# Process Management

Jakob Frank
Dillon Tice

# Linux

- Processes are managed in memory as pointers to structs known as task_structs

- Each struct contains a large amount of information for each process, such as process id, state, name of executable, and links to other processes such as children, parents, and siblings

- fork() and exec() are primary used for the creation of new processes

# Structure

- **task_struct**
  - Differs by architecture
- **Init**
  - Statically allocated task_struct
  - /init/init_task.c

```
56  struct task_struct init_task
57  #ifdef CONFIG_ARCH_TASK_STRUCT_ON_STACK
58          __init_task_data
59  #endif
60  = {
61  #ifdef CONFIG_THREAD_INFO_IN_TASK
62          .thread_info    = INIT_THREAD_INFO(init_task),
63          .stack_refcount = REFCOUNT_INIT(1),
64  #endif
65          .state          = 0,
66          .stack          = init_stack,
67          .usage          = REFCOUNT_INIT(2),
68          .flags          = PF_KTHREAD,
69          .prio           = MAX_PRIO - 20,
70          .static_prio    = MAX_PRIO - 20,
71          .normal_prio    = MAX_PRIO - 20,
72          .policy         = SCHED_NORMAL,
73          .cpus_ptr       = &init_task.cpus_mask,
74          .cpus_mask      = CPU_MASK_ALL,
75          .nr_cpus_allowed= NR_CPUS,
76          .mm             = NULL,
77          .active_mm      = &init_mm,
78          .restart_block  = {
79                  .fn = do_no_restart_syscall,
80          },
81          .se             = {
82                  .group_node     = LIST_HEAD_INIT(init_task.se.group_node),
83          },
84          .rt             = {
85                  .run_list       = LIST_HEAD_INIT(init_task.rt.run_list),
86                  .time_slice     = RR_TIMESLICE,
87          },
88          .tasks          = LIST_HEAD_INIT(init_task.tasks),
```

# task_struct fields

- state
  - A set of bits indicating process state
    - Running, stopped, interruptible, uninterrup
- flags
  - Is a process being created? Exiting? Allocating
- comm
  - Name of the executable without the path
- static_prio
  - Priority
    - Lower priority > higher priority
    - *Actual* priority determined dynamically

# task_struct fields

- mm and active_mm
  - Process memory descriptors
  - Active is the previous process descriptors
    - Context switching
- tasks
  - Used for linked list representation
    - Init is the head
  - Prev and next pointer
- thread_struct
  - Context switch storage
    - Registers, program counter, etc.

```
#ifdef CONFIG_SMP
    struct plist_node          pushable_tasks;
    struct rb_node             pushable_dl_tasks;
#endif

    struct mm_struct           *mm;
    struct mm_struct           *active_mm;

    /* Per-thread vma caching: */
    struct vmacache            vmacache;

#ifdef SPLIT_RSS_COUNTING
    struct task_rss_stat       rss_stat;
#endif
```

```
1276    /*
1277     * New fields for task_struct should be added above here, so that
1278     * they are included in the randomized portion of task_struct.
1279     */
1280    randomized_struct_fields_end
1281
1282    /* CPU-specific state of this task: */
1283    struct thread_struct          thread;
1284
1285    /*
1286     * WARNING: on x86, 'thread_struct' contains a variable-sized
1287     * structure.  It *MUST* be at the end of 'task_struct'.
1288     *
1289     * Do not put anything below here!
1290     */
```

# Overview of Linux Process Management

- With regards to creating and managing process, linux contains the following commands:

    - Commands for creating processes via fork() and exec()
        - Can also use System() but it is slow and a security risk

    - Commands for killing processes such as kill, exit, pkill, and killall

    - Commands for viewing process information such as pgrep(), ps(), and top()

# Process States

- Possible process states:
  - TASK_RUNNING
    - Running or in run queue
  - TASK_INTERRUPTIBLE
    - Sleeping
  - TASK_UNINTERRUPTIBLE
    - Sleeping but unable to be awoken
  - TASK_STOPPED
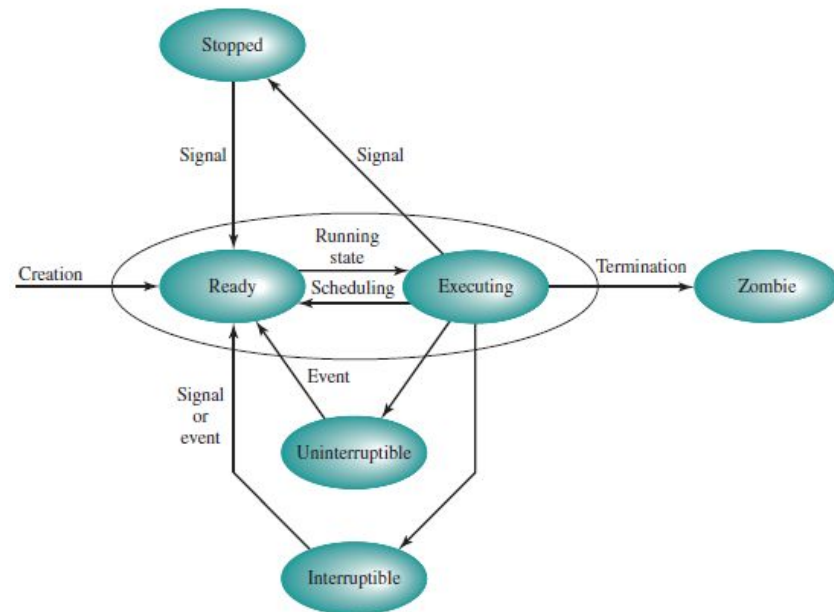    - Task becomes a zombie (stopped but still contains an entry in the process table)



Figure 4.15   Linux Process/Thread Model

# Creation

- fork() -> clone() -> do_fork()
- do_fork() - clone_flags
  - copy_process(..., args)
- copy_process()
  - Security_task_create
    - Further SELinux functions
    - Return value == 0 ?
  - dup_task_struct()
  - copy_creds()
  - CPU time reset
  - CPU time sharing setup

```
1   long _do_fork(struct kernel_clone_args *args)
2   {
3       u64 clone_flags = args->flags;
4       struct completion vfork;
5       struct pid *pid;
6       struct task_struct *p;
7       int trace = 0;
8       long nr;
9
10      /*
11       * Determine whether and which event to report to ptracer.  When
12       * called from kernel_thread or CLONE_UNTRACED is explicitly
13       * requested, no event is reported; otherwise, report if the event
14       * for the type of forking is enabled.
15       */
16      if (!(clone_flags & CLONE_UNTRACED)) {
17          if (clone_flags & CLONE_VFORK)
18              trace = PTRACE_EVENT_VFORK;
19          else if (args->exit_signal != SIGCHLD)
20              trace = PTRACE_EVENT_CLONE;
21          else
22              trace = PTRACE_EVENT_FORK;
23
24          if (likely(!ptrace_event_enabled(current, trace)))
25              trace = 0;
26      }
27
28      p = copy_process(NULL, trace, NUMA_NO_NODE, args);
29      add_latent_entropy();
30
31      if (IS_ERR(p))
32          return PTR_ERR(p);
33      |
```

 _do_fork

# Destruction

POSIX defines two ways regarding how a process can be terminated

- A process can terminate itself, either by calling exit(),
  _exit(), returning from main(), or the last thread of the
  process terminates
- A process can be killed by a signal, possibly sent by the
  kernel, another process, or the process itself

# Alternative Ideas/Strategies

- Different built in Linux scheduling policies
  - SCHED_FIFO
  - SCHED_RR
  - SCHED_OTHER
- Different task manager implementations
  - htop
  - conky
  - pstree
  - GNOME System Monitor