

Trabajos Prácticos con JOS

Federico del Mazo - 100029

Rodrigo Souto - 97649

Trabajos Prácticos con JOS

Respuestas teóricas de los distintos trabajos prácticos/labs de Sistemas Operativos (75.08).

TP1: Memoria virtual en JOS (26/04/2019)

Memoria física: `boot__alloc__pos`

1. Incluir: Un cálculo manual de la primera dirección de memoria que devolverá `boot__alloc()` tras el arranque. Se puede calcular a partir del binario compilado (`obj/kern/kernel`), usando los comandos `readelf` y/o `nm` y operaciones matemáticas.

Truncando la salida de ambos comandos (con `grep`), vemos las siguientes líneas:

```
sisop_2019a_delmazo_souto TP1 % readelf -s obj/kern/kernel | grep end
112: f0117950      0 NOTYPE GLOBAL DEFAULT      6 end
sisop_2019a_delmazo_souto TP1 % nm obj/kern/kernel | grep end
f0117950 B end
```

Como podemos ver en ambos casos, la dirección de memoria que recibe `boot__alloc()` es `f0117950` (en decimal, `4027677008`). A este valor, la función (en su primera llamada) lo redondea a 4096 (`PGSIZE`) llamando a `ROUNDUP(a,n)`. Además de devolver ese valor redondeado, guarda la variable `nextfree` en una página más de lo recibido.

Por ende, el valor devuelto será el de `ROUNDUP(4027677008, 4096)`. Esta función, como indica su documentación, redondea a `a` al múltiplo más cercano de `n`. Este múltiplo será `4027678720`, que está más cerca que el siguiente múltiplo (`4027682816`). Para confirmarlo desde la práctica, se traducen a Python los cálculos que utiliza la función de redondeo:

```
sisop_2019a_delmazo_souto TP1 % python3
Python 3.6.7 |Anaconda, Inc.| (default, Oct 23 2018, 19:16:44)
>>> a = 0xf0117950
>>> n = 4096
>>> def rounddown(a,n): return a - a % n
>>> def roundup(a,n): return rounddown(a + n - 1, n)
>>> res = roundup(int(a),n)
>>> res
4027678720
>>> hex(res)
'0xf0118000'
```

2. Incluir: Una sesión de GDB en la que, poniendo un breakpoint en la función `boot__alloc()`, se muestre el valor de `end` y `nextfree` al comienzo y fin de esa primera llamada a `boot__alloc()`.

```
sisop_2019a_delmazo_souto TP1 % make gdb
gdb -q -s obj/kern/kernel -ex 'target remote 127.0.0.1:26000' -n -x .gdbinit
Reading symbols from obj/kern/kernel...done.
Remote debugging using 127.0.0.1:26000
0x0000fff0 in ?? ()
```

```

(gdb) break boot_alloc
Breakpoint 1 at 0xf0100a58: file kern/pmap.c, line 89.
(gdb) continue
Continuing.
The target architecture is assumed to be i386
=> 0xf0100a58 <boot_alloc>: push    %ebp

Breakpoint 1, boot_alloc (n=4096) at kern/pmap.c:89
89 {
(gdb) print (char*) &end
$1 = 0xf0117950 ""
(gdb) watch &end
Watchpoint 2: &end
(gdb) watch nextfree
Hardware watchpoint 3: nextfree
(gdb) continue
Continuing.
=> 0xf0100aac <boot_alloc+84>: jmp     0xf0100a68 <boot_alloc+16>

Hardware watchpoint 3: nextfree

Old value = 0x0
New value = 0xf0118000 ""
0xf0100aac in boot_alloc (n=4096) at kern/pmap.c:100
100     nextfree = ROUNDUP((char *) end, PGSIZE);
(gdb) continue
Continuing.
=> 0xf0100a81 <boot_alloc+41>: mov     0xf0117944,%edx

Hardware watchpoint 3: nextfree

Old value = 0xf0118000 ""
New value = 0xf0119000 ""
boot_alloc (n=4096) at kern/pmap.c:111
111     if (nextfree >= (char *) (KERNBASE + npages * PGSIZE)) {
(gdb) continue
Continuing.

```

Como se puede ver, se cumple todo lo planteado. `end` comienza en `0xf0117950`, luego `nextfree` se inicializa en el número ya redondeado `0xf0118000`, y finalmente se avanza una página, y queda `0xf0119000`.

Memoria física: `page__alloc`

1. Responder: ¿en qué se diferencia `page2pa()` de `page2kva()`?

Como bien indican sus nombres, `page2pa()` y `page2kva()` se diferencian en el valor de retorno. Ambas reciben una página física, pero `page2pa()` devuelve su dirección física (de tipo `physaddr_t`) mientras que `page2kva()` devuelve la dirección virtual (kernel virtual address), de tipo `void*`.

Incluso, `page2kva()` no es más que un llamado a `page2pa()` y luego a la función del preprocesador `KADDR()` que recibe una dirección física y devuelve la respectiva dirección virtual.

Large pages: `map__region__large`

1. Responder: ¿cuánta memoria se ahorró de este modo? ¿Es una cantidad fija, o depende de la memoria física de la computadora?

Se ahorran 4KB, que es el tamaño de una página, ya que se deja de usar `entry_pgtable` y se mapea la misma cantidad de memoria consecutiva (4MB) directamente con una *large page*.

Debido a que JOS se compila con la arquitectura de 32 bits i386, independientemente de cuál sea la memoria física disponible de la máquina, las páginas tendrán un tamaño de 4KB, y cómo lo que sea ahorra es crear a `entry_pgtable`, que tiene el tamaño de una página, se ahorra esa cantidad de bytes.

TP2: Procesos de usuario (17/5/2019)

Inicializaciones: `env__alloc`

1. Responder: ¿Qué identificadores se asignan a los primeros 5 procesos creados? (Usar base hexadecimal.)

La generación de ids de entornos se logra con la siguiente porción de código:

```
generation = (e->env_id + (1 << ENVGENSHIFT)) & ~(NENV - 1);
e->env_id = generation | (e - envs);
```

Analíticamente, sabiendo que `ENVGENSHIFT` equivale a 12, y `NENV` es `1 << 10` que es 1024, se puede ver que `generation` equivale al id anterior más 4096 y a eso aplicar el AND de bits con `~(1023)`. Luego, en la segunda línea, se le suma a este valor el número de entorno. Con un id de entorno igual a 0, `generation` será 4096 (`0x1000` en hexadecimal).

Efectivamente, gracias a GDB se puede comprobar que el valor de `generation` en la primera ejecución es el `0x1000` (esto es así ya que en la primera corrida los id de los entornos son 0, una vez que se empiecen a reciclar los entornos el valor de `generation` dependerá del id del entorno previo). Por ende, luego de la suma del offset, los primeros 5 entornos serán: `0x1000`, `0x1001`, `0x1002`, `0x1003` y `0x1004`.

2. Responder: Supongamos que al arrancar el kernel se lanzan `NENV` procesos a ejecución. A continuación se destruye el proceso asociado a `envs[630]` y se lanza un proceso que cada segundo muere y se vuelve a lanzar. ¿Qué identificadores tendrá este proceso en sus primeras cinco ejecuciones?

El primer proceso lanzado no será mas que la suma entre el primer `generation` (`0x1000`) y el offset 630 (en hexadecimal, `0x0276`). Es decir, el `0x1276`.

Luego, una vez que este proceso muera y se relance, se utilizara este id como base para el nuevo `generation`. El `generation` nuevo será entonces la suma entre `0x1276` (el previo id), la constante `0x1000`, y a eso el AND con `~(0x03FF)`. Luego, a ese número se le suma nuevamente el `0x0276`.

Para los primeros 5 procesos queda:

```
>>> def generate_id(id_prev): return hex( (id_prev + 0x1000 & ~(0x03FF)) + 0x0276 )
>>> generate_id(0x0)
'0x1276'
>>> generate_id(0x1276)
'0x2276'
>>> generate_id(0x2276)
'0x3276'
>>> generate_id(0x3276)
'0x4276'
>>> generate_id(0x4276)
'0x5276'
```

Inicializaciones: `env__init_percpu`

1. Responder: ¿Cuántos bytes escribe la función `lgdt`, y dónde?

`lgdt` escribe 6 bytes en la Global Descriptor Table.

2. Responder: ¿Qué representan esos bytes?

Estos bytes son el tamaño del GDT (2 bytes) y la dirección de la tabla (4 bytes).

Lanzar procesos: env_pop_tf

1. Responder: ¿Qué hay en (%esp) tras el primer movl de la función?
2. Responder: ¿Qué hay en (%esp) justo antes de la instrucción iret? ¿Y en 8(%esp)?
3. Responder: ¿Cómo puede determinar la CPU si hay un cambio de ring (nivel de privilegio)?

Lanzar procesos: gdb_hello

1. Incluir una sesión de GDB con diversos pasos:
 - paso 1
 - paso 2
 - etc

Interrupts y syscalls: kern_idt

1. Responder: ¿Cómo decidir si usar TRAPHANDLER o TRAPHANDLER_NOEC? ¿Qué pasaría si se usara solamente la primera?
2. Responder: ¿Qué cambia, en la invocación de handlers, el segundo parámetro (istrap) de la macro SETGATE? ¿Por qué se elegiría un comportamiento u otro durante un syscall?
3. Responder: Leer user/softint.c y ejecutarlo con make run-softint-nox. ¿Qué excepción se genera? Si es diferente a la que invoca el programa... ¿cuál es el mecanismo por el que ocurrió esto, y por qué motivos?

Protección de memoria: user_evilhello

Se guarda el siguiente programa en evilesthello.c:

```
#include <inc/lib.h>

void
umain(int argc, char **argv)
{
    char *entry = (char *) 0xf010000c;
    char first = *entry;
    sys_cputs(&first, 1);
}
```

1. Responder: ¿En qué se diferencia el código de la versión en evilhello.c con evilesthello.c?

Como se puede observar, la diferencia esta sencillamente en que el evilhello.c original pasa la dirección de memoria sin modificar, mientras que el modificado antes de eso cambia el puntero (desreferencia, y luego pasa la referencia).

Específicamente al asignar la variable first a entry (es decir, la línea char first = *entry;), se procede a ‘engañar’ al sistema operativo (con nada más que un swap) y se logra acceder a una dirección privilegiada. De ser la ejecución satisfactoria, revelaría una gran vulnerabilidad en el sistema: se puede imprimir todo lo que contenga el kernel!

2. Responder: ¿En qué cambia el comportamiento durante la ejecución? ¿Por qué? ¿Cuál es el mecanismo?

Gracias al swap de direcciones, la version modificada (aun más malvada) de evilhello.c sí logró imprimir el entry point del kernel como cadena (la versión original no imprimió más que simbolos basura). Esto no debería pasar y se tiene que atrapar de alguna manera (con los assertions del user_mem_check), ya que la dirección de memoria a la que se accede es una dirección privilegiada para el usuario y se debe prohibir el acceso a esta.

Ejecución de `evilhello.c` (versión original):

```
[00000000] new env 00001000
Incoming TRAP frame at 0xeffffbc
fr Incoming TRAP frame at 0xeffffbc
[00001000] exiting gracefully
[00001000] free env 00001000
Destroyed the only environment - nothing more to do!
```

Ejecución de `evilesthelloc.c` (versión modificada):

```
[00000000] new env 00001000
Incoming TRAP frame at 0xeffffbc
[00001000] user fault va f010000c ip 00800039
TRAP frame at 0xf01c1000
edi 0x00000000
esi 0x00000000
ebp 0xeebdfdf0
oesp 0xeffffdc
ebx 0x00000000
edx 0x00000000
ecx 0x00000000
eax 0x00000000
es 0x----0023
ds 0x----0023
trap 0x0000000e Page Fault
cr2 0xf010000c
err 0x00000005 [user, read, protection]
eip 0x00800039
cs 0x----001b
flag 0x00000082
esp 0xeebdfdb0
ss 0x----0023
[00001000] free env 00001000
Destroyed the only environment - nothing more to do!
```