

Laborator 3 - Kernel API

Obiectivele laboratorului

- familiarizarea cu API-ul de bază pentru nucleul Linux și nucleul Windows
- descrierea mecanismelor de alocare a memoriei
- descrierea mecanismelor de locking

Cuvinte cheie

- contexte de execuție
- printk
- kmalloc / kfree
- list_head
- spinlock_t
- struct semaphore
- atomic_t
- IoWriteErrorLogEntry
- ExAllocatePool / ExFreePool
- ExAllocatePoolWithTag / ExFreePoolWithTag
- LIST_ENTRY
- KSPIN_LOCK
- KSEMAPHORE
- Interlocked*

Materiale ajutătoare

- [Slide-uri de suport pentru laborator](#)
- [SO2 Reference Card](#)

Noțiuni generale

În cadrul laboratorului curent se prezintă un set de concepte și funcții de bază necesare programării kernel. Este important de reținut faptul că programarea kernel diferă extrem de mult față de programarea în user-space. Kernel-ul este o entitate de sine stătătoare, care nu poate folosi bibliotecile din user-space (nici chiar `libc` în Linux sau `kernel32.dll` în Windows). Drept urmare, funcțiile uzuale utilizate în user-space (`printf`, `malloc`, `free`, `open`, `read`, `write`, `memcpy`, `strcpy` etc.) nu mai pot fi folosite. În concluzie, programarea kernel se bazează pe un API total nou și independent, ce nu are legătură cu API-ul din user-space, fie că ne referim la POSIX, Win32 sau [ANSI C](#) (funcțiile standard de bibliotecă pentru limbajul C).

Accesarea memoriei

O diferență importantă în programarea kernel este modul de accesare și alocare a memoriei. Datorită faptului că programarea kernel se face la un nivel foarte aproape de mașina fizică, există reguli importante în ceea ce privește gestiunea memoriei. În primul rând, se lucrează cu mai multe tipuri de memorie:

- memorie fizică
- memorie virtuală din spațiul de adresare kernel
- memorie virtuală din spațiul de adresare al unui proces
- memorie rezidentă ? știm sigur că paginile accesate sunt prezente în memoria fizică

Memoria virtuală din spațiul de adresare al unui proces nu poate fi considerată rezidentă datorită mecanismelor de memorie virtuală implementate de sistemul de operare: paginile pot fi în swap, sau pur și simplu pot să nu fie prezente în memoria fizică datorită mecanismului de demand paging. Memoria din spațiul de adresare kernel poate fi rezidentă sau nu. Atât segmentele de date și cod ale unui modul cât și stiva kernel a unui proces sunt rezidente (în Windows, dacă se dorește,

și acestea se pot swapa). Memoria dinamică poate fi sau nu rezidentă, în funcție de modul în care se alocă.

Atunci când se lucrează cu memorie rezidentă lucrurile sunt simple: memoria se poate accesa oricând. Dacă se lucrează însă cu memorie nerezidentă, atunci aceasta se poate accesa doar din anumite contexte. Memoria nerezidentă se poate accesa doar din context proces. Accesarea memoriei nerezidente din context întrerupere are rezultate imprevizibile, și din această cauză, atunci când sistemul de operare detectează un astfel de access, va lua măsuri drastice: blocarea sau resetarea sistemului, pentru a preveni coruperi grave.

Memoria virtuală a unui proces nu se poate accesa direct din kernel. În general este descurajată total accesarea spațiului de adresă al unui proces, dar există situații în care un device driver trebuie să o facă. Cazul tipic este cel în care device driver-ul trebuie să acceseze un buffer din user-space. În acest caz, device driverul va trebui să folosească funcții speciale, și nu să acceseze direct bufferul. Acest lucru este necesar pentru a preveni accesarea unor zone invalide de memorie.

O altă diferență față de programarea din userspace, relativ la lucrul cu memoria este datorată stivei; stiva a cărei dimensiune este fixă și limitată. În kernel-ul Linux 2.6.x se folosește implicit o stivă de 4K, iar în Windows se folosește o stivă de 12K. Din această cauză vor trebui evitate alocarea unor structuri de mari dimensiuni pe stivă sau folosirea apelurilor recursive.

Contexte de execuție

Relativ la modul de execuție în kernel, distingem două contexte: context **proces** și context **întrerupere**. Ne aflăm în context proces atunci când rulăm cod ca urmare a unui apel de sistem sau când rulăm în contextul unui kernel thread. Atunci când rulăm în rutina de tratare a unei întreruperi sau a unei acțiuni amânabile, rulăm în context întrerupere.

Unele dintre apelurile din API-ul kernel pot duce la blocarea procesului curent. Exemple comune sunt folosirea unui semafor sau așteptarea unei condiții. În acest caz, procesul este trecut în starea **WAITING** și alt proces este rulat. O situație interesantă apare în momentul în care o funcție ce poate duce la suspendarea procesului curent este chemată din context întrerupere. În acest caz, nu există un proces curent, și din această cauză rezultatele sunt imprevizibile. De câte ori sistemul de operare detectează această condiție va genera o condiție de eroare care va duce la oprirea sistemului de operare.

Locking

Una dintre cele mai importante caracteristici ale programării în kernel este paralelismul. Atât Linux cât și Windows suportă sisteme SMP, cu mai multe procesoare, dar și preemptivitate în kernel. Acest lucru face programarea kernel mai dificilă, deoarece accesul la variabilele globale trebuie sincronizat, fie cu primitive de spinlock, fie cu primitive blocante. Deși este recomandat să se folosească primitive blocante, acestea nu pot fi folosite în context întrerupere, așa că singura soluție de locking în context întrerupere sunt spinlock-urile.

Spinlock-urile sunt folosite pentru realizarea excluderii mutuale. Atunci când nu pot obține accesul la regiunea critică nu suspendă procesul curent, ci folosesc mecanismul de busy-waiting (așteaptă într-un ciclu while eliberarea lock-ului). Codul care se execută în regiunea critică protejată de un spinlock nu are voie să suspende procesul curent (trebuie să respecte condițiile execuției în context întrerupere). Mai mult, nu se va ceda procesorul decât pentru servirea întreruperilor. Datorită mecanismului folosit, este important ca un spinlock să fie deținut cât mai puțin timp posibil.

Preemptivitate

Atât Linux cât și Windows folosesc kernele preemptive. Nu trebuie confundată noțiunea de multitasking preemptiv cu noțiunea de kernel preemptiv. Noțiunea de multitasking preemptiv se referă la faptul că sistemul de operare întrerupe rularea unui proces în mod forțat, atunci când acestuia i-a expirat cuanta de timp și rulează în user-space, pentru a rula alt proces. Un kernel este preemptiv dacă un proces ce rulează în kernel-mode (ca urmare a unui apel de sistem) poate fi întrerupt pentru a rula un alt proces.

Datorită preemptivității, atunci când partajăm resurse între două porțiuni de cod ce pot rula din contexte proces diferite, trebuie să ne protejăm cu primitive de sincronizare, chiar și în cazul uni-procesor.

Linux Kernel API

Convenție indicare erori

În Linux kernel convenția folosită la apelul funcțiilor pentru a indica succes este identică cu cea din programarea UNIX: 0

pentru succes, sau o valoare diferită de 0 pentru insucces. Pentru insucces se returnează valori negative, așa cum este prezentat în exemplul de mai jos:

```
if (alloc_memory() != 0)
    return -ENOMEM;

if (user_parameter_valid() != 0)
    return -EINVAL;
```

Lista exhaustivă a erorilor și o sumară explicație găsiți în [_include/asm-generic/errno-base.h](#) și [_include/asm-generic/errno.h](#).

Șiruri de caractere

În Linux, programatorului de kernel i se pun la dispoziție funcțiile uzuale de lucru pe șiruri: strcpy, strncpy, strlcpy, strcat, strncat, strlcat, strcmp, strncmp, strnicmp, strchr, strnchr, strchr, strstr, strlen, memset, memcpy, memmove, memscan, memcmp, memchr. Aceste funcții sunt declarate în headerul [_include/linux/string.h](#) și sunt implementate în kernel.

printk

Echivalentul printf în kernel este printk, definit în [_include/linux/kernel.h](#). Sintaxa printk seamănă foarte mult cu cea a printf. Primul parametru al printk decide categoria de mesaje în care se încadrează mesajul curent:

```
#define KERN_EMERG    "<0>" /* system is unusable */
#define KERN_ALERT   "<1>" /* action must be taken immediately */
#define KERN_CRIT    "<2>" /* critical conditions */
#define KERN_ERR     "<3>" /* error conditions */
#define KERN_WARNING "<4>" /* warning conditions */
#define KERN_NOTICE  "<5>" /* normal but significant condition */
#define KERN_INFO    "<6>" /* informational */
#define KERN_DEBUG   "<7>" /* debug-level messages */
```

Astfel, un mesaj în kernel de tip warning ar fi trimis cu:

```
printk(KERN_WARNING "my_module input string %s\n", buff);
```

În cazul în care nivelul de logging lipsește din apelul printk, se realizează logging cu nivelul implicit de la momentul apelului. Un lucru ce trebuie reținut este că mesajele trimise cu printk sunt vizibile doar pe consolă ¹¹ și doar dacă nivelul lor depășește nivelul implicit setat pe consolă ²¹.

Alocare memorie

În Linux se poate alocă doar memorie **rezidentă**, cu ajutorul apelului [kmallocc](#). Un apel tipic kmalloc este prezentat în continuare:

```
#include <linux/slab.h>

if (!(string = kmalloc (string_len+1, GFP_KERNEL))) {
    //report error: -ENOMEM;
}
```

După cum se observă, primul parametru indică dimensiunea în octeți a zonei de alocat. Funcția întoarce un pointer către o zonă de memorie ce poate fi folosită direct în kernel, sau NULL dacă nu s-a putut alocă memorie. Cel de-al doilea parametru specifică modul în care se dorește să se facă alocarea, iar cele mai folosite valori sunt:

- GFP_KERNEL - folosirea acestei valori poate duce la suspendarea procesului curent; nu poate fi deci folosit în context de întrerupere;
- GFP_ATOMIC - atunci când se folosește această valoare se garantează ca funcția kmalloc nu suspendă procesul curent; poate fi folosită oricând.

Complementara funcției kmalloc este [kfree](#), funcție ce primește ca argument o zonă alocată de kmalloc. Această funcție nu suspendă procesul curent și, în consecință, poate fi apelată din orice context.

Liste

Pentru că listele înlănțuite sunt deseori folosite, Linux kernel API pune la dispoziție o modalitate unitară de definire și folosire a listelor. Aceasta implică folosirea unui element de tipul struct list_head în cadrul structurii pe care vrem să o folosim ca element al unei liste. Structura list_head este definită în [_include/linux/list.h](#) alături de toate celelalte

funcții ce lucrează pe liste. Codul următor arată definiția structurii `list_head` și folosirea unui element din acest tip într-o altă structură bine cunoscută din kernelul de Linux:

```
struct list_head {
    struct list_head *next, *prev;
};

struct task_struct {
    ...
    struct list_head children;
    ...
};
```

Rutinele uzuale pentru lucrul cu liste sunt următoarele:

- `LIST_HEAD(name)` este folosit pentru a declara santinela unei liste
- `INIT_LIST_HEAD(struct list_head *list)` se folosește pentru a inițializa un element de tipul `struct list_head` prin setarea valorii câmpurilor `next` și `prev` la `list`.
- `list_add(struct list_head *new, struct list_head *head)` adaugă elementul `new` după elementul `head`.
- `list_del(struct list_head *entry)` șterge elementul ce conține `entry` din lista din care face parte.
- `list_entry(ptr, type, member)` întoarce structura de tip `type` care conține elementul `ptr` din listă cu numele `member` în cadrul structurii.
- `list_for_each(pos, head)` iterează o listă, folosind `pos` drept cursor.
- `list_for_each_safe(pos, n, head)` iterează o listă, folosind `pos` drept cursor și `n` cursor temporar. Acest macro este folosit în cazul în care se dorește ștergerea unui element din listă.

Următorul cod arată modul de folosire al acestor rutine:

```
#include <linux/list.h>

struct pid_list {
    pid_t pid;
    struct list_head list;
};

LIST_HEAD(my_list);

static int add_pid(pid_t pid)
{
    struct pid_list *ple = kmalloc(sizeof *ple, GFP_KERNEL);

    if (!ple)
        return -ENOMEM;

    INIT_LIST_HEAD(&ple->list);
    ple->pid = pid;
    list_add(&ple->list, &my_list);

    return 0;
}

static int del_pid(pid_t pid)
{
    struct list_head *i, *tmp;
    struct pid_list *ple;

    list_for_each_safe(i, tmp, &my_list) {
        ple = list_entry(i, struct pid_list, list);
        if (ple->pid == pid) {
            list_del(i);
            kfree(ple);
            return 0;
        }
    }

    return -EINVAL;
}

static void destroy_list(void)
{
    struct list_head *i, *n;
    struct pid_list *ple;

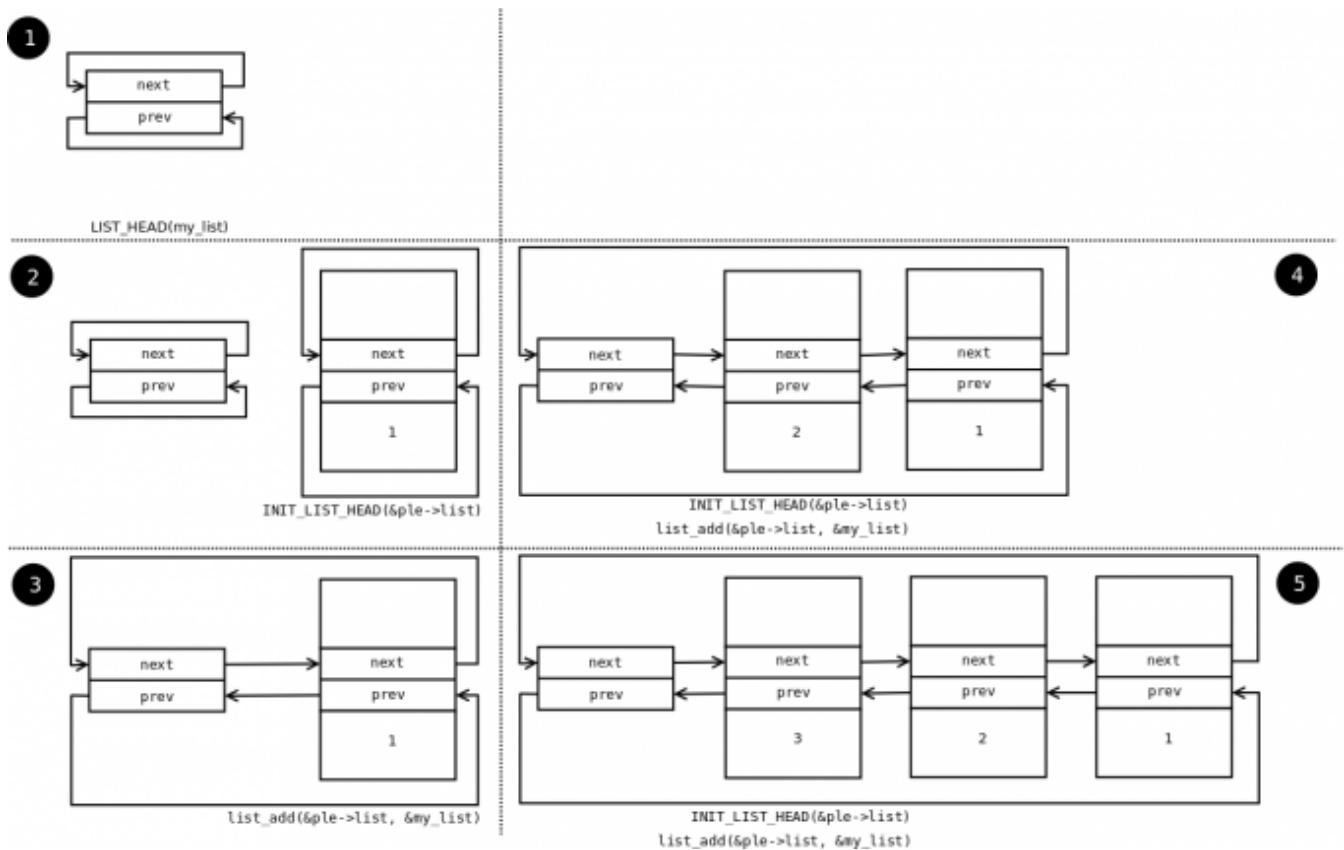
    list_for_each_safe(i, n, &my_list) {
        ple = list_entry(i, struct pid_list, list);
        list_del(i);
        kfree(ple);
    }
}
```

```

    }
}

```

Evoluția listei poate fi văzută în următoarea figură:



Se observă comportamentul de tip stivă introdus de macro-ul `list_add` precum și folosirea unei santinele.

Din exemplul de mai sus se observă că modalitatea de definire și folosire a unei liste (dublu înlănțuite) este generică și în același timp nu introduce un overhead suplimentar. Structura `list_head` este folosită pentru a menține legăturile între elementele listei. Se observă, de asemenea, că iterarea prin listă se face tot cu ajutorul acestei structuri, iar obținerea elementelor din listă se face cu ajutorul `list_entry`. Această idee de implementare și folosire a unei liste nu este nouă, ea fiind descrisă în *The Art of Computer Programming* de Donald Knuth în anii '80.

Mai multe funcții și macrodefiniții de lucru cu liste kernel sunt prezentate și explicate în headerul <include/linux/list.h>.

Locking

Spinlock-uri

`spinlock_t` (definit în <linux/spinlock.h>) este tipul de bază ce implementează conceptul de spinlock în Linux. El descrie un spinlock, iar operațiile asociate cu un spinlock sunt `spin_lock_init`, `spin_lock`, `spin_unlock`. Un exemplu de utilizare este prezentat mai jos:

```

#include <linux/spinlock.h>

DEFINE_SPINLOCK(lock1);
spinlock_t lock2;

spin_lock_init(&lock2);

spin_lock(&lock1);
/* critical region */
spin_unlock(&lock1);

spin_lock(&lock2);
/* critical region */
spin_unlock(&lock2);

```

În Linux se pot folosi spinlock-uri de tip read/write, utile în probleme de genul cititori-scriitor. Aceste tipuri de lockuri sunt

identificate de `rwlock_t`, iar funcțiile cu care se poate opera asupra unui spinlock de tip read/write sunt `rwlock_init`, `read_lock`, `write_lock`. Un exemplu de utilizare:

```
#include <linux/spinlock.h>

DEFINE_RWLOCK(lock);

struct pid_list {
    pid_t pid;
    struct list_head list;
};

int have_pid(struct list_head *lh, int pid)
{
    struct list_head *i;
    void *elem;

    read_lock(&lock);
    list_for_each(i, lh) {
        struct pid_list *pl = list_entry(i, struct pid_list, list);
        if (pl->pid == pid) {
            read_unlock(&lock);
            return 1;
        }
    }
    read_unlock(&lock);

    return 0;
}

void add_pid(struct list_head *lh, struct pid_list *pl)
{
    write_lock(&lock);
    list_add(&pl->list, lh);
    write_unlock(&lock);
}
```

Semafoare

Un semafor este reprezentat de o variabilă de tipul `struct semaphore` (definit în [linux/semaphore.h](#)). Funcțiile și macro-urile pentru lucrul cu semafoare sunt prezentate în continuare:

```
#include <linux/semaphore.h>

/* functii pentru initializarea semaforului */
void sema_init(struct semaphore *sem, int val);
DECLARE_MUTEX(name);
void init_MUTEX(struct semaphore *sem);
void init_MUTEX_LOCKED(struct semaphore *sem);

/* functii pentru achizitionarea semaforului */
void down(struct semaphore *sem);
int down_interruptible(struct semaphore *sem);
int down_trylock(struct semaphore *sem);

/* functie pentru eliberarea semaforului */
void up(struct semaphore *sem);
```

Funcția `down` decrementează valoarea semaforului și se blochează până când aceasta devine iar nenegativă. Funcția `down_interruptible` face același lucru numai că operația poate fi întreruptă. Se recomandă testarea de fiecare dată a valorii întoarse de această funcție deoarece o valoare diferită de 0 înseamnă că operația a fost întreruptă și apelantul nu a obținut semaforul. Funcția `down_trylock` este varianta neblocaută pentru achiziționarea unui semafor, dacă nu se poate lua semaforul se întoarce o valoare diferită de 0.

Trebuie reținut faptul că nu este permisă o operație `down` blocantă în context de întrerupere sau într-o regiune în care se deține un spinlock (context atomic).

După apelul funcției `up` apelantul nu mai deține semaforul.

Variabile atomice

De multe ori este nevoie doar de sincronizarea accesului la o variabilă simplă, de exemplu un contor. Pentru aceasta se poate folosi o variabilă de tip `atomic_t` (definit în [include/asm-generic/atomic.h](#)) care ține o valoare întreagă. Mai jos sunt prezentate unele operații care pot fi efectuate asupra unei variabile `atomic_t`:

```
#include <asm/atomic.h>

void atomic_set(atomic_t *v, int i);
int atomic_read(atomic_t *v);
void atomic_add(int i, atomic_t *v);
void atomic_sub(int i, atomic_t *v);
void atomic_inc(atomic_t *v);
void atomic_dec(atomic_t *v);
int atomic_inc_and_test(atomic_t *v);
int atomic_dec_and_test(atomic_t *v);
```

Operatii atomice pe biti

Kernelul pune la dispozitie un set de functii (in [include/linux/bitops.h](#)) care modifica sau testeaza biti in mod atomic.

```
#include <asm/bitops.h>

void set_bit(int nr, void *addr);
void clear_bit(int nr, void *addr);
void change_bit(int nr, void *addr);
int test_and_set_bit(int nr, void *addr);
int test_and_clear_bit(int nr, void *addr);
int test_and_change_bit(int nr, void *addr);
```

addr reprezinta adresa zonei de memorie ai carei biti se modifica sau testeaza, iar nr reprezinta asupra carui bit din aceasta zona se efectueaza operatia.

Windows Kernel API

Convenție indicare erori

În Windows kernel convenția folosită la apelul funcțiilor pentru a indica succes este următoarea: se întoarce STATUS_SUCCESS pentru succes și o valoare diferită de STATUS_SUCCESS pentru insucces. Aceste coduri pentru succes sau diferite erori sunt definite în [ntstatus.h](#). Semnificația codurilor de eroare este explicată în [NTSTATUS values](#) din MSDN.

șiruri de caractere (ASCII, Unicode)

În Windows, atunci când se lucrează cu funcții uzuale de lucru pe șiruri programatorul are două posibilități: folosește apelurile clasice ([strcpy](#), [strncpy](#), [strcat](#), etc.) sau folosește apeluri specifice programării kernel ([RtlStringCbCat](#), [RtlStringCbCopy](#), [RtlStringCbPrintf](#), etc.). Folosirea primei variante are avantajul simplității dar introduce un overhead, pentru că funcțiile folosite vor fi legate în modul. Atunci când se folosește cea de a doua variantă, în afară de faptul că dispăre overhead-ul prezentat anterior, mai avem un avantaj: se pot folosi atât șiruri ASCII cât și Unicode.

În Windows, pentru a obține efecte similare cu printf, trebuie folosită funcția [IoWriteErrorLogEntry](#), al cărei mod de folosire a fost prezentat sumar într-un exemplu în laboratorul trecut. Rolul acestei funcții este să trimită un mesaj în user-space, unde acest mesaj va fi logat de serviciul Event Viewer.

Din păcate lucrul cu această funcție (și cu alte funcții de altfel) este îngreunat de faptul că șirurile de caractere așteptate sunt Unicode și nu ASCII. Un șir unicode poate fi generat de compilator folosind prefixul L, astfel încât L"string" generează un șir Unicode, iar "string" generează un șir ASCII. Tipul de date al unui caracter unicode este WCHAR sau wchar_t, iar un șir Unicode este de tip wchar_t * sau PWSTR. Constanta de terminare a șirului este UNICODE_NULL (16 biți de zero).

Există o structură predefinită în scopul lucrului mai ușor cu șirurile Unicode: UNICODE_STRING. La fel, kernel-ul pune la dispoziție niște funcții de lucru cu aceste șiruri care să înlocuiască funcțiile din biblioteca standard C. Funcțiile de manipulare de șiruri Unicode includ: [_RtlInitUnicodeString](#), [_RtlAnsiStringToUnicodeSize](#) și [RtlAnsiStringToUnicodeString](#), [RtlIntegerToUnicodeString](#), [RtlAppendUnicodeStringToString](#) și [RtlCopyUnicodeString](#), [RtlUppcaseUnicodeString](#), [RtlCompareUnicodeString](#) și [RtlEqualUnicodeString](#).

Spre exemplu, următoarea funcție transformă un șir de caractere într-un șir unicode:

```
UNICODE_STRING* TO_UNICODE(const char *str, UNICODE_STRING *unicodeStr) {
    ANSI_STRING ansiStr;

    RtlInitAnsiString(&ansiStr, str);
    if (RtlAnsiStringToUnicodeString(unicodeStr, &ansiStr, TRUE)
        != STATUS_SUCCESS)
```

```
    return NULL;
    return unicodeStr;
}
```

Alocare memorie

Alocarea și dezalocarea memoriei în Windows kernel se face cu [ExAllocatePool](#), respectiv [ExFreePool](#). Aceste funcții alocă/dezalocă atât memorie **rezidentă** cât și **nerezidentă**, la cerere. Un exemplu de folosire a acestor funcții este prezentat mai jos:

```
#include <ntddk.h>

struct my_struct *ptr;

ptr = ExAllocatePool(NonPagedPool, sizeof(struct my_struct));
if (!ptr) {
    DbgPrint("Out of memory.\n");
    return STATUS_NO_MEMORY;
}
//...
ExFreePool(ptr);
```

După cum se observă, primul parametru specifică tipul de memorie de alocat: `NonPagedPool` pentru memorie rezidentă, `PagedPool` pentru memorie swapabilă. Pentru a ușura depanarea problemelor cauzate de leak-uri de memorie se pot folosi funcțiile [ExAllocatePoolWithTag](#) și [ExFreePoolWithTag](#) care primesc un parametru suplimentar: un întreg pe 32 de biți. În general acest întreg este scris sub forma `lgat` pentru a ușura procesul de debug (pentru exemplul dat, la debug vom observa că memoria care nu a fost dezalocată are un tag de tipul `tag1` ³¹).

```
#include <ntddk.h>

struct my_struct *ptr;

ptr = ExAllocatePoolWithTag(NonPagedPool, sizeof(struct my_struct), 'lgat');
if (!ptr) {
    DbgPrint("Out of memory.\n");
    return STATUS_NO_MEMORY;
}
//...
ExFreePoolWithTag(ptr, 'lgat');
```

Funcțiile `ExAllocatePool` și `ExFreePool` au fost declarate deprecated astfel încât se recomandă folosirea variantelor cu `Tag`.

Liste

și în Windows kernel există construcții și funcții ajutătoare pentru lucrul cu listele, urmând același model de implementare ca pe Linux. Un exemplu de folosire a acestora este prezentat în continuare:

```
#include <ntddk.h>

struct pid_list {
    SINGLE_LIST_ENTRY lh;
    HANDLE pid;
};

SINGLE_LIST_ENTRY my_list = { NULL };

static NTSTATUS add_pid(HANDLE pid)
{
    struct pid_list *ple;

    if (!(ple = ExAllocatePoolWithTag(NonPagedPool, sizeof(*ple), 'lp1t')))
        return STATUS_NO_MEMORY;
    ple->pid = pid;
    ple->lh.Next = NULL;
    PushEntryList(&my_list, &ple->lh);
    return STATUS_SUCCESS;
}

static NTSTATUS del_pid(HANDLE pid)
{
    SINGLE_LIST_ENTRY *i, *j;
    struct pid_list *ple = NULL;
    NTSTATUS ret = STATUS_INVALID_PARAMETER;

    for(j = &my_list, i = my_list.Next; i; j = i, i = i->Next) {
```



```

        /* e OK sa fie i - primul element e doar head */
        ple = CONTAINING_RECORD(i, struct pid_list, lh);
        if (ple->pid == pid) {
            PopEntryList(j);
            ret = STATUS_SUCCESS;
            break;
        }
    }

    if (ret == STATUS_SUCCESS)
        ExFreePoolWithTag(ple, 'lp1t');
    return ret;
}

static void destroy_list(void)
{
    SINGLE_LIST_ENTRY *i, *j;
    struct pid_list *ple = NULL;
    j = NULL;

    for(i = my_list.Next; i; i = j)
    {
        ple = CONTAINING_RECORD(i, struct pid_list, lh);
        j = i->Next; /* salveaza pointerul catre elementul urmator pentru parcurgere */

        PopEntryList(&my_list);
        ExFreePoolWithTag(ple, 'lp1t');
    }
}

```

După cum se observă din exemplu, lista înlănțuită este una simplu înlănțuită, și nu dublu înlănțuită. Dacă se dorește folosirea unei liste dublu înlănțuite se folosește LIST_ENTRY. În rest modalitatea de folosire și implementare a listei generice este similară, urmând aceleași model al lui Donald Knuth.

Locking

Spinlock

Un spinlock este definit de KSPIN_LOCK, iar funcțiile de lock/unlock sunt [KeInitializeSpinLock](#), [KeAcquireSpinLock](#), [KeReleaseSpinLock](#). Funcțiile KeAcquireInStackQueuedSpinLock, KeReleaseInStackQueuedSpinLock sunt variantele mai performante pentru sistemele multiprocesor. Un exemplu de folosire:

```

#include <ntddk.h>

KSPIN_LOCK lock;
KIRQL      irq;

KeInitializeSpinLock(&lock);

KeAcquireSpinLock(&lock, &irq);
/* critical region */
KeReleaseSpinLock(&lock, irq);

```

Semafoare

În Windows kernel un semafor este reprezentat de o variabilă de tipul KSEMAPHORE care se inițializează cu funcția [KeInitializeSemaphore](#). Pentru achiziționarea semaforului se folosesc funcțiile [KeWaitForSingleObject](#) și [KeWaitForMultipleObjects](#), iar pentru eliberarea semaforului se folosește funcția [KeReleaseSemaphore](#). De asemenea se poate afla valoarea unui semafor folosind funcția [KeReadStateSemaphore](#).

Un exemplu de utilizare:

```

#include <ntddk.h>

NTSTATUS status;
KSEMAPHORE sem;

KeInitializeSemaphore(&sem, 1, 1);

status = KeWaitForSingleObject(&sem, UserRequest, KernelMode, TRUE, NULL);
if (NT_SUCCESS(status)) {
    /* critical region */
    KeReleaseSemaphore(&sem, 0, 1, FALSE);
}

```

```
else {
    /* failed to acquire the semaphore */
}
```

Pentru sincronizare mai pot fi folosite si alte obiecte cum ar fi procese, threaduri, mutexuri, fastmutexuri, timere, evenimente. Detalii despre acestea se gasesc in documentatia de DDK.

Operatii atomice

In Windows kernel pentru a realiza operatii atomice se folosesc functiile Interlocked.

In continuare sunt prezentate cateva dintre acestea:

- [InterlockedCompareExchange](#) compara doua valori de 32 biti si in functie de rezultatul comparatiei schimba cu o a treia valoare de 32 biti
- [InterlockedDecrement](#) decrementeaza cu o unitate o variabila intreaga
- [InterlockedExchange](#) schimba doua valori
- [InterlockedExchangeAdd](#) aduna doua valori si intoarce suma
- [InterlockedIncrement](#) incrementeaza cu o unitate o variabila intreaga

Quiz

Pentru auto-evaluare (de preferat înainte de laborator) răspundeți la întrebările de [aici](#).

Exerciții

- Folosiți [arhiva de sarcini](#) a laboratorului.
- Pe mașina virtuală de Linux recomandăm folosirea wget pentru descărcarea arhivei.
- Mașinile virtuale pot fi accesate, respectiv, prin Multicast DNS, folosind numele spook.local (Linux) și chooch.local Windows.
 - Pentru accesarea mașinilor virtuale puteți folosi SSH (are mai mult sens pentru Linux). Autentificarea se realizează folosind chei (fără parolă):


```
$ ssh -l root spook.local
$ ssh -l Administrator chooch.local
```
 - Conturile mașinilor virtuale sunt:
 - Linux: root/student, student/student
 - Windows: Administrator/student, student/student
 - fiind vorba de kernel programming/driver development veți folosi preponderent conturile privilegiate (root respectiv Administrator)
- Fiecare mașină virtuală are un director partajat prin [SMB/CIFS](#). Un set de intrări în /etc/fstab pe sistemul gazdă permit montarea facilă a acestora. Share-urile sunt:
 - **(Linux)** directorul /home/student de pe mașina virtuală Linux se poate monta în sistemul gazdă folosind comanda:


```
$ mount /home/student/spook-share/
▪ directorul /home/student/spook-share sau icon-ul proaspăt creat pe Desktop permite accesarea share-ului
```
 - **(Windows)** directorul C:\Cygwin\home\Administrator\share se poate monta în sistemul gazdă folosind comanda:


```
$ mount /home/student/chooch-share/
▪ directorul /home/student/chooch-share sau icon-ul proaspăt creat pe Desktop permite accesarea share-ului
```
- Daca nu sunt recunoscute numele spook.local sau chooch.local este posibil ca daemon-ul [Avahi](#) sa fie oprit. Reporniti folosind comanda /etc/init.d/avahi-daemon start.
- Pe Windows, pentru a folosi Vim in prompt-ul de Windows DDK sau Visual Studio, folositi comanda vim-nox.

Linux

- Folosiți directorul lin/ din [arhiva de sarcini](#) a laboratorului.
- În cazul apariției unui oops, reporniți mașina virtuală (sau faceți revert la un snapshot realizat anterior de voi)
- Punctaj total: **7 puncte**

1. **(1 punct)** Creați un modul de kernel în care să afișați conținutul unei zone de memorie alocate cu `kmalloc`

◦ **Hints:**

- Nu uitați de `kfree`
- Puteți porni de la fișierele din directorul `lin/1-2-mem/` din arhiva de sarcini a laboratorului

- Afișați caracterele printabile cu %c, folosiți funcția (macro-ul) `isprint(char c)`
 - Citiți secțiunea [Alocare memorie](#) din laborator.
 - Alocați memorie de ordinul KB
2. (1 punct) Creați un modul de kernel în care să protejați cu spinlock-uri o instrucțiune de alocare de memorie de tip `GFP_KERNEL`.
 - Ce se întâmplă la încărcarea modului?
 - Ce se întâmplă dacă alocați memorie de tip `GFP_ATOMIC`?
 - Înlocuiți alocarea de memorie cu un apel al funcției `schedule_timeout()`
 - Ce face funcția `schedule_timeout`?
 - Funcția `schedule_timeout()` primește ca parametru valoarea pentru timeout în jiffies⁴; puteți seta timeout-ul la valoarea 0.
 - Citiți secțiunile [Alocare memorie](#) și [Spinlock-uri](#) din laborator.
 3. (2 puncte) Creați trei module kernel, după cum urmează:
 - În primul definiți, exportați și inițializați un semafor
 - **Hints:**
 - Pentru a exporta semaforul folosiți macro-ul `EXPORT_SYMBOL`
 - Semaforul trebuie inițializat cu valoarea 0
 - În cel de-al doilea importați semaforul și așteptați la semafor.
 - **Hints:**
 - Pentru a importa semaforul folosiți `extern` la declarare
 - În cel de-al treilea importați semaforul și eliberați semaforul.
 - Operațiile asupra semaforului (inițializare, așteptare, eliberare) se execută în funcțiile `init` ale modulelor.
 - Afișați mesaje sugestive înainte și după fiecare operație asupra semaforului.
 - Folosiți `lsmode` pentru a observa dependențele între module după încărcare.
 - Vizualizați fișierul `modules.order`, observați dependențele
 - **Hints:**
 - Semafoarele le definiți ca "non-pointer", și apoi le apălați cu `&`
 - La inserarea celui de-al doilea modul, acesta va rămâne blocat în execuție; pentru inserarea celui de-al treilea modul, deschideți o nouă consolă în mașina virtuală cu `Alt+F2`; puteți reveni la consola anterioară cu `Alt+F1`
 - Puteți porni de la fișierele din directorul `lin/3-sem/` din arhiva de sarcini a laboratorului; pentru compilarea celor trei submodule simultan, rulați comanda `make` în directorul `3-sem/`.
 - Citiți secțiunea [Semafoare](#) din laborator.
 4. (2 puncte) Creați un modul de kernel pentru lucrul cu liste.
 - În funcția `init` creați o listă kernel de întregi și afișați lista.
 - În funcția `exit` distrugeți lista și afișați elementele pe măsură ce sunt șterse.
 - **Hints:**
 - Puteți porni de la fișierele din directorul `lin/4-list/` din arhiva de sarcini a laboratorului.
 - Citiți secțiunea [Liste](#) din laborator.
 5. (1 punct) Creați două module kernel care comunică folosind o variabilă atomică.
 - Primul modul inițializează variabila și o exportă.
 - La încărcare, cel de-al doilea modul incrementează variabila.
 - La descărcare, primul modul raportează dacă al doilea modul a fost prezent.
 - **Hints:**
 - Puteți porni de la fișierele din directorul `lin/5-atom/` din arhiva de sarcini a laboratorului; pentru compilarea celor două submodule simultan, rulați comanda `make` în directorul `5-atom/`.
 - Citiți secțiunea [Variabile atomice](#) din laborator.

Windows

- Folosiți directorul `win/` din [arhiva de sarcini](#) a laboratorului.
 - Punctaj total: **4 puncte**
1. (1 punct) Creați un modul de kernel în care să afișați conținutul unei zone de memorie alocate cu `ExAllocatePool`.
 - **Hints:**
 - Folosiți [DebugView](#) pentru captura mesajelor afișate din kernel.
 - Nu uitați de `ExFreePool`.
 - Puteți porni de la fișierele din directorul `win/1-mem/` din arhiva de sarcini a laboratorului
 - Citiți secțiunea [Alocare memorie](#) din laborator.
 2. (1 punct) Creați un modul de kernel în care să utilizați un semafor.
 - Inițializați semaforul la valoarea 5.
 - Într-un ciclu `for` cu 5 pași așteptați la semafor și afișați simultan valoarea acestuia.
 - **Hints:**
 - Pentru a lua valoarea semaforului folosiți funcția `KeReadStateSemaphore`.
 - Puteți porni de la fișierele din directorul `win/2-sem/` din arhiva de sarcini a laboratorului

- Citiți secțiunea [Locking](#) din laborator.
- 3. (2 puncte) Creați un modul de kernel pentru lucrul cu liste.
 - În funcția `DriverEntry` creați o listă kernel de întregi și afișați lista.
 - În funcția `DriverUnload` distrugeți lista și afișați elementele pe măsură ce sunt șterse.
 - Nu folosiți variabile auxiliare pentru parcurgerea listei. Folosiți **doar** capatul listei.
 - **Hints:**
 - Puteți porni de la fișierele din directorul `win/3-list/` din arhiva de sarcini a laboratorului
 - Citiți secțiunea [Liste](#) din laborator.

Soluții

- [Soluții exerciții laborator 3](#)

Resurse utile

Linux

1. Linux Device Drivers, 3rd edition
 - [Chapter 5. Concurrency and Race Conditions](#)
 - [Chapter 8. Allocating Memory](#)
 - [Chapter 11. Data Types in the Kernel](#)
2. [The Linux Kernel API](#)
3. [Linux Kernel Linked-List Explained](#)
4. [Unreliable Guide To Locking](#)
5. [Linux Kernel Walkthrough Screencast](#)

Windows

1. The Windows 2000 Device Driver Book, Second Edition ? Chapter 5. General Development Issues
2. [Windows Driver Programming Techniques](#)
3. [Windows Synchronization Techniques](#)
4. [Windows Memory Management](#)
5. [Summary of Kernel-Mode Support Routines](#)
6. [Preventing Errors and Deadlocks While Using Spin Locks](#)
7. [Locks, Deadlocks, and Synchronization](#)
8. [Managing Hardware Priorities](#)
9. [Scheduling, Thread Context, and IRQL](#)
10. [Interlocked Variable Access](#)

¹¹ în Linux consola este terminalul virtual curent; din această cauză atunci când folosiți X Windows, aceste mesaje nu o sa apară în emulorul de terminal `xterm`, însă le puteți afișa folosind comanda `dmesg` sau accesând log-ul asociat `/var/log/syslog`

¹² Pentru mai multe detalii despre configurari pentru logging consultati [Laboratorul 2](#)

¹³ pentru testarea memoriei marcată cu tag, se pot folosi o serie de utilitare ([Testing for Errors in Accessing and Allocating Memory](#)), printre care Driver Verifier ([How to Use Driver Verifier to Troubleshoot Windows Drivers](#))

¹⁴ Pentru mai multe detalii despre jiffies consultati [Laboratorul 7](#)

From:

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

Permanent link:

<http://elf.cs.pub.ro/so2/wiki/laboratoare/lab03>

Last update: **2011/03/02 20:35**