Gabriele Proietti Mattia

# Advanced Operating Systems and Virtualization

[10] Process Management



Department of Computer, Control and Management Engineering "A. Ruberti", Sapienza University of Rome

gpm.name · proiettimattia@diag.uniroma1.it

A.Y. 2020/2021 · V2

# Outline

- 1. Process Control Block
  - 1. Accessing the PCB
- 2. The fork()/exec() model
  - 1. Kernel Threads
- 3. Out Of Memory (OOM) Killer
- 4. Process Starting
  - 1. The ELF Format
  - 2. Dynamic Linking
  - 3. Initial Steps of Programs' Life



10. Process Management

# **Process Control Block**



Advanced Operating Systems and Virtualization

#### **Processes**

The concept of a process is fundamental to any multiprogramming OS. The term process is often used with several different meanings, for us it means an **instance of a program in execution** (or even a set of data structures which describes how far the execution has progressed).

When a process is created is almost **identical** to its parent:

- it receives a logical copy of the parent's address space
- it executes the same code of the parent, at next instruction after the fork()

But the child process has separate copies of the data (stack and heap) so that changes in the child are invisible to the parent and vice versa.

#### **Processes**

While earlier versions of Unix supported this model, modern ones do not. They instead supported multithreaded applications, in which a process is composed of several *user threads* (or simply threads), each of which represents a an execution flow of the process (pthread library).

Older versions of Linux did not support multithreaded applications, so from the kernel point of view a multithreaded application was just a normal process. So threads were created, handled and scheduled in User Mode. Therefore if a thread was blocked for a system call, every other thread would be blocked.

Nowadays Linux uses **lightweight processes** that are independent from each other but at the same time they can share resources (e.g. memory). They are mapped, in the end, to threads. A process in modern versions of Linux is just a group of lightweight processes, also called a **thread group** (created with pthread library).

### **The Process Control Block**

To manage processes, the kernel must have a clear picture of what each process is doing, for instance, the priority, the state, the address space and so on. This is the role of the **process descriptor** (also called Process Control Block - PCB).



### task\_struct

The **struct** task\_struct object represents the Process Control Block within the Linux Kernel. This is declared in **include/linux/sched.h** and it is one of the largest structures in the kernel (almost 600 LOCs). Relevant members are:

- volatile long state
- **struct** mm\_struct \*mm
- **struct** mm\_struct \*active\_mm
- pid\_t pid
- pid\_t tgid
- **struct** fs\_struct \*fs
- **struct** files\_struct \*files
- **struct** signal\_struct \*sig
- struct thread\_struct thread /\* CPU-specific state: TSS, FPU, CR2, perf events, ...
   \*/
- int prio; /\* to implement nice() \*/
- **unsigned long** policy /\* for scheduling \*/
- **int** nr\_cpus\_allowed
- cpumask\_t cpus\_allowed

#### task\_struct



Figure 3-1. The Linux process descriptor

Bovet, Daniel P., and Marco Cesati. Understanding the Linux Kernel: from I/O ports to process management. "O'Reilly Media, Inc.", 2005.

#### **Process State**

The state field of a process descriptor describe what is currently happening to the process. Possible states are always exclusive:

- TASK\_RUNNING the process is either executing on CPU or waiting to be executed
- TASK\_INTERRUPTIBLE the process is sleeping until some condition becomes true
- TASK\_UNINTERRUPTIBLE like TASK\_INTERRUPTIBLE but except that raising a signal to the process will leave the state unchanged
- TASK\_STOPPED process has stopped (after signal SIGSTOP, SIGTSTP)
- TASK\_TRACED process execution has stopped by a debugger

Then to the exit\_state field we can have

- EXIT\_ZOMBIE process terminated but the parent did not issued wait to retrieve the data, so the kernel cannot discard it
- EXIT\_DEAD the parent issued wait

#### The mm member

The field mm points to a mm\_struct defined in include/linux/mm\_types.h. The mm\_struct is used to manage the memory map of the process:

- virtual address of the page table (pgd member)
- a pointer to a list of vm\_area\_struct records (mmap field)

Each record tracks a user-level virtual memory area which is valid for the process. active\_mm is used to "steal" a mm when running in an anonymous process, and mm is set to NULL. Non-anonymous processes have active\_mm == mm.

386	struct	mm_struct {
200		struct {
300		struct vm_area_struct "mmap; /* List of VMAS */
389		struct rb_root mm_rb;
390		ub4 vmacacne_seqnum; /* per-thread vmacacne */
391	#itdet	CONFIG_MMU
392		unsigned long (* <b>get_unmapped_area</b> ) (struct file *filp,
393		unsigned long addr, unsigned long len,
394		unsigned long <b>pgoff</b> , unsigned long flags);
395	#endif	
396		<pre>unsigned long mmap_base;</pre>
397		unsigned long mmap_legacy_base; /* base of mmap area in bottom-up allocations */
398	#ifdef	CONFIG HAVE ARCH COMPAT MMAP BASES
399		/* Base adresses for compatible mmap() */
400		unsigned long mmap compat base;
401		unsigned long mmap compat legacy base:
402	#endif	
403		unsigned long task size: /* size of task vm space */
404		unsigned long highest vm end: /* highest vma end address */
405		ngd + * ndg
100		

https://eli

mm.html

#### vm\_area\_struct

The struct vm\_area\_struct describes a Virtual Memory Area (VMA). It contains:

- **struct** mm\_struct \*vm\_mm: the address space the structure belongs to
- **unsigned long** vm\_start: the start address in vm\_mm
- **unsigned long** vm\_end: the end address
- pgprot\_t vm\_page\_prot: access permissions of this VMA
- **const struct** vm\_operations\_struct \*vm\_ops: operations to deal with this structure
- **struct** mempolicy \*vm\_policy: the NUMA policy for this range of addresses
- **struct** file \*vm\_file: pointer to a memory-mapped file
- struct vm\_area\_struct \*vm\_next, \*vm\_prev: linked list of VM areas per task, sorted by address

#### vm\_operations\_struct

560	/*
561	* These are the virtual MM functions - opening of an area, closing and
562	* unmapping it (needed to keep files on disk up-to-date etc), pointer
563	* to the functions called when a no-page or a wp-page exception occurs.
564	*/
565	<pre>struct vm operations struct {</pre>
566	<pre>void (*open)(struct vm_area_struct * area);</pre>
567	<pre>void (*close)(struct vm_area_struct * area);</pre>
568	/* Called any time before splitting to check if it's allowed */
569	int (* <b>may_split</b> )(struct <b>vm_area_struct</b> * <b>area</b> , unsigned long addr);
570	int (* <b>mremap</b> )(struct <b>vm_area_struct</b> * <b>area</b> , unsigned long flags);
571	/*
572	* Called by mprotect() to make driver-specific permission
573	* checks before mprotect() is finalised. The VMA must not
574	* be modified. Returns 0 if eprotect() can proceed.
575	*/
576	<pre>int (*mprotect)(struct vm_area_struct *vma, unsigned long start,</pre>
577	unsigned long end, unsigned long newflags);
578	vm_fault_t (*fault)(struct vm_fault *vmf);
579	vm_fault_t (*huge_fault)(struct vm_fault *vmf,
580	enum page_entry_size pe_size);
581	void (*map_pages)(struct vm_tault *vmt,
582	<pre>pgott_t start_pgott, pgott_t end_pgott);</pre>
583	unsigned long (*pagesize)(struct vm_area_struct * area);
584	(* petition that a provincely productly area is shout to become
202	/~ notification that a previously read-only page is about to become
500	" willable, II an error is returned it will cause a SIGBUS */
500	vii_rautt_t ('page_iikwirte)(struct Vii_rautt "Viii');
000	

#### int (\*access)(struct vm area struct \*vma, unsigned long addr, 596 void \*buf, int len, int write); 597 598 /\* Called by the /proc/PID/maps code to ask the vma whether it \* has a special name. Returning non-NULL will also cause this \* vma to be dumped unconditionally. \*/ const char \*(\*name)(struct vm area struct \*vma); #ifdef CONFIG NUMA 604 /\* \* set policy() op must add a reference to any non-NULL @new mempolicy 606 \* to hold the policy upon return. Caller should pass NULL @new to \* remove a policy and fall back to surrounding context--i.e. do not \* install a MPOL DEFAULT policy, nor the task or system default \* mempolicy. 610 \*/ int (\*set policy)(struct vm area struct \*vma, struct mempolicy \*new); 611 612 613 /\* 614 \* get policy() op must add reference [mpol get()] to any policy at 615 \* (vma,addr) marked as MPOL SHARED. The shared policy infrastructure 616 \* in mm/mempolicy.c will do this automatically. 617 \* get policy() must NOT add a ref if the policy at (vma,addr) is not 618 \* marked as MPOL SHARED. vma policies are protected by the mmap lock. 619 \* If no [shared/vma] mempolicy exists at the addr, get policy() op \* must return NULL--i.e., do not "fallback" to task or system default 620 621 \* policy. 622 \*/ 623

602

608

624

struct mempolicy \*(\*get policy)(struct vm area struct \*vma, unsigned long addr):

### **Userspace Memory Management**



### **Userspace Memory Management**



https://manybutfinite.com/post/how-the-kernel-manages-your-memory/

### **Userspace Memory Management**

----→ vm\_end: first address outside virtual memory area
→ vm start: first address within virtual memory area



https://manybutfinite.com/post/how-the-kernel-manages-your-memory/

#### From Kernel 2.6

Processes are very dynamic entities whose lifetime ranges from a few milliseconds to months. For this reason the kernel must be able to handle many processes at the same time and process descriptors are allocated in **dynamic memory**, rather than in the memory permanently assigned to the kernel. Therefore PCBs can be dynamically allocated upon request.

For each process, Linux packs **two different data structures** in a single per-process memory area: thread\_info and the Kernel Mode process stack. The length of this area is usually **2 pages** (8,192Kbytes).





Bovet, Daniel P., and Marco Cesati. Understanding the Linux Kernel: from I/O ports to process management. " O'Reilly Media, Inc.", 2005.

#### **union** thread\_union

This union is used to easily allocate thread\_info at the base of the stack, independently of its size. It works as long as its size is smaller than the stack's, of course, this is mandatory.

1732	union thread_union {
1733	<pre>#ifndef CONFIG_ARCH_TASK_STRUCT_ON_STACK</pre>
1734	<pre>struct task_struct task;</pre>
1735	#endif
1736	<pre>#ifndef CONFIG_THREAD_INFO_IN_TASK</pre>
1737	struct thread info thread info;
1738	#endif
1739	<pre>unsigned long stack[THREAD SIZE/sizeof(long)];</pre>
1740	};

https://elixir.bootlin.com/linux/v5.11/source/include/linux/sched.h#L1732

V5.11

struct thread\_info

This is the organization of thread\_info up to version 4.3. Later on, thread\_info has been progressively deprived of most members on x86. Security implications of this struct on the stack have been severe.

55	<pre>struct thread_info {</pre>	
56	struct task_struct	<pre>*task;</pre>
57	u32	<pre>flags;  /* low level flags */</pre>
58	u32	<pre>status;  /* thread synchronous flags */</pre>
59	u32	cpu; /* current CPU */
60	int	<pre>saved_preempt_count;</pre>
61	mm_segment_t	addr_limit;
62	void <u>user</u>	*sysenter_return;
63	unsigned int	<pre>sig_on_uaccess_error:1;</pre>
64	unsigned int	<pre>uaccess_err:1; /* uaccess failed */</pre>
65	};	

https://elixir.bootlin.com/linux/v4.3/source/arch/x86/include/asm/thread\_info.h

struct thread\_info



https://elixir.bootlin.com/linux/v5.11/source/arch/x86/include/asm/thread\_info.h

Where's the task\_struct pointer?

# Virtually Mapped Kernel Stack

Kernel-level stacks have always been the weak point in the system design. This is quite small: you must be careful to avoid overflows. **Stack overflows** (and also recursion overwrite) have been successfully used as attack vectors.

When an overflow occurs, the Kernel **is not easily able to detect it**. Even less able to counteract on it! Stacks are in the ZONE\_NORMAL memory and are contiguous but access is done through the MMU via virtual addresses



# Virtually Mapped Kernel Stack

There is no need to have a physically contiguous stack, so Andy Lutomirski within its <u>patch</u> proposed to allocate stack relying on vmalloc(). This had different benefits:

- resolved the problem of fragmentation (since you do not need anymore contiguous memory for the stack)
- it added graceful handling of overflows, killing the responsible process

But this had a big drawback since it introduced a 1.5µs delay in process creation which was unacceptable. Instead of improving vmalloc subsystem, Linus suggested to add a per-CPU cache of kernel-level stacks getting memory from vmalloc() has been introduced.

In the end it was also decided to move thread\_info completely off the stack and its content was moved to the task\_struct.



**10. Process Management** 1. Process Control Block

# Accessing the PCB



Advanced Operating Systems and Virtualization

#### current

current always refers to the currently-scheduled process, it is therefore architecture-specific. It returns the memory address of its PCB (evaluates to a pointer to the corresponding task\_struct).

On early versions, it was a macro current defined in include/asm-i386/current.h it performed computations based on the value of the stack pointer, by exploiting that the stack is aligned to the couple of pages/frames in memory, therefore changing the stack's size requires re-aligning this macro.

When thread\_info was introduced, masking the stack gave the address to task\_struct. To return the task\_struct, the content of the task member of task\_struct was returned.

Later, current has been mapped to the static \_\_always\_inline struct task\_struct \*get\_current(void) function. It returns the per-CPU variable current\_task declared in arch/x86/kernel/cpu/common.c. The scheduler updates the current\_task variable when executing a context switch. This is compliant with the fact that thread\_info has left the stack



```
struct task struct;
 9
10
11
     DECLARE PER CPU(struct task struct *, current task);
12
13
     static always inline struct task struct *get current(void)
14
     {
15
             return this cpu read stable(current task);
16
     }
17
18
     #define current get current()
```

https://elixir.bootlin.com/linux/v5.11/source/arch/x86/include/asm/current.h#L18

#### Up to 2.6

This function in include/linux/sched.h allows to retrieve the memory address of the PCB by passing the process/thread pid as input.

```
static inline struct task_struct *find_task_by_pid(int pid) {
    struct task_struct *p, **htable = &pidhash[pid_hashfn(pid)];
    for(p = *htable; p && p->pid != pid; p = p->pidhash_next);
    return p;
}
```

Since 2.6

find\_task\_by\_pid has been replaced by:

```
struct task_struct *find_task_by_vpid(pid_t vpid)
```

This is based on the notion of **virtual pid**. It has to do with userspace namespaces, to allow processes in different namespaces to share the same pid numbers.

A **namespace** is a feature of the Linux kernel which partitions the available resources in such a way all the process in the same namespace see the same amount of resources. At boot every process belongs to the same namespace. Namespaces are used for implementing containers. Namespaces are used in conjunction with **cgroups**, another kernel feature that limits the usage of CPU/RAM/IO for a specific set of processes.

Up to 4.14

```
/* PID hash table linkage. */
```

```
struct task_struct *pidhash_next;
```

```
struct task_struct **pidhash_pprev;
```

There is a hash defined as below in include/linux/sched.h

- #define PIDHASH\_SZ (4096 >> 2)
- extern struct task\_struct \*pid\_hash[PIDHASH\_SZ];
- #define pid\_hashfn(x) ((((x) >> 8) ^ (x)) & (PIDHASH\_SZ 1))

#### Today

The hash data structure has been replaced by a radix tree.

- PIDs are replaced with Integer IDs (idr)
- idr is the kernel-level library for the management of small integer ID numbers

An idr is a sparse array mapping integer IDs onto arbitrary pointers. Look back at the data structures lab.

10.2

10. Process Management

# The fork()/exec() model

DIAG

Advanced Operating Systems and Virtualization

# **Creating a new process**

To create a new process, a couple of fork() and exec\*() calls should be issued. In general new process share everything with the parent so it would be inefficient to truly copy all the data. To overcome this the Linux kernel:

- implements the Copy-on-Write that allows both parent and child to read the same physical pages, whenever one tries to write on a physical page the kernel copies its content into a new physical page;
- lightweight processes allow both parent and child to share many kernel data structures, such as the paging tables, open files struct and signals

Not every child need to share everything from the parent, for this reason right after a fork() we can issue an exec\*().



This function **creates a new process**. The return value is zero in the child and the process-id number of the child in the parent, or -1 upon error.

Both processes start executing from the next instruction to the fork() call.



#### **Processes and threads creation**



# Calling sys\_clone() from Userspace

Lightweight processes are created by using a function named clone().

When using sys\_clone(), we must allocate a new stack first. By convention, userspace memory is always allocated from userspace. Indeed, a thread of the same process share the same address space. Also, the TLS (Thread Local Storage) must be allocated in user space, this is architecture-dependent, thus the **unsigned long** type. glibc offers a uniform function but the implementation of the syscall entry points is slightly different on every architecture.

# sys\_fork() and sys\_clone()

```
SYSCALL_DEFINE0(fork)
{
    return _do_fork(SIGCHLD, 0, 0, NULL, NULL, 0);
}
SYSCALL_DEFINE5(clone, unsigned long, clone_flags, unsigned long, newsp, int
__user *, parent_tidptr, int __user *, child_tidptr, unsigned long, tls)
```

{

}

# do\_fork()

The do\_fork() function makes use of an auxiliary function called copy\_process() to set up the process descriptor and any other kernel data structure for child's execution. Here's the main steps:

- 1. allocates a new PID, a new PCB and a new kernel stack
- 2. copies PCB information/data structures by using copy\_process(). The information copied depends on the passed flags. That for example are:

10	#define CSIGNAL	0x000000ff /*	<pre>signal mask to be sent at exit */</pre>							
11	#define CLONE VM	0x00000100 /*	<pre>set if VM shared between processes */</pre>							
12	#define CLONE FS	0x00000200 /*	<pre>set if fs info shared between processes */</pre>							
13	#define CLONE FILES	0x00000400 /*	<pre>set if open files shared between processes */</pre>							
14	#define CLONE SIGHAND	0x0000800 /*	<pre>set if signal handlers and blocked signals shared */</pre>							
15	#define CLONE PIDFD	0x00001000 /*	<pre>set if a pidfd should be placed in parent */</pre>							
16	#define CLONE PTRACE	0x00002000 /*	set if we want to let tracing continue on the child too */							
17	#define CLONE VFORK	0x00004000 /*	set if the parent wants the child to wake it up on mm release */							
18	#define CLONE PARENT	0x00008000 /*	' set if we want to have the same parent as the cloner $\overline{*}/$							
19	#define CLONE THREAD	0×00010000 /*	Same thread group? */							
20	#define CLONE_NEWNS	0x00020000 /*	New mount namespace group */							
21	#define CLONE_SYSVSEM	0x00040000 /*	share system V SEM UNDO semantics */							
22	<pre>#define CLONE_SETTLS</pre>	0x00080000 /*	<pre>create a new TLS for the child */</pre>							
23	#define CLONE_PARENT_SE	ETTID 0×00100000	) /* set the TID in the parent */							
24	#define CLONE_CHILD_CLE	EARTID 0×00200000	) /* clear the TID in the child */							
25	#define CLONE_DETACHED	0×00400000	) /* Unused, ignored */							
26	<pre>#define CLONE_UNTRACED</pre>	0×0080000	) /* set if the tracing process can't force CLONE PTRACE on this							
27	#define CLONE_CHILD_SE	TTID 0×0100000	) /* set the TID in the child */							
28	#define CLONE_NEWCGROUP	• 0×0200000	) /* New cgroup namespace */							
29	<pre>#define CLONE_NEWUTS</pre>	0×0400000	) /* New utsname namespace */							
30	<pre>#define CLONE_NEWIPC</pre>	0×0800000	) /* New ipc namespace */							
31	<pre>#define CLONE_NEWUSER</pre>	0×1000000	) /* New user namespace */							
32	<pre>#define CLONE_NEWPID</pre>	0×2000000	) /* New pid namespace */							
33	<pre>#define CLONE_NEWNET</pre>	0×4000000	) /* New network namespace */							
34	<pre>#define CLONE_IO</pre>	0×8000000	) /* Clone io context */							
	https://elixir.bootlin.com/linux/v5.11/source/include/uapi/linux/sched.h#L11									

### do\_fork()

#### SYSCALL\_DEFINE0(fork) #ifdef CONFIG MMU struct kernel\_clone\_args args = { .exit signal = SIGCHLD, }; return kernel clone(&args); #else /\* can not support in nommu mode \*/ return -EINVAL; #endif } https://elixir.bootlin.com/linux/v5.11/source/kernel/fork.c#L2518

<pre>2414 * 0k, this is the main fork-routine. 2415 * 2416 * It copies the process, and if successful kick-starts 2417 * it and waits for it to finish using the VM if required. 2419 * args-&gt;exit_signal is expected to be checked for sanity by the calle 2420 */ 2421 pid_t kernel_clone(struct kernel_clone_args *args) 2422 { 2423 u64 clone_flags = args-&gt;flags; 2424 struct completion vfork; 2425 struct pid *pid;</pre>	
<pre>2415 * 2416 * It copies the process, and if successful kick-starts 2417 * it and waits for it to finish using the VM if required. 2418 * 2419 * args-&gt;exit_signal is expected to be checked for sanity by the calle 2420 */ 2421 pid_t kernel_clone(struct kernel_clone_args *args) 2422 { 2423 u64 clone_flags = args-&gt;flags; 2424 struct completion vfork; 2425 struct pid *pid;</pre>	
<pre>2416 * It copies the process, and if successful kick-starts 2417 * it and waits for it to finish using the VM if required. 2418 * 2419 * args-&gt;exit_signal is expected to be checked for sanity by the calle 2420 */ 2421 pid_t kernel_clone(struct kernel_clone_args *args) 2422 { 2423 u64 clone_flags = args-&gt;flags; 2424 struct completion vfork; 2425 struct pid *pid;</pre>	
<pre>2417 * it and waits for it to finish using the VM if required. 2418 * 2419 * args-&gt;exit_signal is expected to be checked for sanity by the calle 2420 */ 2421 pid_t kernel_clone(struct kernel_clone_args *args) 2422 { 2423 u64 clone_flags = args-&gt;flags; 2424 struct completion vfork; 2425 struct pid *pid;</pre>	
<pre>2418 * 2419 * args-&gt;exit_signal is expected to be checked for sanity by the calle 2420 */ 2421 pid_t kernel_clone(struct kernel_clone_args *args) 2422 { 2423 u64 clone_flags = args-&gt;flags; 2424 struct completion vfork; 2425 struct pid *pid;</pre>	
<pre>2419 * args-&gt;exit_signal is expected to be checked for sanity by the calle 2420 */ 2421 pid_t kernel_clone(struct kernel_clone_args *args) 2422 { 2423 u64 clone_flags = args-&gt;flags; 2424 struct completion vfork; 2425 struct pid *pid;</pre>	
<pre>2420 */ pid_t kernel_clone(struct kernel_clone_args *args) 2422 { 2423 u64 clone_flags = args-&gt;flags; 2424 struct completion vfork; 2425 struct pid *pid;</pre>	er.
<pre>2421 pid_t kernel_clone(struct kernel_clone_args *args) 2422 { 2423 u64 clone_flags = args-&gt;flags; 2424 struct completion vfork; 2425 struct pid *pid;</pre>	
2422 { 2423 u64 clone_flags = args->flags; 2424 struct completion vfork; 2425 struct pid *pid;	
2423 u64 clone_flags = args->flags; 2424 struct completion vfork; 2425 struct pid *pid;	
2424 struct completion vfork; 2425 struct pid *pid;	
2425 struct pid *pid;	
2426 struct task struct *p;	
2427 int trace = $0$ ;	
2428 <b>pid t nr</b> ;	
2429	
2462 D = CODY PROCESS (NULL trace NUMA NO NODE args):	
2463 add latent ontrony():	
2404 if (TS EDD(n))	
$2465$ II (IJ_ERA( $\mu$ ))	
2466 return <b>Pik_ekk</b> ( <b>p</b> );	
<pre>2474 pid = get task pid(p, PIDTYPE PID);</pre>	
2475 $\mathbf{nr} = \mathbf{pid} \mathbf{vnr}(\mathbf{pid})$ :	
2476	
2490	
2497 <b>put pid(pid)</b> ;	
2498 return <b>nr</b> ;	
2499 }	

https://elixir.bootlin.com/linux/v5.11/source/kernel/fork.c#L2518
## copy\_process()

The function implements several checks on namespaces. Pending signals are processed immediately in the parent process.

- p = dup\_task\_struct(current, node);
- setup\_thread\_stack(tsk, orig);
- copy\_creds(p, clone\_flags);
- copy\_files(clone\_flags, p);
- copy\_fs(clone\_flags, p);
- copy\_mm(clone\_flags, p); -> dup\_mm()

## dup\_mm()

	1335	/**
	1336	<pre>* dup_mm() - duplicates an existing mm structure</pre>
	1337	* @tsk: the task_struct with which the new mm will be associated.
	1338	* @oldmm: the mm to duplicate.
	1339	*
	1340	* Allocates a new mm structure and duplicates the provided @oldmm structure
	1341	* content into it.
	1342	*
	1343	* Return: the duplicated mm or NULL on failure.
	1344	*/
	1345	<pre>static struct mm_struct *dup_mm(struct task_struct *tsk,</pre>
	1346	struct <b>mm_struct</b> *oldmm)
	1347	{
	1348	struct mm_struct *mm;
	1349	int err;
	1350	
	1351	<pre>mm = allocate_mm();</pre>
	1352	lT (!mm)
	1353	goto fall_nomem;
	1354	
Allocates a new PGD	1355	memcpy(mm, olamm, sizeor(*mm));
Allocates a new I GD	1257	$\rightarrow$ if (1mm init(mm tak mm $\rightarrow$ usor no))
	1358	
	1350	goto rate_noment,
	1360	err = dun mman(mmoldmm);
	1361	if (err)
	1362	acto free pt:
	1363	5 <b>-</b>
	1364	<pre>mm-&gt;hiwater rss = get mm rss(mm);</pre>
	1365	mm->hiwater vm = mm->total vm;

https://elixir.bootlin.com/linux/v5.11/source/kernel/fork.c#L1345

10.2.1

10. Process Management
2. The fork()/exec() model

## **Kernel Threads**

Advanced Operating Systems and Virtualization

### **Kernel Thread Creation API**

```
struct task struct *kthread create on node(int (*threadfn)(void *data),
This is seen as
                      12
                                                                       void *data,
                      13
any other task
                                                                       int node,
                      14
15
                                                                       const char namefmt[], ...);
by the scheduler
                      16
                           /**
                      17
                            * kthread create - create a kthread on the current node
                      18
                            * @threadfn: the function to run in the thread
                      19
                            * @data: data pointer for @threadfn()
                            * @namefmt: printf-style format string for the thread name
                      21
                            * @arg...: arguments for @namefmt.
                      22
                            *
                      23
                            * This macro will create a kthread on the current node, leaving it in
                      24
                            * the stopped state. This is just a helper for kthread create on node();
                      25
                            * see the documentation there for more details.
                      26
                            */
                      27
                           #define kthread create(threadfn, data, namefmt, arg...) \
                                   kthread create on node(threadfn, data, NUMA NO NODE, namefmt, ##arg)
                      28
                      29
                      30
                      31
                           struct task struct *kthread create on cpu(int (*threadfn)(void *data),
                      32
                                                                      void *data.
                      33
                                                                      unsigned int cpu,
                      34
                                                                      const char *namefmt);
```

https://elixir.bootlin.com/linux/v5.11/source/include/linux/kthread.h#L27

Kthreads are always stopped upon creation, they must be activated by calling wake\_up\_process().

V5.11

### **Kernel Thread Daemon**

```
330
      struct task struct * kthread create on node(int (*threadfn)(void *data),
331
                                                            void *data, int node,
332
                                                             const char namefmt[],
333
                                                             va list args)
334
      {
335
              DECLARE COMPLETION ONSTACK(done);
               struct task struct *task;
336
337
               struct kthread create info *create = kmalloc(sizeof(*create),
338
                                                              GFP KERNEL);
339
340
              if (!create)
341
                       return ERR PTR(-ENOMEM);
342
              create -> threadfn = threadfn:
343
              create->data = data;
344
              create->node = node:
345
              create -> done = \& done;
346
347
               spin lock(&kthread create lock);
               list add tail(&create->list, &kthread create list);
348
349
               spin unlock(&kthread create lock);
350
351
              wake up process(kthreadd task);
```

https://elixir.bootlin.com/linux/v5.11/source/kernel/kthread.c#L330

#### **Kernel Thread Daemon**

630	<pre>int kthreadd(void *unused)</pre>
631	{
632	<pre>struct task_struct *tsk = current;</pre>
633	
634	/* Setup a clean context for our children to inherit. */
635	<pre>set_task_comm(tsk, "kthreadd");</pre>
636	<pre>ignore_signals(tsk);</pre>
637	set_cpus_allowed_ptr(tsk, housekeeping_cpumask(HK_FLAG_KTHREAD));
638	<pre>set_mems_allowed(node_states[N_MEMORY]);</pre>
639	
640	<pre>current-&gt;flags  = PF_NOFREEZE;</pre>
641	cgroup_init_kthreadd();
642	
643	for (;;) {
644	<pre>set_current_state(TASK_INTERRUPTIBLE);</pre>
645	if ( <b>list_empty</b> (&kthread_create_list))
646	schedule();
647	<pre>set_current_state(TASK_RUNNING);</pre>
648	
649	<pre>spin_lock(&amp;kthread_create_lock);</pre>
650	while (! <b>list_empty</b> (&kthread_create_list)) {
651	struct <b>kthread_create_info</b> * <b>create</b> ;
652	
653	<pre>create = list_entry(kthread_create_list.next,</pre>
654	struct kthread_create_info, list);
655	<pre>list_del_init(&amp;create-&gt;list);</pre>
656	<pre>spin_unlock(&amp;kthread_create_lock);</pre>
657	
658	<pre>create_kthread(create);</pre>
659	
660	<b>spin_lock</b> (&kthread_create_lock);
662	}
662	spin_unicck(@kthread_create_tock);
664	ł
665	
666	
000	J https://elixir.bootlin.com/linux/v5.11/source/kernel/kthread.c#L630



10. Process Management

# Out Of Memory (OOM) Killer



Advanced Operating Systems and Virtualization

## **The OOM Killer**

It is implemented in mm/oom\_kill.c. This module is activated (if enabled) when the system runs out of memory.

There are three possible actions:

- kill a random task (bad)
- let the system crash (worse)
- try to be smart at picking the process to kill

The OOM Killer picks a "good" process and kills it in order to reclaim available memory.

## **The OOM Killer**

Entry point of the system is out\_of\_memory(). It tries to select the "best" process checking for different conditions:

- if a process has a pending SIGKILL or is exiting, this is automatically picked (check done by task\_will\_free\_mem())
- Otherwise, it issues a call to select\_bad\_process() which will return a process to be killed:
  - the picked process is then killed
  - if no process is found, a panic() is raised

## select\_bad\_process()

This iterates over all available processes calling oom\_evaluate\_task() on them, until a killable process is found. Unkillable tasks (i.e., kernel threads) are skipped, oom\_badness() implements the heuristic to pick the process to be killed by computing the "score" associated with each process, the higher the higher the score the higher the probability of getting killed.

#### oom\_badness()

A score of zero is given if:

- the task is unkillable
- the mm field is NULL
- if the process is in the middle of a fork

The score is then computed proportionally to the RAM, swap, and pagetable usage:

10.4

10. Process Management

# **Process Starting**



Advanced Operating Systems and Virtualization

### How a Program is Started?

We all know how to compile a program:

gcc program.c -o program

We all know how to launch the compiled program:

./program

The question is: why does all this work? What is the *convention* used between kernel and user space?



## Starting a program from bash

From the bash shell source - https://github.com/bminor/bash/blob/master/execute\_cmd.c

## Starting a program from bash

```
pid_t make_child (char *command, int async_p) {
    pid t pid;
    int forksleep;
    start pipeline();
    forksleep = 1;
    while ((pid = fork ()) < 0 && errno == EAGAIN && forksleep < FORKSLEEP MAX) {</pre>
        sys error("fork: retry");
        reap zombie children();
        if (forksleep > 1 && sleep(forksleep) != 0)
            break:
        forksleep <<= 1;</pre>
    }
    /* ... */
    return (pid);
}
```

## Starting a program from bash

```
int shell execve (char *command, char **args, char **env) {
    execve (command, args, env);
    READ SAMPLE BUF (command, sample, sample len);
    if (sample len == 0)
         return (EXECUTION SUCCESS);
    if (sample len > 0) {
         if (sample len > 2 && sample[0] == '#' && sample[1] == '!')
             return (execute shell script(sample, sample len, command, args, env));
         else if (check binary file (sample, sample len)) {
             internal_error (_("%s: cannot execute binary file"), command);
             return (EX BINARY FILE);
                                                    * These are extern so execute_simple_command can set them, and then
         }
                                                      longimp back to main to execute a shell script, instead of calling
                                                      main () again and resulting in indefinite, possibly fatal, stack
                                                      arowth. */
                                                   procenv t subshell top level;
                                                                                   https://github.com/bminor/bash/blob/master/shell.c
    longjmp(subshell top level, 1)
```



**exec\*()** changes the program file that an existing process is running:

- it first **wipes** out the memory state of the calling process
- it then goes to the filesystem to **find** the program file requested
- it **copies** this file into the program's memory and initializes register state, including the PC
- It **doesn't alter** most of the other fields in the PCB. The process calling exec\*() (the child copy of the shell, in this case) can, e.g., change the opened files

Let's see how exec\*() is implemented.

### struct linux\_binprm

The **struct** linux\_binprm is in charge of keeping information about a binary file.

```
struct linux binprm {
    char buf BINPRM BUF SIZE;
    struct page *page MAX ARG PAGES];
    unsigned long p; /* current top of mem */
    int sh bang;
    struct file* file:
    int e uid, e gid;
    kernel cap t cap inheritable, cap permitted, cap effective;
    int argc, envc;
    char *filename; /* Name of binary */
    unsigned long loader, exec;
};
```

```
In kernel 5.11 do execve atcommon()
do_execve()
                                               ttps://elixir.bootlin.com/linux/v5.11/source/fs/exec.c#L1855
      int do execve(char *filename, char **argv, char **envp, struct pt_regs *regs) {
            struct linux_binprm bprm;
            struct file *file:
            int retval:
            int i:
            file = open exec(filename);
            retval = PTR ERR(file);
            if (IS ERR(file))
                  return retval;
            bprm.p = PAGE SIZE*MAX ARG PAGES-sizeof(void *);
            memset(bprm.page, 0, MAX ARG PAGES*sizeof(bprm.page[0]));
            bprm.file = file;
            bprm.filename = filename;
                                                                              count(argv, max) counts
            bprm.sh bang = 0;
                                                                              the number of strings
            bprm.loader = 0;
            bprm.exec = 0:
            if ((bprm.argc = count(argv, bprm.p / sizeof(void *))) < 0) {</pre>
                  allow write access(file);
                  fput(file);
                  return bprm.argc;
```

```
do_execve()
                     if ((bprm.envc = count(envp, bprm.p / sizeof(void *))) < 0) {</pre>
                           allow write access(file);
                           fput(file);
                           return bprm.envc;
                     }
                     retval = prepare binprm(&bprm);
                     if (retval < 0)
                           goto out:
                     retval = copy strings kernel(1, &bprm.filename, &bprm);
                     if (retval < 0)
                           goto out:
                     bprm.exec = bprm.p;
                     retval = copy strings(bprm.envc, envp, &bprm);
                     if (retval < 0)
                           goto out;
                     retval = copy strings(bprm.argc, argv, &bprm);
                     if (retval < 0)
                           doto out:
                     retval = search binary handler(&bprm,regs);
                     if (retval >= 0)
                           /* execve success */
                           return retval;
```

## do\_execve()

#### out:

```
/* Something went wrong, return the inode and free the argument pages*/
allow_write_access(bprm.file);
if (bprm.file)
    fput(bprm.file);
for (i = 0 ; i < MAX_ARG_PAGES ; i++) {
    struct page * page = bprm.page[i];
    if (page)
        ___free_page(page);
}</pre>
```

return retval;

## search\_binary\_handler()

The function scans a list of binary file handlers registered in the kernel. If no handler is able to recognize the image format, syscall returns the ENOEXEC error ("Exec Format Error").

For ELF files we have in fs/binfmt\_elf.c:

- load\_elf\_binary(), the function:
  - loads image file to memory using mmap;
  - reads the program header and sets permissions accordingly
  - elf\_ex = \*((struct elfhdr \*)bprm->buf);

10.4.1

**10. Process Management** 4. Process Starting

## **The ELF Format**

DIAG

Advanced Operating Systems and Virtualization

## **ELF: Executable and Linking Format**

The Executable and Linking Format was originally developed and published by UNIX System Laboratories (USL) as part of the Application Binary Interface (ABI). The Tool Interface Standards committee (TIS) has selected the evolving ELF standard as a portable object file format that works on 32-bit Intel Architecture environments for a variety of operating systems.

The ELF standard is intended to streamline software development by providing developers with a set of binary interface definitions that extend across multiple operating environments. This should reduce the number of different interface implementations, thereby reducing the need for recoding and recompiling code.

## **The Compiling Process**



#### **ELF Format**

ELF defines the format of binary executables. There are four different categories:

- **Relocatable**, created by compilers and assemblers. Must be processed by the linker before being run.
- **Executable**, all symbols are resolved, except for shared libraries' symbols, which are resolved at runtime.
- **Shared object**, a library which is shared by different programs, contains all the symbols' information used by the linker, and the code to be executed at runtime.
- **Core file,** a core dump.

ELF files have a twofold nature

- compilers, assemblers and linkers handle them as a set of logical sections;
- the system loader handles them as a set of **segments**.

#### **ELF Format**







https://upload.wikimedia.org/wikipedia/commons/e/e4/ELF Executable and Linkable Format diagram by Ange Albertini.png

#### **Relocatable File**

A relocatable file or a shared object is a collection of sections. Each section contains a single kind of information, such as executable code, read-only data, read/write data, relocation entries, or symbols.

Each symbol's address is defined in relation to the section which contains it. For example, a function's entry point is defined in relation to the section of the program which contains it.

#### **Section Header**

#### typedef struct {

```
Elf32 Word sh_name; /* Section name (string tbl index) */
    Elf32 Word sh_type; /* Section type */
    Elf32 Word sh flags; /* Section flags */
    Elf32 Addr sh addr; /* Section virtual addr at execution */
    Elf32 Off sh offset; /* Section file offset */
    Elf32 Word sh size; /* Section size in bytes */
    Elf32 Word sh link; /* Link to another section */
    Elf32 Word sh info; /* Additional section information */
    Elf32 Word sh addralign; /* Section alignment */
    Elf32 Word sh entsize; /* Entry size if section holds table */
} Elf32 Shdr;
```

## **Section Header**

#### Types and Flags

Types:

- PROGBITS: The section contains the program content (code, data, debug information).
- NOBITS: Same as PROGBITS, yet with a null size.
- SYMTAB and DYNSYM: The section contains a symbol table.
- STRTAB: The section contains a string table.
- REL and RELA: The section contains relocation information.
- DYNAMIC and HASH: The section contains dynamic linking information.

Flags:

- WRITE: The section contains runtime-writeable data.
- ALLOC: The section occupies memory at runtime.
- EXECINSTR: The section contains executable machine instructions.

## Sections

#### Examples

- .text: contains program's instructions
  - Type: PROGBITS
  - Flags: ALLOC + EXECINSTR
- .data: contains pre-initialized read/write data
  - Type: PROGBITS
  - Flags: ALLOC + WRITE
- .rodata: contains pre-initialized read-only data
  - Type: PROGBITS
  - Flags: ALLOC
- .bss: contains uninitialized data. Will be set to zero at startup.
  - Type: NOBITS
  - Flags: ALLOC + WRITE

#### **Executable Files**

Usually, an executable file has only few segments:

- A read-only segment for code.
- A read-only segment for read-only data.
- A read/write segment for other data.

Any section marked with flag ALLOC is packed in the proper segment, so that the operating system is able to map the file to memory with few operations.

If .data and .bss sections are present, they are placed within the same read/write segment.

## **Program Header**

#### typedef struct {

```
Elf32 Word p type; /* Segment type */
    Elf32 Off p offset; /* Segment file offset */
    Elf32 Addr p vaddr; /* Segment virtual address */
    Elf32 Addr p paddr; /* Segment physical address */
    Elf32 Word p filesz; /* Segment size in file */
    Elf32 Word p memsz; /* Segment size in memory */
    Elf32 Word p_flags; /* Segment flags */
    Elf32 Word p align; /* Segment alignment */
} Elf32 Phdr;
```

## Linker's Role



## **Static Relocation**



 data section

 732e
 6d79
 6174
 0062
 732e
 7274
 6174
 0062

 732e
 7368
 7274
 6174
 0062
 742e
 7865
 0074

 642e
 7461
 0061
 622e
 7373
 6174
 0062
 7865



**info** tells you the index in the symbol table and the type of the symbol



name is the position of the name in the string table
#### Symbols Visibility

A symbol can be:

- **strong**, a strong symbol replaces a weak one and if two strong symbols have the same name the linker resolves in favour of the first; by default every symbol is strong
- weak, more modules can have a symbol with the same name of a weak one, the declared entity cannot be overloaded by other modules; It is useful for libraries which want to avoid conflicts with user programs.

gcc version 4.0 gives the command line option -fvisibility:

- default: normal behaviour, the symbol is seen by other modules;
- hidden: two declarations of an object refer the same object only if they are in the same shared object;
- internal: an entity declared in a module cannot be referenced even by pointer;
- protected: the symbol is weak;

#### Symbols Visibility

int variable \_\_attribute\_\_ ((visibility ("hidden")));



10.4.2

**10. Process Management** 4. Process Starting

# **Dynamic Linking**



Advanced Operating Systems and Virtualization

#### **Program Entry Point**

The main() function is not the actual entry point for the program. glibc inserts auxiliary functions. The actual entry point is called \_start.

The Static Relocation works at linking time but you obviously do not want to include all the libraries that you use in your program in your executable file, this because eats up memory and almost all the programs use the same set of libraries (e.g. the stdlib). Symbols that are not included in the final executable file are resolved with the Dynamic Linking that is performed by the kernel when the program starts.

The Kernel starts the dynamic linker which is stored in the .interp section of the program (usually /lib/ld-linux.so.2). If no dynamic linker is specified, control is given at address specified in e\_entry.

Initialization steps:

- Self initialization
- Loading Shared Libraries
- Resolving remaining relocations
- Transfer control to the application

The most important data structures which are filled are:

- **Procedure Linkage Table** (PLT), used to call functions whose address isn't known at link time
- Global Offsets Table (GOT), similarly used to resolve addresses of data/functions

#### **Data Structures**

- .dynsym: a minimal symbol table used by the dynamic linker when performing relocations
- .hash: a hash table that is used to quickly locate a given symbol in the .dynsym, usually in one or two tries.
- .dynstr: string table related to the symbols stored in .dynsym

These tables are used to fill the GOT table, that is populated upon need (lazy binding).

#### Steps

The first PLT entry is special. Other entries are identical, one for each function needing resolution.

- 1. A jump to a location which is specified in a corresponding GOT entry
- 2. Preparation of arguments for a resolver routine
- 3. Call to the resolver routine, which resides in the first entry of the PLT

The first PLT entry is a call to the *resolver* located in the dynamic loader itself.

#### Steps

When func is called for the first time:

- 1. PLT[n] is called, and jumps to the address pointed to it in GOT[n]
- 2. This address points into PLT[n] itself, to the preparation of arguments for the resolver.
- 3. The resolver is then called, by jumping to PLT[o]
- The resolver performs resolution of the actual address of func, places its actual address into GOT[n] and calls func.



Steps after the first resolution





**10. Process Management** 4. Process Starting

#### **Initial Steps of Programs' Life**



Advanced Operating Systems and Virtualization

#### **Initial Steps**

So far the dynamic linker has loaded the shared libraries in memory. GOT is populated when the program requires certain functions. Then, the dynamic linker calls \_start



#### **Userspace Life of a Program**



#### Stack Layout at Program Startup

local variables of main saved registers of main return address of main argc argv envp stack from startup code argc argv pointers NULL that ends argv[] environment pointers NULL that ends envp[] ELF Auxiliary Table argv strings environment strings program name NULL

main()

\_\_libc\_start\_main()

#### kernel

# Advanced Operating Systems and Virtualization

[10] Process Management

LECTURER Gabriele **Proietti Mattia** 

BASED ON WORK BY http://www.ce.uniroma2.it/~pellegrini/



gpm.name · proiettimattia@diag.uniroma1.it

