delay(int t) {
    ENTERKERNEL
    if (t > 0) {
        int t1 = t;
        Process p0 = delayQ, p1 = delayQ;
        while (p1 != NULL && t1 >= p1->pdelay) {
            t1 -= p1->pdelay;
            p0 = p1;
            p1 = p0->next;
            if (p0 == p1) /* == delayQ */
                delayQ = running;
            else
                p0->next = running;
                running->next = p1;
                running->pdelay = t1;
                if (p1 != NULL) {
                    p1->pdelay -= t1;
                    dispatch();
                }
            EXITKERNEL
        }
    }
}
```

*p0 is the current position in delayQ
p1 is next process in delayQ*

If there is another process, p1, after the newly added process, then its pdelay must be adjusted.

Interrupts First level Interrupt Handler, int_flih

- Want to use a high-level language to write interrupt routines.
- Actual "first-level" interrupt handler is general interface for the compiled high-level routines.
- It saves all the registers, pushes the interrupt number onto the stack and calls our interrupt routine.
- When that returns we restore the registers and execute an iret instruction.

```
#define INT_FLIH_CPU(inum) \
\ .align 4 \
\ .global int_flih_##inum \
int_flih_##inum##: \
    Push %es \
    Push %ds \
    Push %eax \
    Push %ebx \
    Push %ecx \
    Push %edx \
    Push %esi \
    Push %edi \
    Push %ebp \
    Push %ebp \
    call int_interrupt \
    addl $4,%esp \
    pop %ebp \
    pop %edi \
    pop %esi \
    pop %edx \
    pop %ecx \
    pop %ebx \
    pop %eax \
    pop %ds \
    pop %es \
    iret
```

6 – Simple Kernel

25

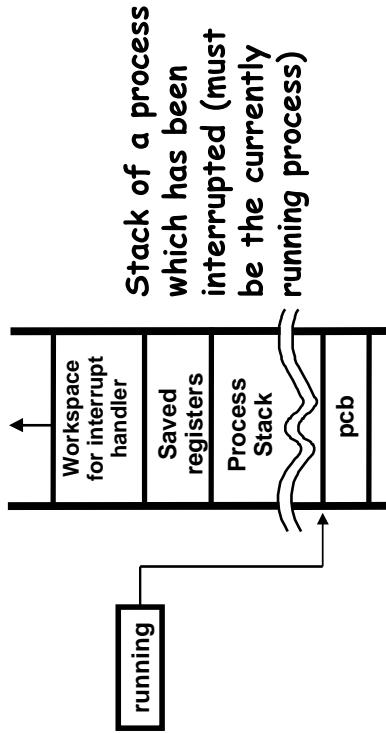
6 – Simple Kernel

26

6 – Simple Kernel

Interrupts (cont.)

- The assembly code routine int_flih also saves and restores the registers since these will be used by the interrupt handler procedure.
- Interrupt handlers are not processes. They use the stack of the currently running process.



Interrupt Operations



```
typedef void Proc();
Proc *handlers[256];
bool started = false;

void startInt() {
    started = true;
}

void install(Proc *P, int vector) {
    ENTERKERNEL
    handlers[vector] = P;
    // must also initialize interrupt vector (machine dependent)
    EXITKERNEL
}
```

```
INTERRUPT OPS
    handlers
    started
```

```
install
intDisp
startInt
```

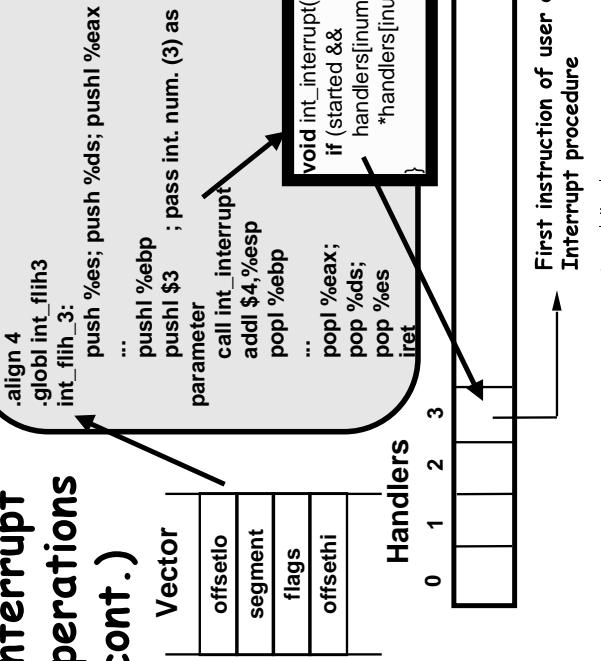
6 – Simple Kernel

27

6 – Simple Kernel

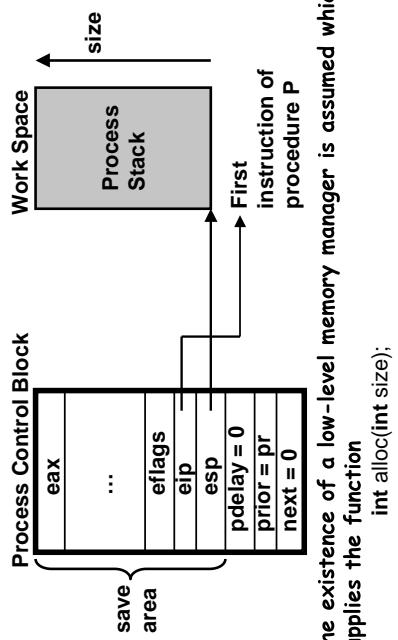
28

Interrupt operations (cont.)



Process Creation

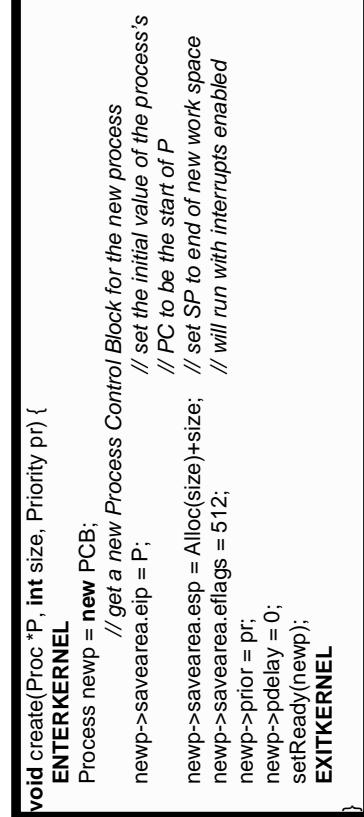
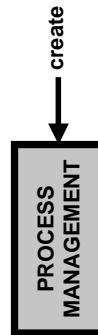
Processes are created in the following state:



6 - Simple Kernel

30

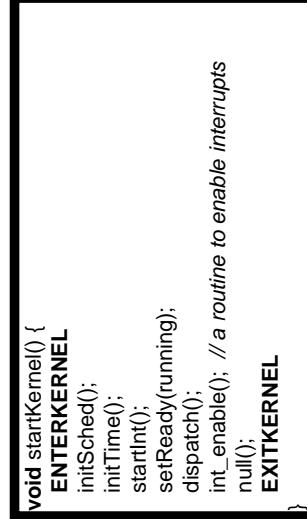
Process Creation Implementation



6 - Simple Kernel

29

Kernel Initialisation



6 - Simple Kernel

32