

# Implementación Sencilla de un Hypervisor con soporte Intel VT-x

## Creación de un entorno invitado

**sys\_env\_mkguest()** [kern/syscall.c] comprueba si el procesador tiene soporte de VMX y crea un nuevo entorno invitado.

En **kern/env.c** se implementan las funciones de gestión de los entornos: **envid2env()**, **env\_init()**, **env\_init\_percpu()**, **env\_setup\_vm()**, **env\_guest\_alloc()**, **env\_guest\_free()**, **env\_alloc()**, **region\_alloc()**, **load\_icode()**, **env\_create()**, **env\_free()**, **env\_destroy()**, **env\_pop\_tf()** y **env\_run()**.

Función **vmx\_check\_support()** [vmm/vmx.c] para comprobar si la cpu tiene soporte VMX. Para ello, hay que leer el bit 5 del registro ECX (si está a 1, hay soporte) a través de la función CPUID.

Función **vmx\_check\_ept()** [vmm/vmx.c] para comprobar si está soportada la paginación extendida (EPT). Para ello, hay que comprobar en el registro MSR IA32\_VMX\_PROCBASED\_CTLs el bit 63 que indica si los controles VMX secundarios están activados (bit a 1). Además, se comprueba en el registro MSR IA32\_VMX\_PROCBASED\_CTLs2 el bit 33 que indica si el EPT está disponible.

Función **sched\_yield()** [kern/sched.c] para que un hilo renuncie a la CPU y de paso al siguiente. Función **vmxon()** [kern/sched.c] para arrancar el modo de ejecución: vmx root.

Doble traducción de dirección:

app -> VMM -> host

va -> guest\_pa -> pa

## Mapeo del bootloader y kernel del huésped

Función **sys\_ept\_map()** [kern/syscall.c] para mapear una página del VMM a una MV. Para ello, hay que comprobar que:

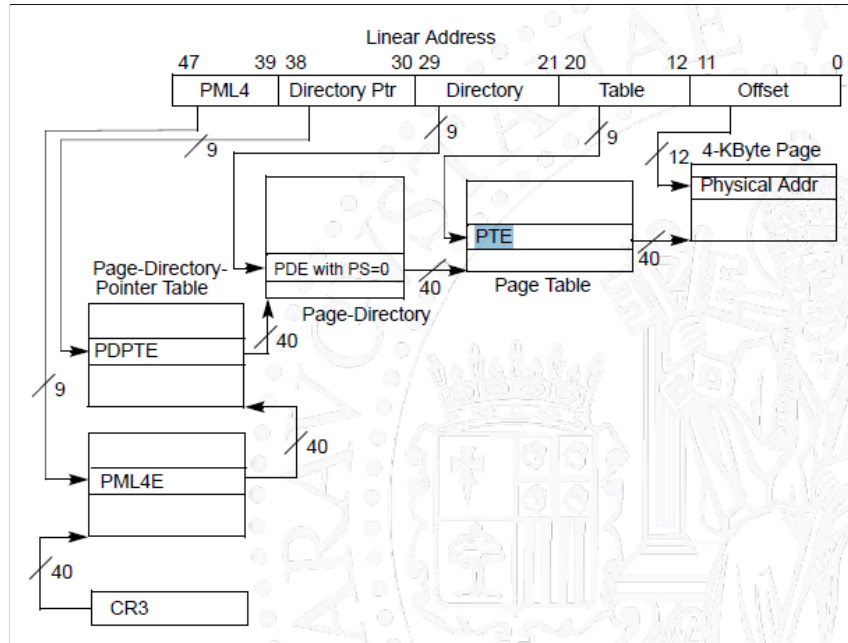
1. Existen los entornos
2. La alineación de las páginas y si superan los límites
3. Existe la página a mapear
4. Los permisos son correctos
5. Se mapea y se comprueba si ha ido bien

Función **ept\_page\_insert()** [vmm/ept.c] añade la entrada de la página mapeada en el EPT. Para ello, hay que:

1. Buscar el PTE final a partir de la dirección gpa (linear address).
2. Si
  - a. está ocupada la entrada PTE por una página:
    - i. Si `overwrite` es 0, no se sobrescribe la página.
    - ii. Si `overwrite` es 1, se sobrescribe la página.
  - b. no está ocupada, se escribe la página.

Función **ept\_lookup\_gpa()** [vmm/ept.c] devuelve el PTE final a partir de una gpa (guest physical address). Para ello, hay que:

1. Comprobar que exista eptrt (EPT root)
2. Con la función pml4e\_walk() invoca el Page Table Walker desde el nivel más alto (pml4, page map level 4) que se encarga de obtener a partir de la gpa (linear address) el PTE final.
3. Antes de devolver el PTE final en epte\_out, hay que establecer los permisos de las entradas EPT intermedias (PML4E, PDPTE y PDE) a **\_\_EPT\_FULL** (porque lo pone en la cabecera).



Función **map\_in\_guest()** [usr/vmm.c] para mapear una región de un fichero (ejecutable) en la MV. Para ello, hay que comprobar que:

1. Con PGOFF se comprueba que la dirección gpa está alineada.
  - a. Si no está alineada, se alinea con ROUNDDOWN
2. En un bucle se recorre toda la región a mapear (de tamaño filesz) saltando de page en page:
  - a. Se reserva espacio en UTEMP (espacio de memoria temporal)
  - b. Se busca la dirección base de la región en el file.
  - c. Se lee una página de la región del file.
  - d. Se escribe la página leída en UTEMP.
  - e. Se escribe la página de UTEMP en el guest.
  - f. Se elimina la página de UTEMP (haciendo que solo haya una página mapeada en UTEMP en un instante dado).
3. En otro bucle se recorre el resto de memoria (hasta memsz) saltando de page en page:
  - a. Se coge una page vacía de UTEMP (todo a 0's).
  - b. Se escribe la página vacía de UTEMP en el guest.
  - c. Se elimina la página de UTEMP.

Función **copy\_guest\_kern\_gpa()** [**usr/vmm.c**] para mapear un fichero (ejecutable) en la MV, es decir, que lee la cabecera del ELF y después invoca recursivamente a **map\_in\_guest()** para mapear todas las regiones. Para ello, hay que:

1. Abrir el fichero ELF con permisos de lectura.
2. Leer los 512 bytes que ocupa el gestor de arranque (según el enunciado).
3. Se comprueba que el número mágico de la cabecera corresponda con el número mágico de los ficheros ELF (ELF\_MAGIC). Este número sirve para asegurar que el fichero es de tipo ELF.
4. *ph* es el iterador que recorre las regiones del fichero a mapear. *eph* es la última región que debe ser mapeada (límite del bucle for).
5. Para cada región se comprueba que sea de tipo ELF\_PROG\_LOAD, lo que significa que esa región debe ser mapeada en memoria.
  - a. Si es así, se mapea la región con la función **map\_in\_guest()**.

## Implementación de **vmlaunch** y **vmresume**

Función **asm\_vmrun()** [**vmm/vmx.c**] para arrancar una MV, es decir, copiar los registros del trapframe a los registros físicos y lanzar la MV (con **vmlaunch** si es la primera vez que se ejecuta la MV, y con **vmresume** en caso contrario). Para ello, en ensamblador:

1. Cargamos los registros físicos con el contenido del trapframe de la MV.
2. Lanzamos la MV con **vmlaunch** o **vmresume** según corresponda.

## Gestionando la terminación de una MV

Función **vmexit()** [**vmm/vmx.c**] para sacar de ejecución a una MV, se comporta de una u otra manera en función de la forma en la que la MV dejó de ejecutar. Para ello:

1. Se comprueba en la estructura VMCS (Virtual-Machine Control Structure) la razón de salida de ejecución.
  - a. Con la función **vmcs\_read32()** se accede a un campo de dicha estructura.
  - b. El campo que nos interesa es **VMCS\_32BIT\_VMEXIT\_REASON**.

## Mapa multi-arranque (e820)

Función **handle\_vmcall()** [**vmm/vmexits.c**] que es el manejador que captura la excepción producida por la instrucción **vmcall** (solicita un mapa de memoria “falso” para la MV) desde el guest. Para ello:

1. Si el trap es **VMX\_VMCALL\_MBMAP**
  - a. Ese mapa de memoria multiboot se divide en tres segmentos: low memory, IO memory, high memory
    - i. Características de la low memory
      - Size = 20 porque es el tamaño mínimo
      - Base\_addr = 0x0 porque comienza al principio
      - Length = IOPHYSMEM porque termina cuando comienza la IO memory
      - Type = MB\_TYPE\_USABLE porque puede usarse
    - ii. Características de la IO memory
      - Size = 20 porque es el tamaño mínimo
      - Base\_addr = IOPHYSMEM porque comienza tras la low memory

- Length = EXTPHYSMEM - IOPHYSMEM porque termina cuando comienza la high memory
- Type = MB\_TYPE\_RESERVED porque no es accesible
- iii. Características de la highmemory
  - Size = 20 porque es el tamaño mínimo
  - Base\_addr = EXTPHYSMEM porque comienza al principio
  - Length = glInfo->phys\_sz - EXTPHYSMEM porque termina al final de la memoria física asignada a la MV
  - Type = MB\_TYPE\_USABLE porque puede usarse
- b. Previo al mapa de memoria multiboot está una pequeña cabecera (multiboot\_info\_t)
  - i. Se indica el tamaño del mapa: tamaño total de los tres segmentos
  - ii. Se indica la dirección base del mapa: la dirección base que ya nos daban sumado al offset de la cabecera
  - iii. Se indican los flags: MB\_FLAG\_MMAP
- c. Se pide una página de memoria libre con page\_alloc()
- d. Se traslada el mapa de memoria multiboot creado anteriormente a la página obtenida del paso anterior. Se emplea la función memmove() primero para la cabecera y segundo para el mapa.
- e. Se deja constancia de la nueva página en el EPT añadiendo una nueva entrada con ept\_page\_insert().
- f. Tal y como indica el comentario de la función, se devuelve un puntero de la nueva región en el registro %rbx
- 2. Si el trap es VMX\_VMCALL\_IPCSEND
  - a. Se invoca a la syscall para llamar a SYS\_ipc\_try\_send() **[explicado más abajo]**
- 3. Si el trap es VMX\_VMCALL\_IPCRECV
  - a. Se invoca a la syscall para llamar a SYS\_ipc\_rcv() **[explicado más abajo]**

Función **handle\_cpuid()** [vmm/vmexits.c] para cargar en el TrapFrame de la MV las características de la CPU (será empleada para comprobar el soporte a *long mode*). Para ello:

1. Se toma el registro RAX del TrapFrame
2. Se invoca a cpuid(), pasándole como primer parámetro el registro anterior, para conocer las características de la CPU
3. Se actualiza el TrapFrame con la información obtenida de la invocación anterior

Función **bc\_pgfault()** [fs/bc.c] que trae un sector de disco a un bloque de memoria:

1. Obtenemos la nueva dirección de memoria donde se va a escribir el sector del disco
2. Se reserva una página de memoria
3. Se obtiene el sector de disco a traer a la memoria
  - a. Si es el SO host se utiliza la función ide\_read()
  - b. Si es el SO guest se utiliza la función host\_read()
    - i. En ambas funciones, se pasa el sector de disco que queremos leer (sec\_no), la dirección de memoria donde queremos escribirlo (new\_addr) y el tamaño del sector (nsecs)

Función **flush\_block()** [fs/bc.c] que expulsa un bloque de memoria a un sector de disco:

1. Obtenemos la dirección de memoria donde se encuentra el bloque a expulsar
2. Se comprueba
  - a. Si el bloque existe y está sucio, hay que escribirlo en disco
    - i. Si es el SO host se utiliza la función `ide_write()`
    - ii. Si es el SO guest se utiliza la función `host_write()`
      1. En ambas funciones, se pasa el sector de disco que queremos escribir (`sec_no`), la dirección de memoria del bloque (`addr`) y el tamaño del sector (`nsecs`)
  - b. Si el bloque existe y no está sucio, no hace falta escribirlo en disco
  - c. Si el bloque no existe, no tiene sentido hacer un flush

Función **`ipc_host_send()`** [**`lib/ipc.c`**] en vez de realizar la llamada de `sys_ipc_try_send()` directamente lo hace por medio de la instrucción en ensamblador `vmcall`, para que se produzca excepción y la capture `handle_vmcall()` entrando en el case `VMX_VMCALL_IPCRECV`:

1. Ejecuta `vmcall` con el parámetro `VMX_VMCALL_IPCRECV`

Función **`ipc_host_recv()`** [**`lib/ipc.c`**] en vez de realizar la llamada de `sys_ipc_recv()` directamente lo hace por medio de la instrucción en ensamblador `vmcall`, para que se produzca excepción y la capture `handle_vmcall()` entrando en el case `VMX_VMCALL_IPCSEND`:

1. Ejecuta `vmcall` con el parámetro `VMX_VMCALL_IPCSEND`

Función **`sys_ipc_try_send()`**

Función **`sys_ipc_recv()`**

## Errores solucionados

Error del `spin_unlock()` solucionando invocando `lock_kernel()` en la función `i386_init()` de `kern/init.c`