

## Monthly Planner

January

We	1	i. Interactive process: They instantly
Th	2	interacts with their users, therefore
Fr	3	they spend lot of time with keyboard
Sa	4	and mouse operations.
Su	5	• The average delay must be falls on
Mo	6	50 and 150 milliseconds.
Tu	7	• The variance of the average delay of
We	8	the system must be cratic.
Th	9	• Typically interactive programs are
Fr	10	command shell, text editors, Graphical
Sa	11	- CS applications.
Su	12	
Mo	13	ii. Batch process: They do not use user
Tu	14	interaction, but they often run in
We	15	the background.
Th	16	• Batch programs are programming lan
Fr	17	- guage compilers, database search en
Sa	18	- gines, scientific computations, data
Su	19	- backup, software updates etc.
Mo	20	
Tu	21	iii. real-time process: This process is
We	22	— stringent stable process.
Th	23	• real-time process are video and au
Fr	24	- dio and robot controllers, the pro
Sa	25	- cess is collection of data to the
Su	26	structures and sen
Mo	27	- sors etc.
Tu	28	
We	29	★. <u>System calls to the scheduler</u>
Th	30	
Fr	31	



4

Saturday

JANUARY 2020

Appointments - queue list. process is a fairness assign  
 - grant of cpu time all sched - RR used  
 - in a same priority of a process.

iii. Sched - normal % it is a conventional process.

- Scheduled algorithm is a differently depend on the behaviour of the conventional process.
- Every conventional process is having own in its static priority.
- The kernel of the scheduled algorithm of the number ranging from 100 to 139.

5

Sunday

Week - 1 C (005-361)

- The new process is inherits from its parent.
- process is changes from the system calls `nice()` and `setpriority()`.
- The real-time process is having a real-time priority() value changing from 1 to 99.
- The real-time priority changes from the set priority and `sched.setscheduler()` system calls.

★ Basic time quantum

- The static priority is essentially determines the basic time quantum process. quantum time duration is assigned to the process it has exhausted previous time of quantum.

"Risk more than others think is safe. Dream more than others think is practical." - Howard Schultz



- Static priority and base time quantum are related by the following formula:-

$$\text{base time quantum (in milliseconds)} = \begin{cases} (140 - \text{Static Priority}) \times 20 & \text{if Static priority} < 120 \\ (140 - \text{Static Priority}) \times 5 & \text{if Static priority} \geq 120 \end{cases}$$

### \* Dynamic priority

- Static priority, a conventional process also dynamic priority which value ranging from 100 to 139.
- It is related to the static priority by following empirical formula

$$\text{dynamic priority} = \max(100, \min(\text{static priority} - \text{bonus} + 5, 139))$$

### 3. Data structure used by the scheduler

- run queue is most important data structure in a Linux 2.6 scheduler.
- Each of the CPU system having own

"Design is not just what it looks like and feels like. Design is how it works." - Steve Jobs



## Appointments

I. Process in Linuxi. Process ~~monitor~~ Light weight processes and threads

- The term "process" often used with several different meanings. and process is a instance of program in execution.
- processes are like human beings? they generated, they have a more or less significant life, they optionally generate one or more child process, and eventually they die.
- From the kernel point of view, the purpose of the process is act as an entity to which system resource (cpu, time, memory etc.) are allocated.
- when a process is created, it is almost identical to its parent.
- It receives the copy of the parent address - its space and executes the same code as parent, beginning at the next instruction following the process creation system call.
- it supports multithreaded application use - programs having many relatively independent execution flows sharing large portion of the application data structure.
- In most of the multithreaded applic

"Do one thing every day that scares you." - Eleanor Roosevelt



- Applications are written using standard set of library functions called pthread (POSIX thread) libraries.
- Older version of the Linux kernel offered no support for a multithreaded application.
  - From kernel point of view, multithreaded application was a just normal process.
  - The multiple execution flows of a multithreaded application was created, handled and scheduled entirely user mode.
  - For instance, suppose a chess program uses two threads: one of them control a Graphical chess boards, waiting for the moves from the human player and showing the moves of the computers, while other thread ponders the next move of the game.
  - While first thread wait for the human move, the second thread should run continuously thus exploiting thinking of the human player.
  - A straight forward way to implement multithreaded applications is to associate a lightweight process with each thread.
  - Examples of the POSIX-compliant pthread libraries that use light-weight process are Linux threads, native POSIX Thread library, IBM next Generation POSIX threading package.



11 Saturday

JANUARY 2020

Week - 2 ○ (011-355)

Appointments

- Posix Compliant multithreaded applications are best handled by kernels that support "thread groups".
- In Linux a thread group is basically a set of lightweight processes that implement a multithreaded application and act as whole with regards to some system calls such as `getpid()`, `kill()` and `exit()`.

## 2. Process Switch

- i. process switching :- In order to control the execution of the processes, the kernel must be able to suspend processes running on the CPU and resume

- execution of the processes previously suspended. This is called a process switch.
- hence, task switching and context switching.

12 Sunday

Week - 2 ○ (012-354)

- ii. Hardware context :- While each process can have its own address space, all processes have to share the CPU registers. So before resuming the execution of the process, the kernel must ensure each register is loaded with the value.

- Set of data must be loaded into the registers before the process resumes execution on the CPU is called a

"A successful man is one who can lay a firm foundation with the bricks others have thrown at him." - David Brinkley



hardware context.

- we will assume the prev local variable refers to the process descriptor of the process being switched out and next refers to the one being switched in to replace it.

• we can define process switch as the activity consisting of saving the hardware context prev and replace with the hardware context of next.

- older version of the linux took advantage of the hardware support offered by the x86 architecture and performed a process switch through a far jmp instructions to the selector of the Task state segment descriptor to the next process.
- Linux 2.6 Software to perform a process switch for the following process:

- For far jmp instructions modify both the CS and EIP registers, while simple jmp instructions modify only EIP.
- Step-by-step switching performed through a sequence of mov instructions.
  - allows better control over the validity of data being loaded.
- The amount of time required old approach and new approach become same.
- process switching occurs only in the kernel mode.

ii. Task State Segment:

Challenges are what make life interesting and overcoming them is what makes life meaningful - Joshua J. Marine



Appointments V. Switch - to() function.

- Switching to new task: we can switch from one task to another, we call Switch - to() function.
- Saving Registers: Switch - to starts by saving the registers (like `eax`, `edx` etc) of the current task (Task A). it's like taking the snapshot of everything the task is currently holding.
- Updating the task: now update the task Switch - to (Task B). it load the Task B registers values. it's like Task B starts with all right information.
- Restoring Registers: when we switch back to A later, Switch - to() restore the stored register values, Task A continue where it left off.
- Control Transfer: it also take care of the transferring control to the new task (Task B). Then it means `cpu` start executes a B instructions.
- Completion: make everything is ready for the next task switch.
- Ending the function: finally it ends its job, and OS continues running on chosen task.

Vi. Saving the floating point registers



17 Friday

Appointments

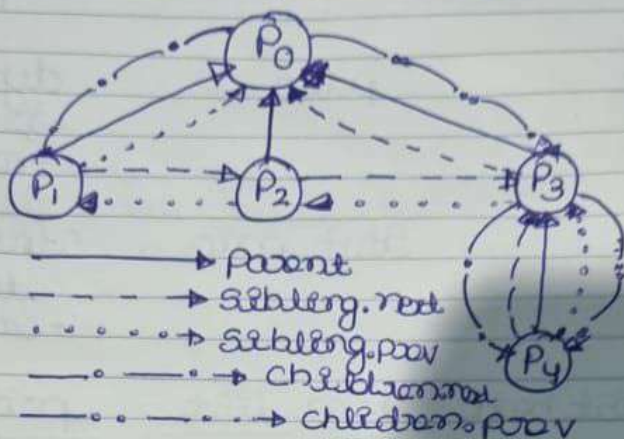
- Starting with the Intel 80486 processors, a part of the CPU is dedicated to performing complex math operations called Floating-point unit (FPU).
  - Intel introduces a MMX instruction set to speed up multimedia applications.
  - The advantage of this choice is the operating system does not need to handle MMX instructions differently.
  - MMX instructions speed up multimedia applications, because they introduce a single-instruction multiple data for 128-bit registers called the XMM registers.
  - The 80x86 microprocessors do not automatically save the FPU, MMX, and XMM registers in TSS.
  - To manage this selective saving & restoring of floating-point, MMX and XMM registers there are data structures. They stored part of the TSS called tss.1387.
- ```

if (prev->flags & PF_USED_FPU) {
    asm ("fsave %0" : "=m" (prev->tss.1387));
    asm ("fwait");
    prev->flags &= ~PF_USED_FPU;
    stts();
}

```
- The stts() macro sets the TS flag of CR0.



Appointments

i. Relationships among Processes• parenthood relationships

i. real-parent : points to the process descriptor of a process that created  $P$ .

ii. parent : points to the current parent of  $P$ . its value coincides with real-parent.

iii. children : head of the list containing all children created by  $P$ .

iv. Sibling : The pointers to the next and previous elements of the sibling list of process that have same parent as  $P$ .

• Non-parenthood relationships

i. group-leader : process descriptor pointer of group leader  $P$



Appointments

```

6 Declare - wait - queue - head (name)
  int - waitqueue - head ( )
  int - waitqueue - entry (q, p) {
8     q -> flags = 0;
9     q -> task = p;
10    q -> func = default - wake - func;
11  }
12 Define - wait ( )
13 add - wait - queue ( )
  add - wait - queue - exclusive ( )
  remove - wait - queue ( )
  wait queue - active ( )

```

### iii. Process Resource Limits

- The resource limits are stored in the current -> signal -> rlim field

```

17 struct rlimit {
  unsigned long rlim_cur;
  unsigned long rlim_max;
18 };

```

Rlimit - AS -> maximum size of address space  
 " - cpu -> maximum cpu time of process  
 " - data -> maximum heap size in bytes  
 20 " - core -> maximum core dump file size  
 " - fsize -> maximum file size allowed  
 " - locks -> maximum number of file locks  
 " - stack -> maximum stack size in bytes  
 " - rss -> maximum number of page frames owned by process



## II Process scheduling

### 1. Scheduling policy :

- Set of rules used to determine when and how to run new process is called a scheduler policy.
- Linux is based on time-sharing technique and several process to run on "time-multiplexing" because CPU time is divided in to two slices. One of the runnable process of course of single process is on one process at given instant of time.
- Scheduling policy is based on ranking process according to their priority.
- Linux process priority is dynamic. Scheduling policy keep track of what doing adjust to their according priority periodically.

### A. Class of process

- There are 3 types of class of process
  - i. Interactive process
  - ii. Batch process
  - iii. Real-time process



Wednesday

JANUARY 2020

1

APPENDIX

TASKS

i. nice()

Description  
change of the static  
priority to conventio-  
-nal process.

ii. Setpriority()

Set the static priori-  
-ty of group of conve-  
-ntional process.

iii. getpriority()

get the static priori-  
-ty of maximum of  
group of conventional  
process.

iv. sched\_setsche-  
-duler()

Set the scheduler poli-  
-cy.

v. sched\_getsche-  
-duler()

get the scheduler policy  
of real-time priority  
of a process.

vi. sched\_setparam()

Set the real-time pri-  
-ority of a process.

vii. sched\_getparam()

get the real-time pri-  
-ority of a process.

viii. sched\_yield()

Relinquish of a con-  
-defined process.

ix. sched\_getpriori-  
-ty\_min()

get the minimum  
real-time value of  
a process.

Your future is hidden in what you do daily - Mike Murdock



x. Sched - get - priority  
max()

get the maximum  
real-time value  
of a process.

xi. Sched - sched - get -  
priority()

get the real-time  
round-robin  
policy.

xii. Sched - setaffinity()

Set the cpu affi-  
-nity of mask of  
process.

xiii. Sched - getaffinity()

get the cpu affi-  
-nity of mask of  
process.

## ★. Process preemption

- Linux process is preemptable.
- A process is preempted when its quantum time expires.
- preemption is a fundamental multi-tasking system that ensures fairness, responsiveness and efficient resource utilization.

## ★. How long quantum last?

- quantum duration is a critical for system.
- performance: it should be neither too short or long.

Thinking is the hardest work there is, which is probably the reason why so few engage in it. - Henry Ford



3

Friday

JANUARY 2020

• The average quantum duration is too short.

## 2. The scheduling algorithm

- Scheduling algorithm used in the early version of Linux was quite simple and straightforward. Each process is switch to kernel scheduled list of compatible priorities and select best process to run.
- Scheduling algorithm used Linux 2.6 was sophisticated.
- Every scheduling process according to their following of scheduling class

i. Sched-FIFO : First - In first out of real-time process when scheduled it assigns to CPU to process and it leaves process descriptor and running its time of current position of the runqueue list.

• At the highest priority of the real-time process. and the process uses CPU as long as wishes. even if other real-time process is having same priority.

ii. Sched-RR : Round-robin policy process when scheduled assigns CPU to the process it puts in process descriptor at the end of the process in the run

The starting point of all achievement is desire. - Napoleon Hill-



7

Tuesday

JANUARY 2020

Appointments

-runqueue.

- runqueue list has been stored in the run-queue cpe variable.

### Types of runqueue data structure

| Types              | name                     | Description                                                       |
|--------------------|--------------------------|-------------------------------------------------------------------|
| i. Spinlock        | lock_t                   | protect the list of process.                                      |
| ii. Unsigned-long  | nr-runnable              | number of runnable process in the list                            |
| iii. Unsigned-long | cpu_load                 | cpu load factor of a runnable process.                            |
| iv. Unsigned-long  | nr-switches              | number of switches for the cpu - 0.                               |
| v. Unsigned-long   | expected-uninterruptable | number of switches for the cpu is to sleep in the uninterruptable |
| vi. Unsigned-long  | executable-time stamp    | inter relation - ships of the runnable process                    |

"Always deliver more than expected." - Larry Page

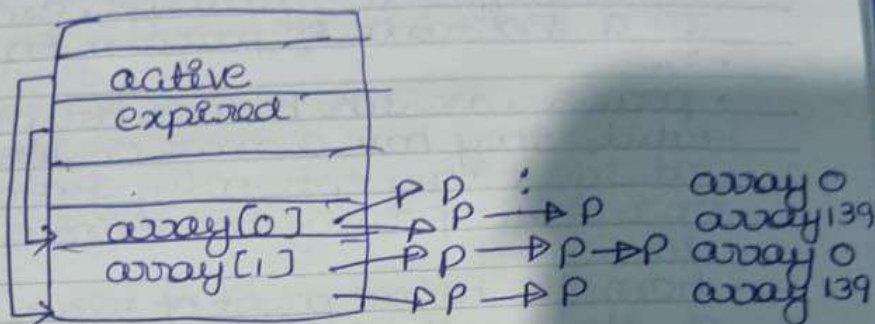


vi. unsigned long

timestamp -  
last tick

timestamp  
of the value  
time.

- most important of the runqueue data structure is a runnable process.



- The runqueue list is having a two arrays. of a each runnable process is having a dynamic priority.
- The active field of prio\_struct array data structure is corresponds to each runnable process in the active field.
- The expired field of prio\_struct array data structure is corresponds to each runnable process in the expired field.



## Appointments

- The Intel 80x86 architecture includes a specific segment type called task state - all segment to store a hardware context - S.
- There are two main reasons:
  - i. when 80x86 cpu switches from user mode to kernel mode it fetches the address of the kernel mode stack from the TSS.
  - ii. when a user mode process attempts to access I/O port by means of in or out instructions.
- Each TSS has its own 8-byte Task State Segment descriptor (TSSD). This descriptor includes a 32-bit base field pointing to the TSS.
- Type field is set to 11 if the TSSD refers to the TSS of the process currently running on the cpu is set to 9.

iv. Switch to macro

- The Switch to macro typically does on an 80x86 microprocessor by using standard assembly language.
- Save the values of the prev and next in the eax and edx registers respectively:
 

```
movl prev, %eax
movl next, %edx
```
- Saving Flags and EBP: pushfl (to save flags) and pushl %ebp (to save EBP)

The Pessimist Sees Difficulty In Every Opportunity. The Optimist Sees Opportunity In Every Difficulty. - Winston Churchill

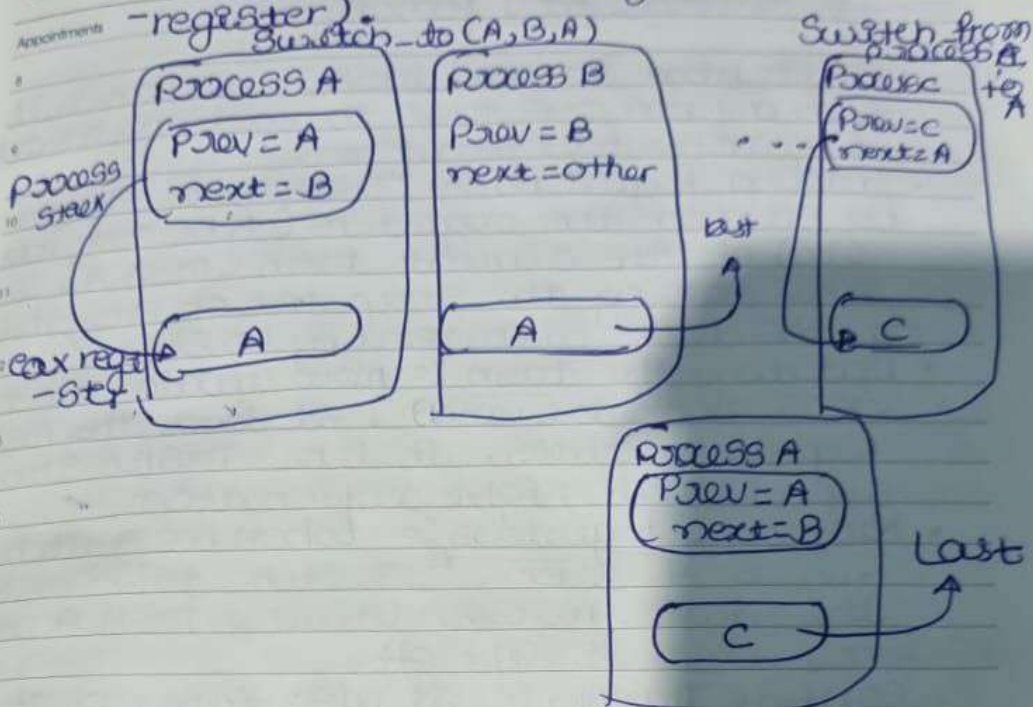


15

• The Switch-macro performs a process switch. The Switch-macro consists of the following parameters, prev; next & last

Wednesday - rec parameters, prev; next & last

macro is invoked by the scheduler



- updating the stack pointer to prev's  
- use process: `movl %esp, 484(%eax)`
- Loading stack pointer to the next process: `movl 480(%edx), %esp`
- pushing the next process instruction pointer: `pushl 480(%edx)`. This prepares the process to start executing code of the next process.
- jumping to "-Switch-to" function:  
Imp-Switch-to it essentially hands to control over the next process.
- Return the original process: `popl %ebp` (to restore EBP) and `popl C` (to restore flags)
- `movl %eax, last`

"Opportunities don't happen. You create them." - Chris Grosser



18 Saturday

JANUARY 2020

Week - 3 (018-340)

Appointments

i. `movl %cro, %eax` ; Load the value of `cro` to `EAX`

ii. `orb $8, %al` ; Set the <sup>bit</sup> (bit 3) in `AL`

iii. `movl %eax, %cro` ; Store the modified value back to `cro`

\*. `void math_state_restore(void)`

{  
    `asm ("alt5")`

    if (current -> used\_math)

`asm ("frstor %0" : : "m" (current -> tss.1387)) ;`

19 Sunday else {

Week - 3 (019-347)

`asm ("fncinit") ;`

    current -> used\_math = 1 ;

current -> flags |= PF\_USGDFPU ;

### 3. Process Descriptor

- Every process descriptor includes several fields related to scheduling.

\* Types

• Unsigned long

name

thread-info  
-> flags

Description

Stores the TIF\_NEED\_RESCHEDULED flag

Great minds discuss ideas, mediocre minds discuss events, small minds discuss people



20

Monday

JANUARY 2020

Munir Lutfi King in Day (USA)

Appointments

- unsigned\_int thread\_info → cpu logical num  
- ber of cpu  
owning the  
runqueue.
- int prio dynamic pri  
- ority of a pro  
- cess.
- int stat\_prio static prio  
- ty of a pro  
- cess.
- struct list\_head run\_list pointers to the  
next previous  
elements in  
the runqueue  
list.
- unsigned\_long sleep\_avg average sleep ti  
- me of the pro  
- cess.

- when a new process is created, Sched\_fork() invoked by copy\_process(), set the timeslice - ce field of both current process in the following way:-  

$$P \rightarrow \text{time\_slice} = (\text{current} \rightarrow \text{time\_slice} + 1) > 1;$$

$$\text{current} \rightarrow \text{time\_slice} > 1;$$

You can tell whether a person is clever by his answers. You can tell whether a man is wise by his questions - Naguib Mahfouz

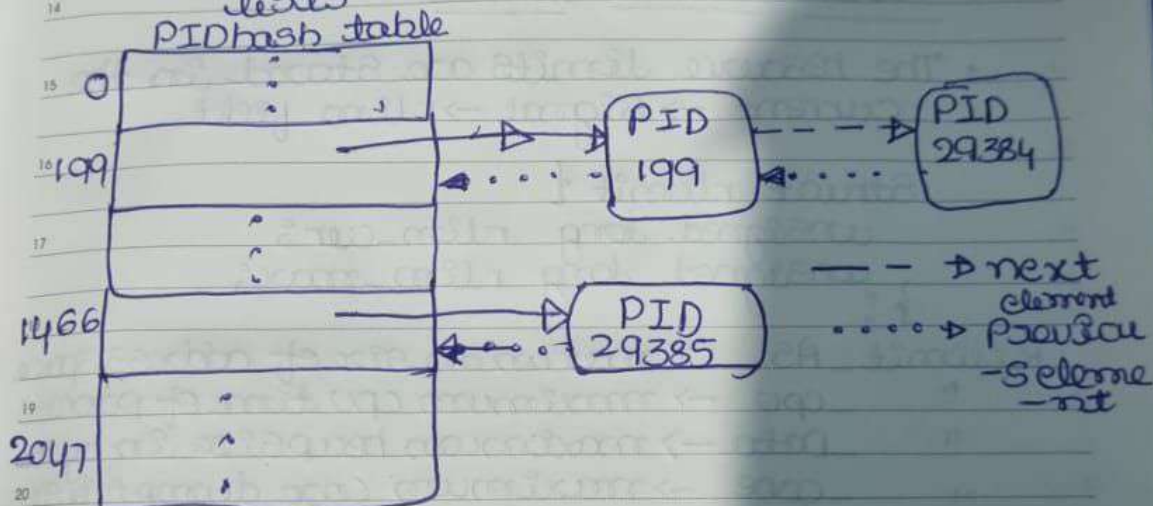


Appointments

- signal  $\rightarrow$  prog P : PID of the group leader P.
- process descriptor is having different types of PID and each type PID requires its own hash table.

| Hash table type | Field name | Description                        |
|-----------------|------------|------------------------------------|
| PIDTYPE - PID   | Pid        | pid of the process                 |
| PIDTYPE - PGID  | Pgrp       | pid of the group leader            |
| PIDTYPE - TGID  | tgid       | pid of the thread group leader     |
| PIDTYPE - SID   | session    | pid of the session leader process. |

- A simple PID hash table and chained lists



## ii. How processes are organized

- Handling wait queues



24

Friday

JANUARY 2020

#### 4. Destroying process

- process terminates the execution of the code they were supposed to run.
- process terminate is to invoke the `exit()` function.

#### i. Process Termination

- In Linux 2.6 there are two system calls that terminates the usermode application.

i. `exit_group()` System call: The main kernel function that implements the system call is called `do_group_exit()`.

ii. `exit()` System call: The main kernel function that implements this system call is called `do_exit()`.

iii. The `do_group_exit()` function: Kills all the processes belonging to the thread group of current after receiving parameter.

• Steps:

i. checking `SIGNAL_GROUP_EXIT` flag.  
If not zero, `do_exit()`

ii. Else setting the `SIGNAL_GROUP_EXIT` flag.

iii. Invoking `zap_other_threads()` function

iv. Invoking `do_exit()` function

"The secret to success is look for the people who want to change the world." - Marc Benioff



25 Saturday

Chinese Lunar New Year (China, Indonesia, Malaysia, Hong Kong and Singapore)

Appointments

- ~~Step~~ The `do_exit()` function: Removes most references to the terminating process from the kernel data structures. Receives process termination code as parameter.

- Steps:

- Setting `PF_EXITING` flag.
  - calling `del_timer_sync()` function
  - Detaching data structures
  - Decreasing usage counters
  - Setting `exit_code` field
- Invoking `exit_notify()` function
    - updating parent-child relationships
    - checking `exit_signal` field and if process is lost in thread group.

26 Sunday

Week - 4 • (329-340)

Republic Day (India), National Day (Australia)

- checking `exit_signal` field and if group has other processes.
  - Setting `exit_state` field to `EXIT_DEAD` and invoking `release_task()` if process is not being traced.
  - Setting `exit_state` field to `EXIT_ZOMBIE` if `exit_signal` not equal to 1.
  - Setting `PF_DEAD` flag.
- Invoking `Schedule()` function.

## ii. Process Removal

### Steps in `release_task()` function

"It is not the strongest of the species that survive, nor the most intelligent, but the one most responsive to change." - Charles Darwin



27

Monday

JANUARY 2020

- Decrease a number of processes belonging to user owner of terminated process.
- Removing process from `ptrace-children` list and assigning parent if being traced.
- Invoking `--exit-signal()` and `exit-sig-misc()`
- Invoking `--exit-sig-handl()`
- Invoking `--unhash-process()`
  - i. `nr-threads --`
  - ii. Invoking `detach-pid()`
  - iii. Invoking `detach-pid()` again if process is thread group leader
  - iv. Using `Remove-LINK` macro.
- Sending signal to notify process death to parent.
- Invoking `sched-exit()` function
- Invoking `put-task-struct()`: if count `-er=0`, dropping remaining reference
  - Decrease usage counter and release `-ng` data structure if usage counter=0
  - Releasing process descriptor and memory area.

## 5. creating Process

- Unix operating systems do, process creation to satisfy user requests. For example, the shell creates a new process that executes another copy of the shell whenever the user enters a command

"Fall seven times and stand up eight" - Anonymous



i. Modern Unix Kernel solve this problem  
 - by introducing three different  
 - kernel mechanisms

• The copy on write technique allows both the parent and the child to read the same physical pages

• Lightweight process: allows both the parent and the child to share many processes  
 - Kernel data structures, such as the paging tables, the open file tables and the signal dispositions.

• vfork() system call creates process that shares memory address space to parent

• To prevent the parent from the over-  
 - ting data needed by the child, the  
 parent execution is blocked until child exists

ii. The clone(), fork() and vfork() system calls

• fn : specific function to be executed by the new process

• arg : points to data passed to the fn() function.

• child\_stack : specified the usermode stack pointer to be assigned



to the esp register of the child process.

- `tls` :• meaningful only if the `clone-  
set_tls_flag` is set.
- `child` :• meaningful only if the `clone-  
child_settid` flag is set.
- `clone()` :• is actually wrapper function  
defined in C library.

### iii. The `do_fork()` function

- The `do_fork()` function, which handles the `clone()`, `fork()`, and `vfork()` system calls, acts on the following parameters:
  - i. `clone_flags` :• Same as the flags parameter of the `clone`
  - ii. `stack_start` :• Same as the child stack parameter of the `clone`.
  - iii. `regs` :• pointers value of the general purpose registers saved into the kernel mode stack when switching from the usermode to kernel mode.

iv. `stack_size` :• Unused (always set to 0).



30

Thursday

JANUARY 2020

Appointments

- parent - tidptr, child - tidptr  
same as the corresponding ptid  
and child parameters of the  
abnec).

"Choose the vision, not the money, the money will end up following you." - Tony Hsieh