

# 4

## Apeluri de sistem (II)

19 martie 2009

- Care este un dezavantaj al transmiterii argumentelor unui apel de sistem prin registre?
- Pe un sistem Linux se realizeaza apelul `read(fd, buf, num)`. Explicati in pseudo-assembly cum se va realiza acest apel de sistem pe un sistem Linux care foloseste `sysenter`.
- Creați stack frame-ul pentru o funcție ce are un parametru și două variabile locale.

- Apeluri de sistem
- GAS, stack frame
- Accesarea spațiului utilizator în Linux
- GCC extended ASM
- Implementarea apelurilor de sistem în Windows

- UTLK: capitolul 10
- LKD: capitolul 5
- WI: capitolul 3 (System Service Dispatching)

- Problema: folosirea `int 0x80/sysenter` depinde de versiunea procesorului, versiunea kernelului și versiunea libc
- Soluția: Virtual Dynamic Shared Object:
  - kernelul decide de instrucțiune de trap să se folosească
  - Pentru genericitate, kernelul ca genera chiar el secvența de instrucțiuni ce trebuie folosită
  - Respectiva secvență va fi mapată în spațiul utilizator în zona VDSO

```
$ cat /proc/$$/maps
```

```
...
```

```
bfac5000-bfada000 rw-p bffeb000 00:00 0 [stack]
```

```
ffffe000-ffffff00 r-xp 00000000 00:00 0 [vdso]
```

```
$ dd if=/proc/self/mem of=linux-gate.so bs=4096  
skip=${0xffffe} count=1
```

```
$ objdump -d linux-gate.so --no-show-raw-insn
```

```
ffffe400 <__kernel_vsyscall>:  
ffffe400:      push    %ecx  
ffffe401:      push    %edx  
ffffe402:      push    %ebp  
ffffe403:      mov     %esp,%ebp  
ffffe405:      sysenter  
ffffe407:      nop
```

- „Apeluri de sistem” care rulează direct din user-space
- În zone VSDO kernelul mapează cod ce rulează direct din user-space și accesează direct date kernel, date mapate și ele în VDSO (dar R/O)
- Vgettimeofday, Vgetcpu (x86\_64)



- Primitive specializate: `get_user`, `put_user`, `copy_from_user`, `copy_to_user`
- Verificarea pointerului și tratarea eventualului page-fault la acces invalid se face intern

```
/* daca s-a generat un page fault datorita unui access
invalid primitivele intorc o valoare diferita de zero */
if (copy_from_user(kernel_buffer, user_buffer, size))
    return -EFAULT;
```

```

#define get_user(x,ptr) \
({ \
    int __ret_gu; \
    unsigned long __val_gu; \
    switch(sizeof (*(ptr))) { \
    case 1: \
        __get_user_x(1,__ret_gu,__val_gu,ptr); \
        break; \
    case 2: \
        __get_user_x(2,__ret_gu,__val_gu,ptr); \
        break; \
    ... \
    __ret_gu; \
})

```

```
#define __get_user_x(size,ret,x,ptr) \
__asm__ __volatile__( \
    "call __get_user_" #size \
    : "=a" (ret), "=d" (x) : "0" (ptr) \
)
```

```
ENTRY(__get_user_1)
    GET_THREAD_INFO(%edx)
    cmpl TI_addr_limit(%edx),%eax
    jae bad_get_user
1:    movzbl (%eax),%edx
    xorl %eax,%eax
    ret
ENDPROC(__get_user_1)
```

- Folosirea de cod ASM împreună cu cod C
- Se pastrează claritatea codului
- Se “optimizează de mâna”
- Se accesează resurse altfel neaccesibile din compilator (regiștri speciali, instrucțiuni speciale)

```
asm(instruction_template
    : constraints_1 output_operand_1,
      constraints_2 output_operand_2,
      ...,
      constraints_n output_operand_n
    : constraints_n+1 input_operand_n+1,
      constraints_n+2 input_operand_n+2,
      ...,
      constraints_n+m input_operand_n+m,
    : clobbered_regmem_1, clobbered_regmem_2,
      ...,
      clobbered_regmem_n);
```

```
instruction_mnemonic [ operand_asm | operand_no ], ...
```

- operand\_asm = „%%regstru”
- operand\_no = „%x” unde x este numărul operandului

```
unsigned long address;  
asm(“mov %%cr2, %0” : (address));  
printk(“Adress: 0x%x”, address);
```

- “m” = orice fel de operator memorie
- “r” = orice fel de operator registru
- “i” = un intreg imediat (intreg cu valoarea cunoscuta la momentul asamblarii)
- “0”, “1”, ..., “9” = impune sa se foloseasca acelasi operand ca cel indicat
- “=” = operandul este write-only; se foloseste pentru operanzii de output

```
void printk_cr2()  
{  
    unsigned long address;  
    asm("mov %%cr2, %0" : "=m" (address));  
    printk("0x%x", address);  
}
```



```

00000000 <printk_cr2>:
  0:   push    %ebp
  1:   mov     %esp,%ebp
  3:   sub    $0x18,%esp
  6:   mov    %cr2,%eax
  9:   mov    %eax,-0x4(%ebp)
  c:   mov    -0x4(%ebp),%eax
  f:   mov    %eax,0x4(%esp)
 13:  movl   $0x0,(%esp)
      16:  R_386_32    .rodata
 1a:  call   1b <x+0x1b>
      1b:  R_386_PC32  printk
 1f:  leave
 20:  ret

```

- Unele instrucțiuni pot modifica și alți regiștri / zone de memorie decât operanzii specificați
- Operanzii clobbered\_regmem sunt folosiți pentru a indica compilatorului care regiștri / zone de memorie sunt afectate de blocul asm
- (blocul asm și codul generat de compilator pot folosi același registru zonă de memorie)

```

#define __get_user_x(size,ret,x,ptr) \
  __asm__ __volatile__ ( \
    "call __get_user_" #size \
    : "=a" (ret), "=d" (x) : "0" (ptr) \
  )

```

Se verifică ca pointerul să fie în user-space

```

ENTRY(__get_user_1)
  GET_THREAD_INFO(%edx)
  cmpl TI_addr_limit(%edx),%eax
  jae bad_get_user
1:
  movzbl (%eax),%edx
  xorl %eax,%eax
  ret
ENDPROC(__get_user_1)

```

Instrucțiunea de copiere și adresa ei

```

bad_get_user:
    xorl %edx,%edx
    movl $-14,%eax
    ret
END(bad_get_user)

```

Adresa instrucțiunii ce  
face accesul în  
userspace din cadrul  
\_\_get\_user\_1

```

.section __ex_table,"a"
    .long 1b,bad_get_user
    .long 2b,bad_get_user
    .long 3b,bad_get_user
.previous

```

- Un vector de perechi (fault instruction address, fix-up code address)
- Generat la compilare în cadrul secțiunii `__ex_table`
- Rutina de tratare a page fault-ului va căuta adresa instrucțiunii ce a generat fault-ul în tabelă, și dacă o găsește va sări la adresa asociată

- Scriptul de link editare (arch/\*/kernel/\*.lds.S) gardează secțiunea cu simboluri de geneul `__start / __stop`

- Exemplu:

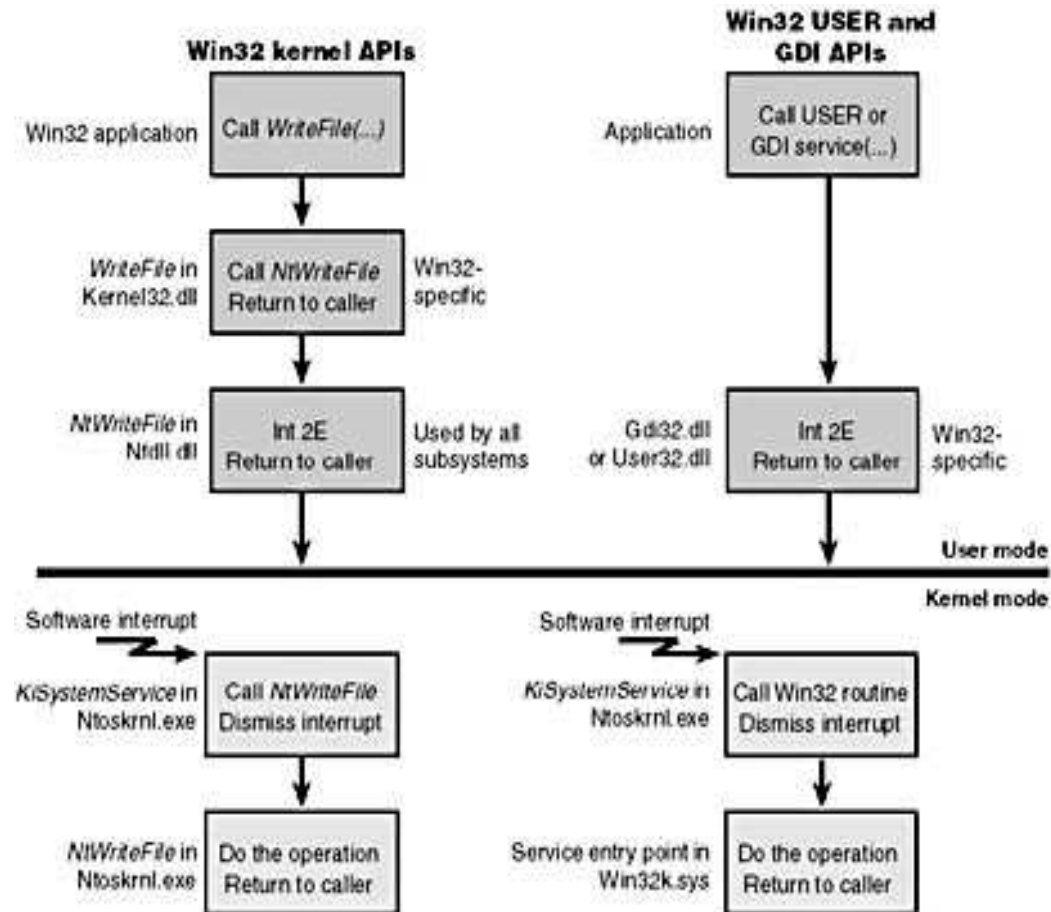
```
. = ALIGN(16); /* Exception table */
__ex_table : AT(ADDR(__ex_table) - LOAD_OFFSET) {
    __start__ex_table = .;
    *(__ex_table)
    __stop__ex_table = .;
}
```

```
int fixup_exception(struct pt_regs *regs)
{
    const struct exception_table_entry *fixup;

    fixup = search_exception_tables(regs->eip);
    if (fixup) {
        regs->eip = fixup->fixup;
        return 1;
    }

    return 0;
}
```

Setăm instrucțiunea de la care se continuă execuția după ce ieșim din handler-ul de tratare al page-fault-ului





- Fluxul de execuție:
  - **WriteFile, ReadFile** (kernel32.dll) sunt transformate în apeluri mai complexe
  - **NtWriteFile, NtReadFile** (ntdll.dll) efectuează apelul de sistem
- Numărul apelului de sistem se pune în `eax`
- Se pun pe stiva user-space parametrii
- `ebx` <- pointer către parametri; `int 0x2e`
- `edx` <- pointer către parametri; `sysenter`

ntdll!NtReadFile:

```
77f5bfa8 mov    eax,0xb7
```

```
77f5bfad mov    edx,0x7ffe0300
```

```
77f5bfb2 call  edx
```

```
77f5bfb4 ret    0x24
```

...

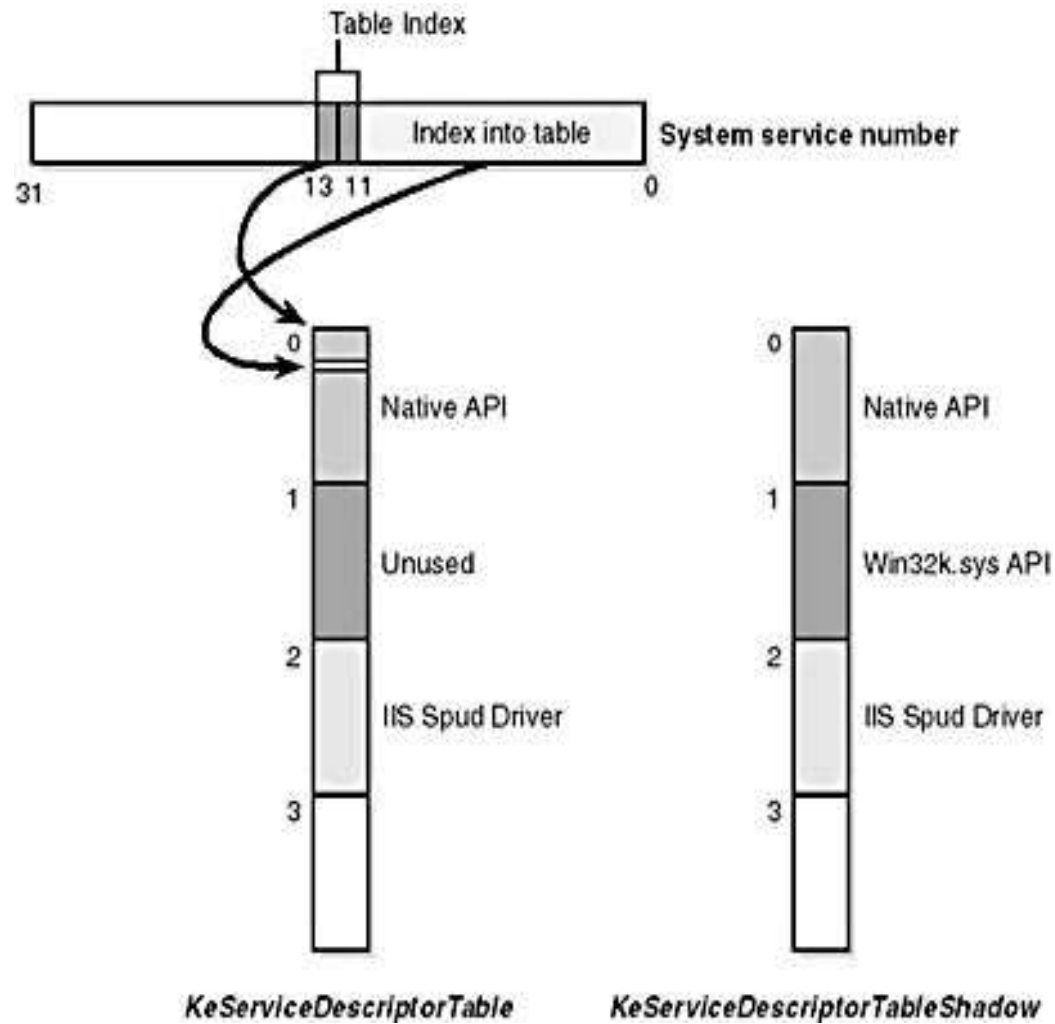
SharedUserData!SystemCallStub:

```
7ffe0300 mov    edx,esp
```

```
7ffe0302 sysenter
```

```
7ffe0304 ret
```

- Rutina de dispatch
  - identifică tabela și indexul în tabela de servicii sistem
  - determină numărul de parametri ai serviciului sistem
  - copiază parametrii pe stiva kernel
  - identifică și rulează rutina de tratare a serviciului
- Serviciile sunt grupate în tabele
  - numărul serviciului este compus din
    - 2 biți pentru tabele: 4 tabele
    - 12 biți pentru index: maxim 4096 intrări
  - fiecare proces poate avea teoretic tabele proprii



```
struct service_descriptor_table {  
    void **service_table;  
    int *counter_table;  
    int services_no;  
    unsigned char *service_parameters_table;  
} KeServiceDescriptorTable[4];
```

- Există două tabele de descriptori: cea principală și cea shadow
- Cea shadow folosește tabela unu pentru apeluri de sistem GDI (și zero pentru apelurile Win32)
- W2k3 SP1 nu mai permite înregistrarea de noi tabele de apeluri de sistem în afară de cele două

- **service\_table** – un vector de pointeri către adresele funcțiilor de tratare ale apelurilor de sistem
- **counter\_table** – un vector ce conține numărul de apeluri de sistem efectuate; este folosit doar atunci când se rulează în mod debug
- **services\_no** – numărul de elemente din `service_table` și `service_parameters_table`
- **service\_parameters\_table** – un vector ce conține dimensiunea pe stivă a parametrilor apelului de sistem

Scrieți în C & ASM secvența de cod ce poate fi folosită dintr-un modul kernel pentru afișarea numărului apelului de sistem și a codului de eroare la fiecare apel de sistem X.

```
void (*f)();

void intercept(int syscall)
{
    int syscall_table, syscall_index;

    syscall_table=syscall>>12;
    syscall_index=syscall&0x0000FFF;
    f=KeServiceDescriptorTable[syscall_table]
->st[syscall_index];
    KeServiceDescriptorTableShadow[syscall_table]
->st[syscall_index]=interceptor;
}
```



```
NTSTATUS interceptor()  
{  
    int syscall, params, syscall_table,  
        syscall_index, r;  
    void *old_stack, *new_stack;  
  
    asm mov syscall, eax  
  
    syscall_table=syscall>>12;  
    syscall_index=syscall&0x0000FFF;  
    params=KeServiceDescriptorTable[syscall_table]  
->spt[syscall_index];
```

```
mov old_stack, ebp
add old_stack, 8
sub esp, params
mov new_stack, esp
memcpy(new_stack, old_stack, params);

r=f();
DbgPrint(„%d: %d\n“, syscall, r);
return r;

}
```

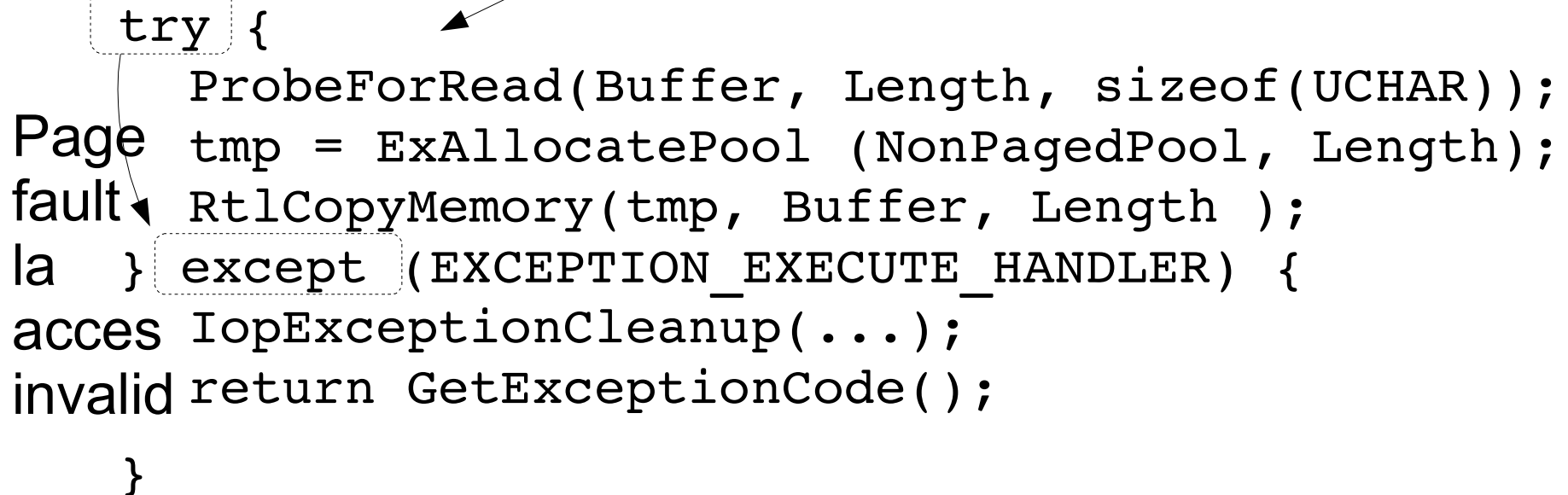
- Este doar un exercițiu didactic, în practică nu se recomandă interceptarea apelurilor de sistem
  - Probleme: race-uri, apelurile de sistem in Windows nu sunt backwards compatible
  - Folosiți filesystem filter drivers

Se verifică ca pointerul să fie în user-space

```

try {
    ProbeForRead(Buffer, Length, sizeof(UCHAR));
Page tmp = ExAllocatePool (NonPagedPool, Length);
fault RtlCopyMemory(tmp, Buffer, Length );
la } except (EXCEPTION_EXECUTE_HANDLER) {
acces IopExceptionCleanup(...);
invalid return GetExceptionCode();
}

```



```

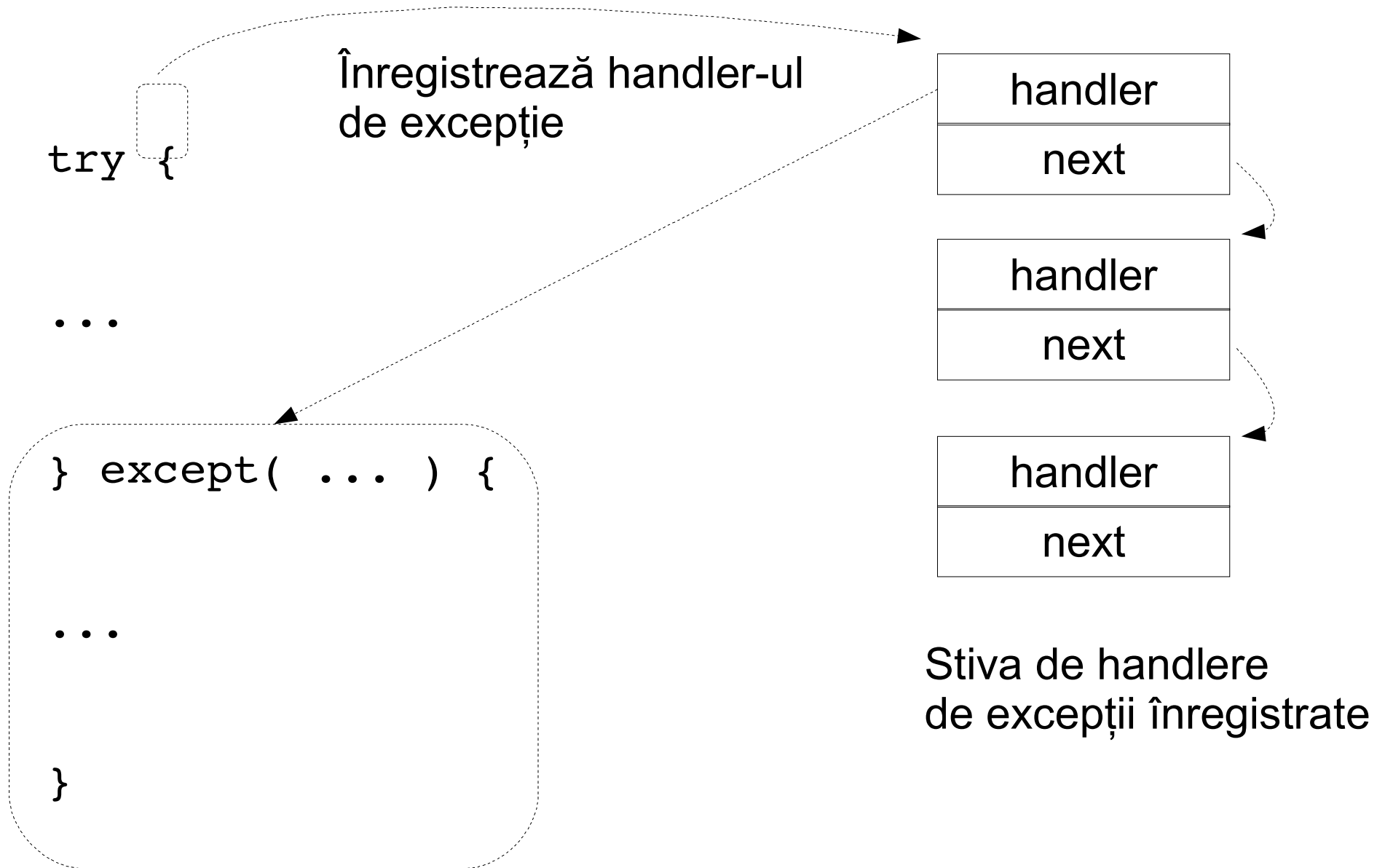
try {
    ...
} except ( ... ) {
    ...
}

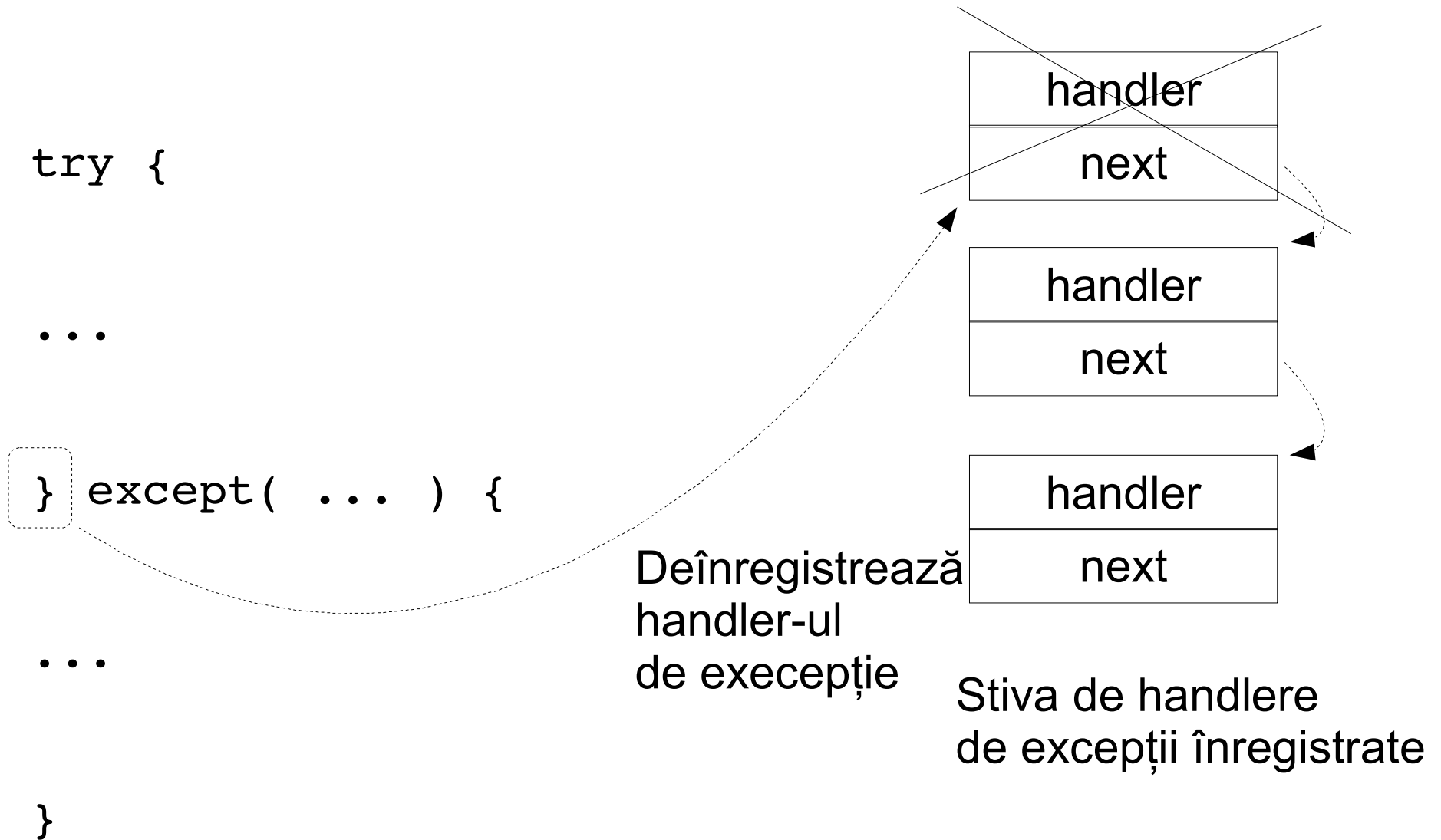
```

Dacă se generează page-fault în această secțiune ...

... și expresia din această secțiune întoarce o valoare pozitivă

... page-fault-ul este tratat în această secțiune









- Merge înapoi pe stiva de handler-e de excepții înregistrate și
  - Deînregistrează handler-ul de tratare a excepției
  - Curăță stack frame-ul asociat cu handler-ul de tratare al excepției (trebuie să ajungem să executăm rutina de tratare în stack frame-ul original)



?

