

# Tema 1 - Interceptare apeluri de sistem

Termen de predare: **joi, 17 Martie 2011, ora 23:00**

## Scopul temei

Monitorizarea apelurilor de sistem venite din partea proceselor, prin crearea unui interceptor de apeluri de sistem în kernel.

## Enunț

Să se scrie un modul de kernel **sci** care să intercepteze, la cerere, apeluri de sistem. Pentru apelurile de sistem interceptate, modulul trebuie să poată monitoriza anumite procese ce fac apeluri și să le jurnalizeze în sysklog.

La inserare, modulul va rezerva un apel de sistem nefolosit (`MY_SYSCALL_NO`) și va aștepta următoarele instrucțiuni (comenzi) din userspace prin respectivul apel:

- `REQUEST_SYSCALL_INTERCEPT` pentru a intercepta un apel de sistem
- `REQUEST_SYSCALL_RELEASE` pentru a renunța la interceptarea unui apel de sistem
- `REQUEST_START_MONITOR` pentru a porni monitorizarea unui apel de sistem (ce a fost în prealabil interceptat) pentru un anumit proces
- `REQUEST_STOP_MONITOR` pentru a opri monitorizarea unui apel de sistem (ce a fost în prealabil interceptat) pentru un anumit proces

Rutina de tratare a apelului rezervat va avea următoarea semnătură:

```
long my_syscall(int cmd, int syscall, int pid);
```

unde:

- **cmd** este comanda așteptată din userspace
- **syscall** este numărul apelului de sistem la care se referă operația
- **pid** este PID-ul procesului pentru care se dorește pornirea sau oprirea monitorizării. PID-ul este necesar doar pentru operațiile de monitorizare și va fi ignorat pentru operațiile de interceptare.

Interceptarea unui apel de sistem se referă la faptul că modulul va înlocui intrarea din tabela de apeluri de sistem, astfel încât toate apelurile ulterioare să "trecă prin modul" (și să ajungă în cele din urmă la apelul de sistem original).

Monitorizarea se referă la faptul că modulul va loga în userspace informații despre proces și apel de sistem: numărul apelului de sistem, parametrii apelului de sistem, codul de eroare întors de apelul de sistem (cel original), PID-ul procesului.

Monitorizarea se poate face pentru un proces (**pid** primește un PID) sau pentru toate procesele (**pid** primește 0).

Înainte de a efectua operațiile cerute din userspace, modulul va trebui să facă verificări de consistență, drepturi, stare și să semnalizeze, dacă este cazul, condițiile de eroare:

- verificări asupra drepturilor [eroare: permisiune/access refuzat]:
  - doar procesele ce aparțin utilizatorului privilegiat vor putea intercepta sau deintercepta un apel de sistem
  - doar procesele ce aparțin utilizatorului privilegiat pot porni sau opri monitorizarea unui apel de sistem pentru toate PID-urile
  - procesele ce aparțin de un utilizator neprivilegiat vor putea opri sau porni monitorizarea unui apel de sistem doar pentru procesele sale
- verificări de consistență [eroare: parametru nevalid]:
  - numărul apelului de sistem (se consideră nevalid apelul `MY_SYSCALL_NO` sau `__NR_exit_group`)
  - PID-ul procesului
  - un apel de sistem ce nu a fost interceptat nu poate fi deinterceptat
  - nu se poate monitoriza un apel de sistem care nu a fost interceptat
  - monitorizarea pentru un proces/apel de sistem nu poate fi oprită dacă, în prealabil, nu a fost pornită
- verificări de stare [eroare: ocupat]:
  - un apel de sistem deja interceptat nu poate fi interceptat din nou
  - monitorizarea nu poate fi pornită de două ori pentru același apel de sistem și pid

- alte verificări:
  - modulul trebuie să verifice codul de eroare al funcțiilor apelate și, în cazul în care nu poate continua, să transmită utilizatorului cauza (exemplu: nu se mai poate alocă memorie)
- cazuri speciale:
  - modulul trebuie să oprească monitorizarea pentru procesele ce s-au terminat

Se impune folosirea unei **liste pentru menținerea informațiilor despre procesele monitorizate**. Folosirea unui vector de dimensiune fixă se consideră neadecvată pentru că timpul de căutare va fi similar, complexitatea implementării va fi similară, vectorul va avea un număr limitat de intrări. Se recomandă folosirea structurilor existente în kernel.

Pentru că numărul de apeluri de sistem este în general mic (250-300) apeluri de sistem, dar și pentru o latență cât mai mică introdusă de interceptare, se recomandă ca **informațiile despre apelurile de sistem interceptate să se țină într-un vector**.

Pentru interceptarea de apeluri de sistem revedeți informațiile prezentate în [Cursul 2](#).

## Precizări Linux

- tabela cu apelurile de sistem este exportată pentru module doar în 2.4; pentru 2.6 va trebui să exportați manual această tabelă, adăugând în fișierul `arch/x86/kernel/i386_ksyms_32.c` următoarele linii:

```
extern void* sys_call_table[];
EXPORT_SYMBOL(sys_call_table);
```

De asemenea, va trebuie să modificați `arch/x86/kernel/entry_32.S` pentru a muta tabela din segmentul `.rodata` în segmentul `.data`:

```
.section .data,"a"
#include "syscall_table_32.S"
```

În plus, pentru că macro-ul `NR_syscalls` nu mai este definit în nucleu, va trebui să definiți o variabilă specială care să conțină numărul de apeluri de sistem în `arch/x86/kernel/entry_32.S`. Nucleul mașinii virtuale definește variabila `my_nr_syscalls` în acest sens.

```
ENTRY(my_nr_syscalls)
.long nr_syscalls
```

De asemenea, trebuie exportată variabila de mai sus în `arch/x86/kernel/i386_ksyms_32.c`:

```
extern long my_nr_syscalls;
EXPORT_SYMBOL(my_nr_syscalls);
```

Imaginea de pe site are făcute aceste modificări;

- pentru a folosi în modulul vostru cele două variabile exportate (`sys_call_table` și `my_nr_syscalls`) va trebui să le marcați cu `extern`:

```
extern void *sys_call_table[];
extern long my_nr_syscalls;
```

- cum kernel-ul 2.6 este preemptiv, trebuie să protejați accesul la datele comune; este indicat să folosiți spinlock-uri
- se impune să se folosească apelul de sistem 0 pentru `my_syscall`, care în prezent nu mai este folosit
- utilizatorul privilegiat va fi `root` (`uid 0`)
- pentru determinarea uid-ului unui proces folosiți câmpul `cred` din structura ce identifică procesul (`struct task_struct`)
- `task_struct`-ul procesului curent este dat de macro-ul `current`
- pentru `task_struct`-ul unui proces oarecare folosiți funcțiile `pid_task()` și `find_vpid()`
- logarea apelului de sistem se va face cu macro-ul `log_syscall` definit în `sci_lin.h` (fișierul se găsește în arhiva de teste)
- pentru eliminarea proceselor ce s-au terminat va trebui să interceptați apelul de sistem `exit_group`

- codurile de eroare ce trebuie întoarse sunt:
  - - EINVALID pentru parametru invalid
  - - EBUSY pentru ocupat
  - - EPERM pentru access refuzat
  - - ENOMEM pentru memorie insuficientă
  - 0 pentru succes

## Precizări Windows

- datorită faptului că accesul la tabelele de apeluri de sistem este dependent de versiunea/build-ul sistemului de operare, tema trebuie rezolvată pentru mașina virtuală de pe site; este foarte probabil ca funcțiile helper puse la dispoziție în `sci_win.h` (fișierul se găsește în arhiva de teste) să nu funcționeze pe alte versiuni

- tabelele de apeluri de sistem pot fi accesate din descriptorii asociați, definiți de

```
struct std KeServiceDescriptorTable[2];
struct std *KeServiceDescriptorTableShadow;
```

structura `std` este definită în `sci_win.h`

- la inițializarea modului, va trebui să chemați funcția **get\_shadow()** care va inițializa descriptorul pentru tabela shadow; asta pentru că tabela shadow nu este exportată de către kernel
- pentru că nu există apeluri de sistem nefolosite în Windows, va trebui să înlocuiți tabela de apeluri de sistem cu una nouă, în care să adăugați apelul de sistem propriu; atenție, trebuie să înlocuiți tabelele de apeluri de sistem atât în descriptorul principal cât și în shadow
- numărul apelului de sistem **my\_syscall** a fost definit în `sci_win.h` la **MY\_SYSCALL\_NO**; puteți presupune că această valoare este întotdeauna mai mare decât numărul de apeluri de sistem prezente
- trebuie tratate doar apelurile de sistem din tabela 0
- pentru determinarea utilizatorului unui process folosiți funcțiile **GetUserOf(pid)** și **GetCurrentUser**, funcții definite în `sci_win.h`; pentru a verifica dacă doi utilizatori sunt identici folosiți funcția **CheckUsers**
- utilizatorul privilegiat se consideră utilizatorul cu dreptul de încărcat / scos module din kernel; puteți verifica dacă procesul curent aparține de un utilizator privilegiat cu funcția **UserAdmin**, definită în `sci_win.h`
- logarea apelului de sistem se va face prin logarea de pachete `struct packet_log` definite în `sci_win.h`, pachete ce vor fi încapsulate în câmpul **DumpData** al unei intrări de log (vezi **IoWriteErrorLogEntry**) cu codul de eroare 0
- pentru că W2K, WXP, W2k3 sunt preemptive în kernel, trebui să folosiți spinlock-uri pentru a vă proteja de race-uri; folosiți pe cât posibil funcții **Interlocked\***
- de la XP se marchează tabela de apeluri de sistem ca read-only; folosiți macrourele **WPOFF** și **WPON** (definite în `sci_win.h`) pentru a opri, respectiv porni protecția la scriere
- pentru eliminarea proceselor ce s-au terminat din structurile folosite vedeți **PsSetCreateProcessNotifyRoutine**
- codurile de eroare ce trebuie întoarse sunt:
  - STATUS\_INVALID\_PARAMETER pentru parametru invalid
  - STATUS\_DEVICE\_BUSY pentru ocupat
  - STATUS\_ACCESS\_DENIED pentru access refuzat
  - STATUS\_NO\_MEMORY pentru memorie insuficientă
  - STATUS\_SUCCESS pentru succes

## Exemplu de cod pentru interceptarea apelurilor de sistem în Windows

```
void (*f)();
void intercept(int syscall)
{
```

```
int syscall_table, syscall_index;
syscall_table=syscall>>12;
syscall_index=syscall&&0x0000FFFF;
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&&0x0000FFFF;
    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;
}
```

## Testare

Pentru simplificarea procesului de corectare al temelor, dar și pentru a reduce greșelile temelor trimise, corectarea temelor se va face automat cu ajutorul unor teste publice ([linux md5:adc363c555b81d063fa61bdd7949d62c](https://md5.adc363c555b81d063fa61bdd7949d62c), [windows md5:dd79fd6529513fa46020087ba50c2295](https://md5.dd79fd6529513fa46020087ba50c2295)). Testele presupun că numele modului de kernel este `sci` atât pe Linux cât și pe Windows.

### Linux

Testul se folosește de utilizatorul `nobody` care ar trebui să existe pe toate sistemele Linux. Din această cauză testul trebuie să fie plasat într-un director public.

### Windows

Utilizatorul cu care veți rula testul trebuie să aibă drepturile "Act as part of the operating system", "Create a token object" și "Replace a process level token". Pentru a configura aceste drepturi: Start Administrative Tools Local Security Policy Local Policies User Rights Assignments.

### Depunctări

Depunctările generale pentru teme se găsesc pe pagina de [Indicații generale](#).

Depunctări pentru probleme constatate în implementarea temei prezente:

- -0.5 implementare incorectă a cerințelor temei
- -0.2 nu se implementează eliminarea proceselor care s-au terminat
- -1.0 nu se folosesc liste în implementare
- -0.2 lipsa sincronizării sau sincronizare incompletă / incorectă
  - - 0.1 lipsa sincronizare acces lista procese
  - - 0.1 lipsa sincronizare vector de apeluri de sistem
- -0.1 utilizare funcții la nivele IRQL necorespunzătoare
- -0.0 observații
- -0.1 pentru fiecare test picat

În cazuri excepționale (tema trece testele prin nerespectarea cerințelor) și în cazul în care tema nu trece toate testele se poate scădea mai mult decât este menționat mai sus.

## Întrebări

Pentru întrebări legate de temă puteți consulta [arhivele](#) listei de discuții sau puteți trimite un [e-mail](#) (trebuie să fiți [înregistrați](#)).

## FAQ

Q: Ce ar trebui să se întâmple în următoarele situații:

1. s-a cerut monitorizarea tuturor proceselor, și apoi se cere demonitorizarea unui proces specific.
2. s-a cerut monitorizarea tuturor proceselor, și apoi se cere monitorizarea unui proces specific.
3. s-a cerut monitorizarea unui proces specific, și apoi se cere demonitorizarea tuturor.

A: Nu se va testa asemenea condiții. Este recomandată o abordare simplă, în care pid=0 este considerat la fel ca PID -ul celorlalte procese.

( <http://cursuri.cs.pub.ro/pipermail/ps0/2009-March/003030.html> )

Q: La tema 1 pe windows, structura std este un pic ambigua. Ce inseamna fiecare camp ?

```
struct std {\n    void **st;    /* service table */\n    int *ct;      /* counter table */\n    int ls;       /* last service no */\n    unsigned char *spt; /* service parameter table */\n};
```

A:

```
void **st; /* service table */
```

Un pointer la service call table [un vector de adrese ale handlerelor pentru fiecare system call, indexat dupa id-ul system call-ului]

```
int *ct; /* counter table */
```

Pentru profiling, folosit doar in debug builds. Un pointer la un vector de countere, indexat tot dupa id-ul system call-ului, care numara de cate ori a fost apelat fiecare system call.

```
int ls; /* last service no */
```

Numarul de intrari din service table / counter table / service parameter table.

```
unsigned char *spt; /* service parameter table */
```

Numarul de bytes ocupati de argumentele fiecarui system call.

Pentru dispatcherul vostru nu veti folosi decat service table si last service no, dar cand creati noua tabela este nevoie sa duplicati si service parameter table.

Pentru mai multe informatii vezi [1], descrierea KeAddSystemServiceTable [2] si trap.asm, sursa dispatcherului de system calls pentru windows [3].

[1] <http://www.honeynet.org/node/438>

[2] <http://koala.cs.pub.ro/lxr/wrk/ntos/ke/misc.c#L729>

[3] <http://koala.cs.pub.ro/lxr/wrk/ntos/ke/i386/trap.asm>

From:  
<http://elf.cs.pub.ro/so2/wiki/> - **Sisteme de Operare 2**

Permanent link:  
<http://elf.cs.pub.ro/so2/wiki/teme/tema1>

Last update: **2011/04/09 15:25**