

## Ausgewählte Betriebssysteme

### Processes and Threads

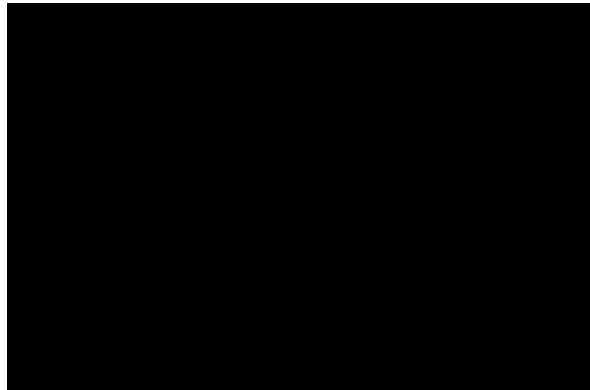
1

## What is a process

- Fundamental concept for multiprogramming
- Instance of program in execution
  - Sequential control flow
- Entity to which system resources are allocated
  - Might be shared among processes (threads)

2

## task\_struct



3

## Process state

- Field in task\_struct
- Currently available
  - TASK\_RUNNING - executing or ready for execution
  - TASK\_INTERRUPTIBLE - suspended
  - TASK\_UNINTERRUPTIBLE - suspend, no signals
  - TASK\_STOPPED - execution has been externally stopped
  - TASK\_ZOMBIE - terminated

4

## Process descriptor handling

- Processes are dynamic entities
  - Dynamic allocation
  - Half of all physical memory might be used for PCB
    - `max_threads = mempages / (THREAD_SIZE/PAGE_SIZE) / 2`
    - `/proc/sys/kernel/threads-max`
- Two different data structures per process
  - Process descriptor
  - Kernel stack

5

## task\_union

```
movl $0xffffe000  
andl $%esp, %ecx  
movl %ecx, p
```

```
union task_union {  
    struct task_struct task;  
    unsigned long stack  
    [2048];  
};
```

6

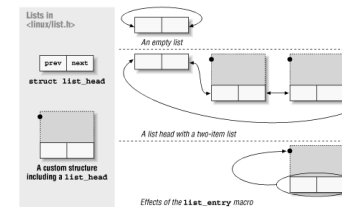
## Process List

- Linux keeps list of processes for different purposes
  - Special properties (e.g. runnable)
- Process List
  - All processes in the system
  - Circular double linked list
  - SET\_LINKS/REMOVE\_LINKS macros ensure consistency
  - next\_task, prev\_task field in task\_struct

7

## Doubly linked lists (implementation)

- Often used
- Reusable implementation
  - Access functions and macros



8

## Run queue

- Scheduling needs only to consider runnable processes
- Linked through `struct list_head run_list`
- Select most viable process to run next

```
|schedule
|do_softirq // manages post-IRQ work
|for each task
|  |calculate counter
|prepare_to_switch // does anything
|switch_mm // change Memory context (change CR3 value)
|switch_to (assembler)
|SAVE ESP
|RESTORE future_ESP
|SAVE EIP
|push future_EIP *** push parameter as we did a call
|  |jmp __switch_to (it does some TSS work)
|  |__switch_to()
|  |
|ret *** ret from call using future_EIP in place of call address
new_task
```

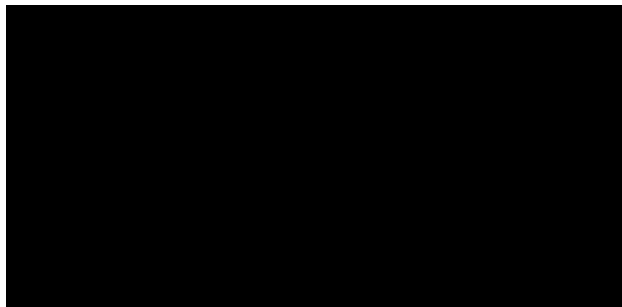
9

## Process identification

- Address of PCB is unique in kernel address space
- PID used at user level
- Process list traversal to slow
- Hash table for fast lookup

10

## PCB linking



11

## Process management

- Process queue
- TASK\_RUNNABLE
  - Run queue
- TASK\_STOPPED, TASK\_ZOMBIE
  - Not grouped
- TASK\_(UN)INTERRUPTIBLE
  - Subdivided into many classes, each of which correspondents to a specific event
  - State alone does not provide enough information
  - Specific lists of processes called wait queues

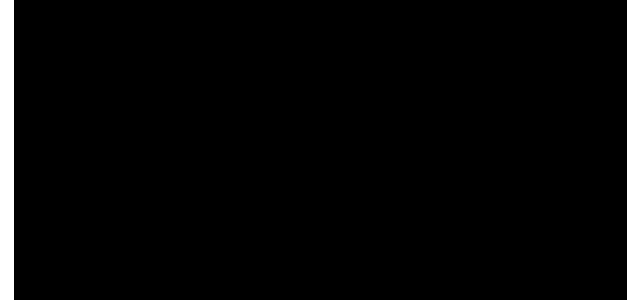
12

## Wait Queues

- Define a new wait queue if needed
  - DECLARE\_WAIT\_QUEUE\_HEAD(...)
- Functions
  - add\_wait\_queue(...), remove\_wait\_queue(...)
  - sleep\_on
  - wake\_up

13

## Wait queue structure



14

## Process creation

- fork syscall
  - Copy process
  - Independent new execution context
- clone syscall
  - Share resources with the new context
  - lightweight

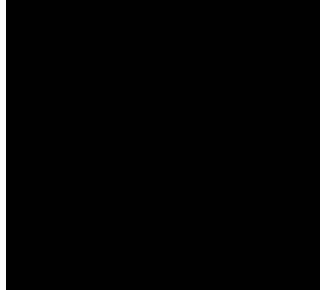
15

## Forking

```
|sys_fork
|do_fork
|alloc_task_struct
|__get_free_pages
p->state = TASK_UNINTERRUPTIBLE
copy_flags
p->pid = get_pid
copy_files
copy_fs
copy_sighand
copy_mm // should manage CopyOnWrite (1 part)
|allocate_mm
|mm_init
|pgd_alloc -> get_pgd_fast
|get_pgd_slow
|dup_mmap
|copy_page_range
|pte_set_wprotect
|clear_bit // set page to read-only
|copy_segments // For LDT
|copy_thread
|childregs->eax = 0
p->thread.esp = childregs // child fork returns 0
p->thread.eip = ret_from_fork // child starts from fork exit
retval = p->pid // parent fork returns child pid
SET_LINKS // insertion of task into the list pointers
nr_threads+ // Global variable
wake_up_process(p) // Now we can wake up just created child
return retval
```

16

## Process relationship



17

## Kernel threads

- Critical tasks implemented as intermittently running processes
  - Flushing disk caches
  - Swapping out unused page frames
- Regular scheduling
  - No unbound kernel activities
- Special characteristics
  - Mostly only one single kernel function
  - No user mode part

18

## Kernel thread creation

```
int kernel_thread(int (*fn)(void *), void * arg, unsigned long flags)
{
    long retval, d0;
    __asm__ __volatile__(
        "movl $temp,$temp\n\t"
        "int $0x80\n\t" /* Linux/i386 system call */
        "cmpl $temp,$temp\n\t" /* child or parent? */
        "je 1f\n\t" /* parent - jump */
        /* Load the argument into eax, and push it. That way, it does
         * not matter whether the called function is compiled with
         * -msregsparm or not... */
        "movl $4,$temp\n\t"
        "pushl $temp\n\t"
        "call *%0\n\t" /* call fn */
        "movl $3,$0\n\t" /* exit */
        "int $0x80\n\t"
        "1:\n\t"
        : "=a" (retval), "=a" (d0)
        : "0" (__NR_clone), "i" (__NR_exit),
          "r" (arg), "r" (fn),
          "b" (flags | CLONE_VM)
        : "memory");
    return retval;
}
```

19

## Kernel threads in action

```
init,1
|-(bdflush,6)
|-(keventd,2)
|-(khubd,53)
|-(kjournald,10)
|-(kjournald,89)
|-(kjournald,90)
|-(kjournald,1969)
|-(ksoftirqd_CPU0,4)
|-(kswapd,5)
|-(kupdated,7)
|-(lockd,19499)
```

20

## Context switch

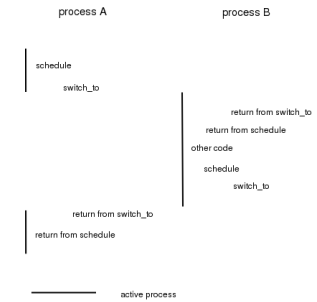
- Transfer control between contexts
  - Save state of current context
  - Load state of next context and resume execution
- Execution context
  - Architectural (user level) cpu state
  - Virtual memory

21

## Context switch (2)

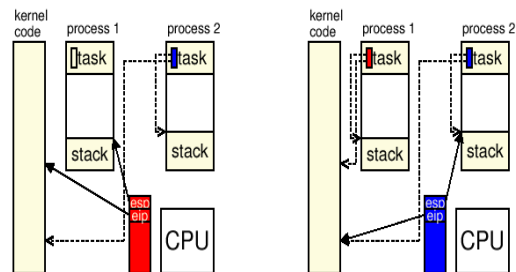
- For the kernel programmer context switching looks like a ordinary function call.
- Interleaved activities of other processes are transparent

```
void schedule(void){
    .
    /* calc next process */
    .
    switch_to(..., next, ...)
    .
    return from switch_to
    return from schedule
}
```



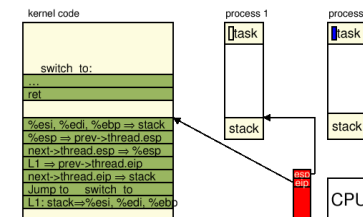
22

## Switch 1



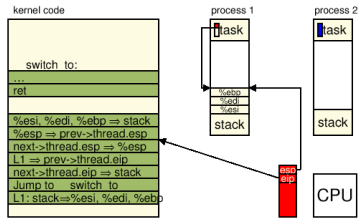
23

## Switch (2)



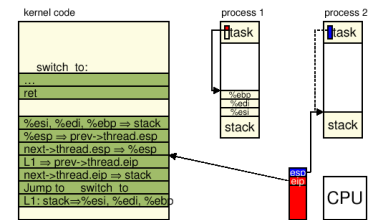
24

### Switch (3)



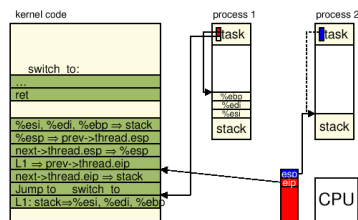
25

### Switch (4)



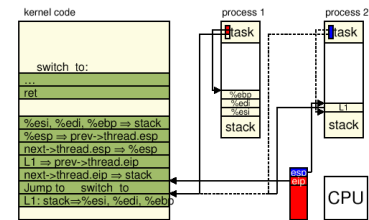
26

### Switch (5)



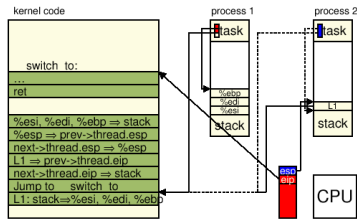
27

### Switch (6)



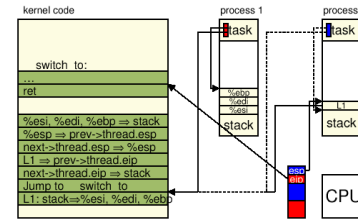
28

### Switch (7)



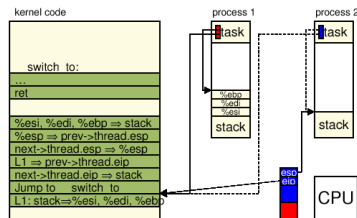
29

### Switch (8)



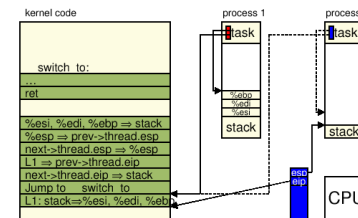
30

### Switch (9)



31

### Switch (10)



32



## Process switch - Code

```
#define switch_to(prev,next,last) do { \
    unsigned long esi,edi; \
    asm volatile("pushfl\n\t" \
        "pushl %%ebp\n\t" \
        "movl %%esp,%0\n\t" /* save ESP */ \
        "movl %5,%%esp\n\t" /* restore ESP */ \
        "movl $1f,%1\n\t" /* save EIP */ \
        "pushl %6\n\t" /* restore EIP */ \
        "jmp __switch_to\n\t" \
        "1:\n\t" \
        "popl %%ebp\n\t" \
        "popfl" \
        : "=m" (prev->thread.esp), "=m" (prev->thread.eip), \
        "=a" (last), "=S" (esi), "=D" (edi) \
        : "m" (next->thread.esp), "m" (next->thread.eip), \
        "2" (prev), "d" (next)); \
} while (0)
```

33

## Threads

- LinuxThreads is the standard POSIX thread library for Linux (1996)
- Based on principles of kernels of that time
  - Cheap kernel thread switches
  - Missing thread aware ABI
    - Thread local data with fixed relation to stack
  - Management thread necessary for creation etc.
  - No adequate kernel synchronization support
    - Signals abused
- Kernel is not aware of threads
  - Processes cooperate

## LinuxThreads problems

- Signal handling is not POSIX compliant
- Extra management thread
- `ps` shows all threads in a process, procs littered
- Core dumps do not contain the stack and machine registers for all threads
- `getpid()` returns different results for each thread
- Threads cannot wait for threads created by another thread
- Parent-child relationship instead of being peers
- **Threads do not share user and group ids**

35

## Kernel support added

- TLS (thread local storage) support in the kernel
- `clone` syscall extensions
  - Flag indicates that thread is created
- POSIX signal handling in the kernel
  - SIGSTOP forwarded to all threads of a process
- `exit` in two flavors for thread and process
- User level synchronization support
  - `futex` (**f**ast **u**ser **m**utex)

36

## Native POSIX thread library (NPTL)

- Better POSIX compliance
- Low startup/teardown costs
- Scalability
  - Enormous (100000) number of threads supported
- NUMA support
  - Node aware memory allocation
- Integration with C++