



Univerzita Hradec Králové  
Fakulta informatiky a managementu

# Systemová volání

Mgr. Josef Jan Horálek, Ph.D. & Ing. Tomáš Svoboda, Ph.D.





- = Systemová volání = volání jádra
  - = základní komunikace aplikačních programů s jádrem
- = Tvůrce programu obvykle oddělen vrstvou standardní knihovny
  - = nepotřebuje tedy znát přesné chování jádra
  - = je však dobré je pochopit



- = Jak vypadá vztah uživatelský proces – jádro ?
- = Jedná se o dva rozdílné světy:
  - = uživatelský proces běží ve svém adresném prostoru (instrukce i data) a vykonává kód programu,
  - = někdy potřebuje službu, kterou pro něj vykoná jen jádro – systémové volání.

= Příklad: Otevření a zapsání do souboru v jazyce C

```
FILE * fopen (const char * filename, const char * mode);
```

```
int main () { //Velmi zjednodušeno, chybí jakákoliv kontrola chyb FILE *  
pFile;
```

```
    pFile = fopen("myfile.txt","w"); //Syscall volání v implementaci  
    fputs("fopen example",pFile); //Syscall volání v implementaci  
    fclose(pFile); // Syscall volání v implementaci
```

```
    return 0;
```

```
}
```



= Příklad 2:

= Zjištění a nastavení parametrů zařízení – funkce `ioctl ()` .

```
int ioctl (int fd, unsigned long request, ...);
```

= Podle konkrétní hodnoty požadavku má volání dva či tři argumenty.

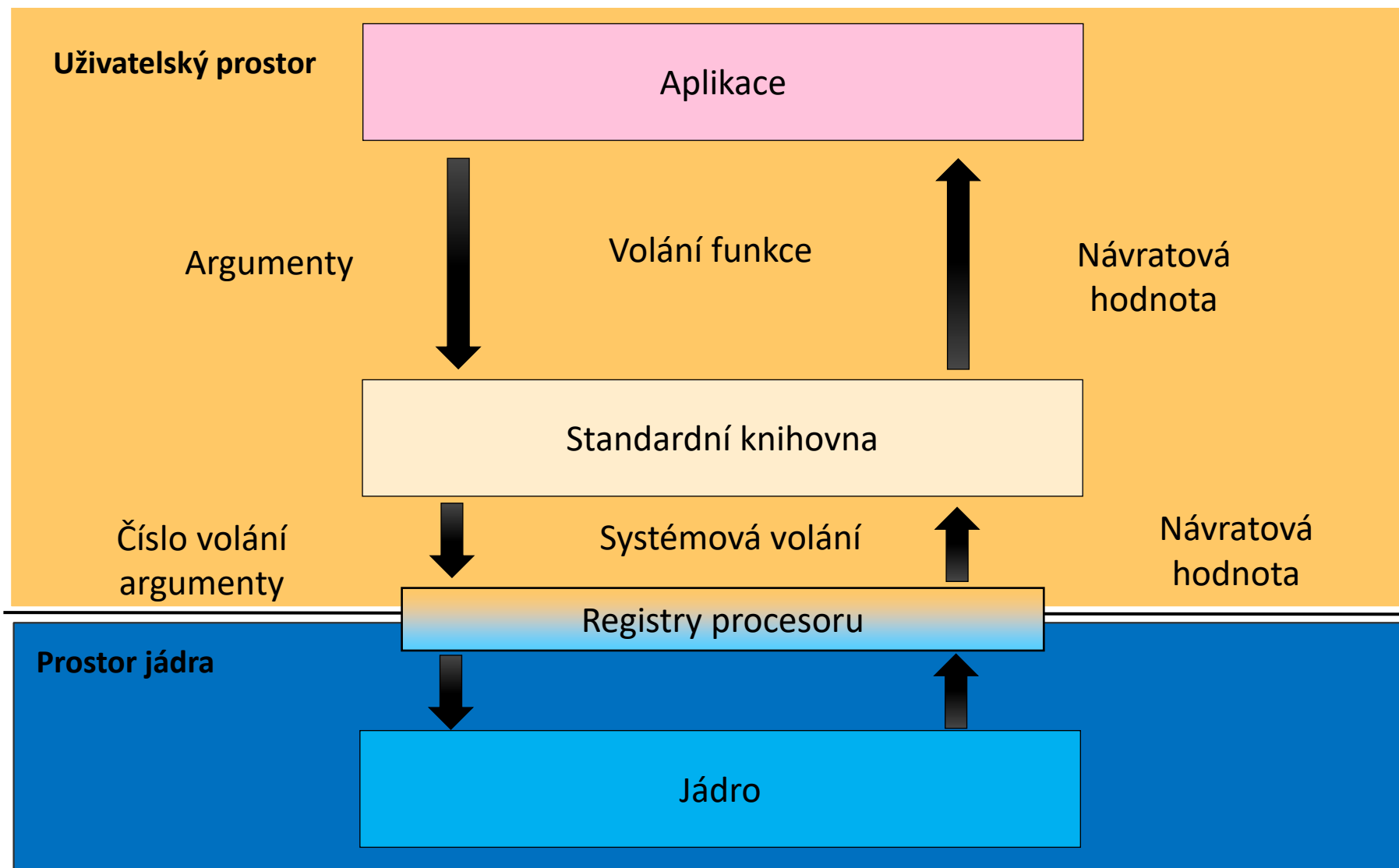
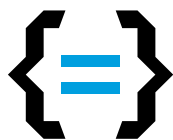
= V jádře pak vypadá prototyp funkce např. takto:

```
long (*unlocked_ioctl) (struct file *filp, unsigned int  
cmd, unsigned long arg);
```



- = Můžeme nyní zanedbat odlišné datové typy a další detaily
  - = pak se jedno volání funkce **transformuje** na jiné.
- = V čem se tedy systémové volání liší od běžné funkce?
  - = leží právě na hranici této transformace.

1. Uloží se aktuální obsah registrů
2. Do registrů se uloží číslo systémového volání a parametry
3. Předá se řízení jádru (přepnutí do režimu jádra a začne vykonávat jeho kód)
4. Podle čísla volání se vyhodnotí, jaká funkce v jádře se zavolá
5. Protože jde o volání spojené s otevřeným souborovým deskriptorem, najde se odpovídající objekt souboru
6. Podle objektu souboru se zavolá příslušná výkonná funkce v ovladači
7. Ve výkonné funkci se provede vše potřebné a na konci se vrátí výsledek
8. Do registru se uloží výsledek operace
9. Řízení se předá zpět programu resp. Standardní knihovně (uživatelský režim)
10. Výsledek se přenesení z registru do paměti
11. Obnoví se uložený stav registrů







- = Jak vlastně probíhá předání řízení jádru a opětovné předání zpět?
- = Dvě metody (platformě závislé):
  - = speciální instrukce
  - = přerušení



- = Starší metoda
- = Široké použití
  - = U jednočipů, kde neběží žádný OS, se používá stále
  - = U netypických architektur procesorů také
  - = Obecně pomalejší u mikroprocesorů
- = Volající proces způsobí softwarové přerušeni (platforma x86 nejčastěji 0x80).
- = Začne se obsluhovat, procesor přepnut do režimu jádra – provádí kód jádra.
- = Po skončení je obnoven stav procesoru – přepnut do uživatelského režimu.

- = Nevyužívá přerušení.
- = Instrukce **sysenter** sama zajistí přepnutí do režimu jádra.
- = Nastaví registry (CS, EIP atd.) a předá řízení kódu jádra.
- = Na závěr se vykoná instrukce **sysexit**, která připraví registry pro návrat a přepne do uživatelského režimu.
- = Na AMD64 se jedná o instrukce **syscall** a **sysret**

# The operating systems

Code you know about	<div data-bbox="715 454 1126 544">Ring 3 (User)</div> <div data-bbox="715 551 1126 641">Ring 0 (Linux)</div> <div data-bbox="715 648 1126 738">Ring -1 (Xen etc.)</div>	
<b>Code you don't know about</b>	<div data-bbox="715 753 1370 903">Ring -2 kernel and ½ kernel Control all CPU resources. Invisible to Ring -1, 0, 3</div> <div data-bbox="772 939 1302 1068">SMM ½ kernel. Traps to 8086 16-bit mode.</div> <div data-bbox="772 1075 1302 1189">UEFI kernel running in 64-bit paged mode.</div>	<div data-bbox="1403 753 2063 875">Ring -3 kernels</div> <div data-bbox="1403 889 2063 1189">Management Engine, ISH, IE. Higher privilege than Ring -2. Can turn on node and reimage disks invisibly. Minix 3.</div>
X86 CPU you know about	X86 CPU(s) you don't know about	

- = Systémové volání končí opuštěním jádra a předáním výsledku přes registr.
- = Konkrétní hodnoty závisí na potřebách volající aplikace, existuje však několik základní možností:
  - = úspěšně zakončené volání, které nepřenáší zpět žádnou konkrétní hodnotu, vrací nulu (0) a příslušná hodnota se ukládá do proměnné v paměti `time` (`()`)
  - = úspěšně zakončené volání, které přenáší návratovou hodnotu, vrací tuto hodnotu (obecně jakákoli hodnota obvykle nezáporná)
  - = neúspěšně zakončené volání vrací zápornou hodnotu odpovídající kódu chyby (`errno`)



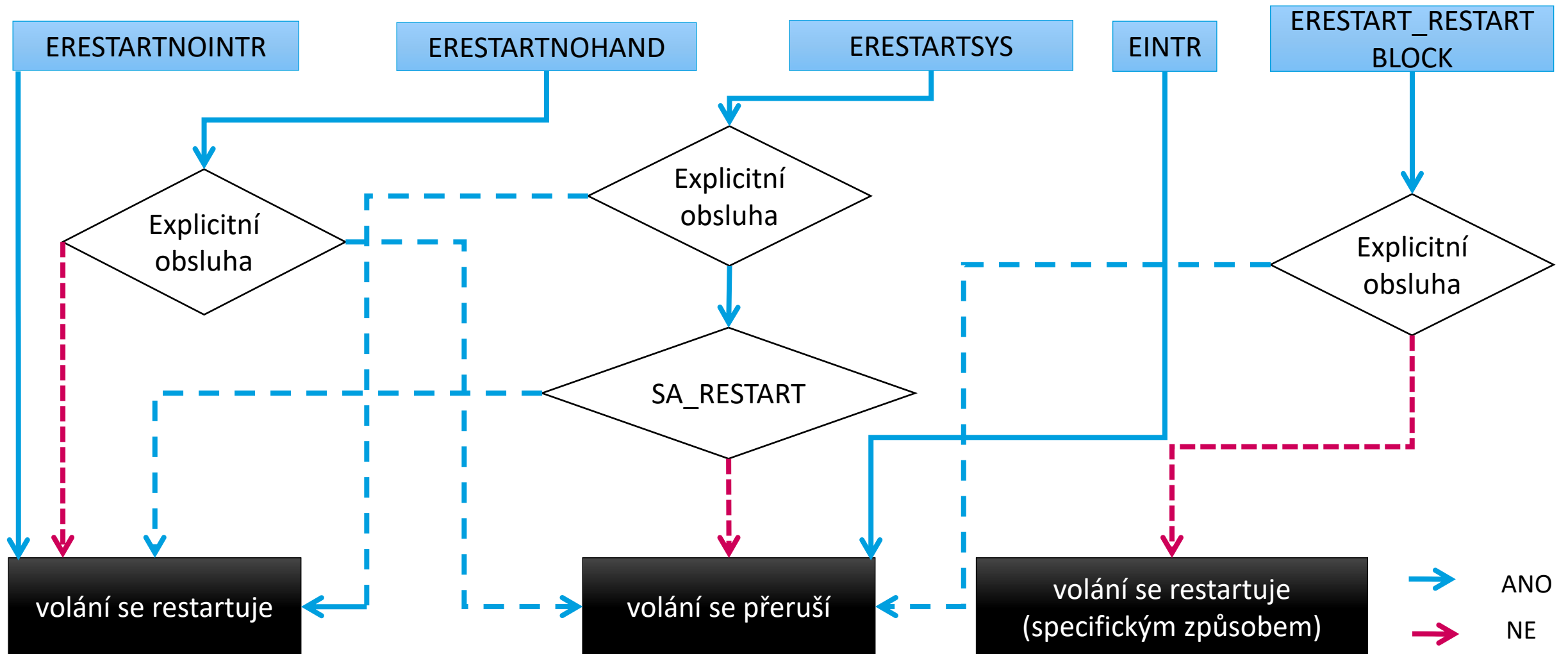
- = Často nastává situace, kdy je vykonáván kód uvnitř implementace systémového volání a v tu chvíli přijde signál, na který se musí reagovat.
- = V tuto chvíli může nastat jedna ze tří základních situací:
  - = 1) Systémové volání je buď krátké, nebo se nachází již blízko svého konce, tedy nemusí již na nic čekat. Pak se volání normálně dokončí a vrátí standardní výsledek.

- = 2) Volání bylo přerušeno na začátku, resp. v situaci, kdy buď nebyly v systému provedeny žádné změny, nebo byly provedeny bezpečně reverzibilní změny. Tato situace se řeší opuštěním jádra a vrácením hodnoty.
- = ERESTARTSYS – volání bude opět restartováno; výjimkou je případ, kdy byl signál explicitně obsloužen a neměl nastaven příznak SA\_RESTART. Tehdy nastane přerušování volání a návrat do programu
- = ERESTARTNOHAND – volání bude restartováno vždy, když není signál explicitně obsloužen uživatelskou rutinou
- = ERESTART\_RESTARTBLOCK – podobná situace s použitím jiného volání
- = ERESTARTNOINTR – volání se restartuje za každých okolností.

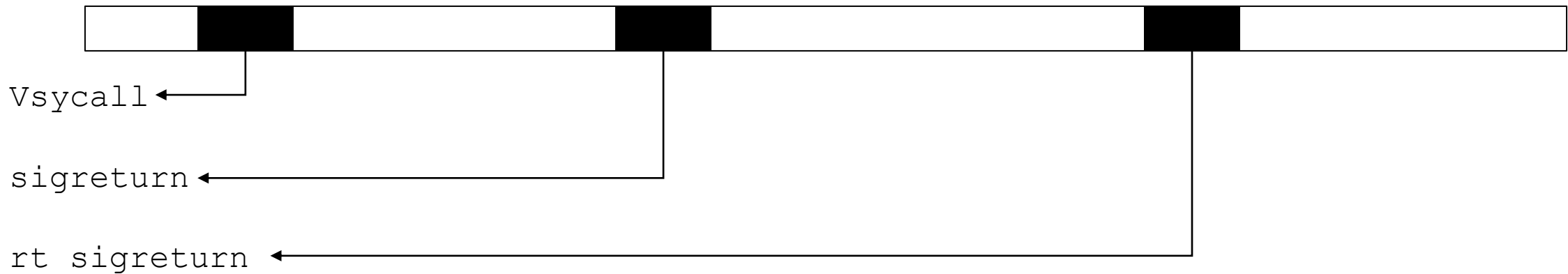
- = 3) K přerušení volání došlo v situaci, kdy již byly provedeny nevratné změny v systému, například byl do zařízení odeslán manipulační příkaz.
- = Situace se řeší opuštěním jádra s návratovou hodnotou `-EINTR`.
- = Nejdříve se samozřejmě musí systém dostat do konzistentního stavu.
- = `EINTR` – volání se vždy přeruší.

```
/* These should never be seen by user programs. To return one of ERESTART*  
 * codes, signal_pending() MUST be set. Note that ptrace can observe these  
 * at syscall exit tracing, but they will never be left for the debugged  
user  
  
 * process to see. */
```





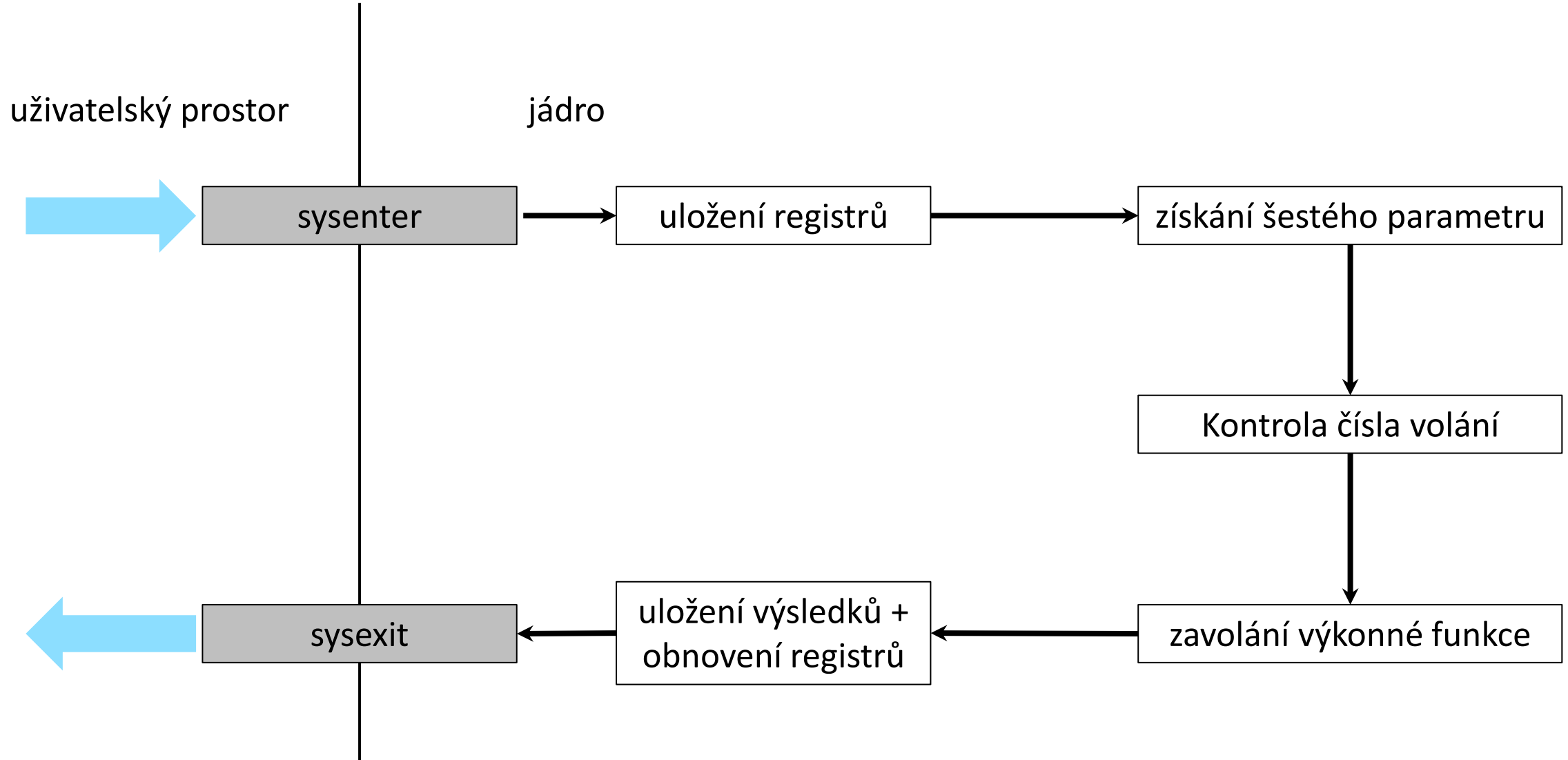
- = Knihovna – `linux-gate.so` se v systému nevyskytuje
- = Programy ji přesto vyžadují
- = Místo ní je na dané adrese **vsdo** (Virtual Dynamic Shared Object)
- = Obsahuje symboly pro virtuální systémová volání
- = Dva symboly pro návrat z obsluhy signálu



- = Souvislost s restartem systémových volání po jejich přerušení a šestým parametrem volání
  - = Procesory x86 poskytují málo registrů proto:
    - = Parametr ukládáme na zásobník (v user mode)
    - = Pro následné získání parametru využíváme registr **EBP**, tam před jeho použitím ukládáme obsah registru **ESP** (ukazatel zásobníku)
    - = Při využití speciální instrukce pro přepínání režimů dochází k přepsání **EBP**
      - = Nutno zajistit obnovení novým zkopírováním **ESP**
    - = Virtuální knihovna tak vytváří vstupní bod do jádra
  - = Umožňuje využívat tzv. virtuálních volání – např. pro pouhé čtení hodnoty

- = Případy průchodu volání jádra
  - = Běžné systémové volání
    - = uživatelská aplikace využila určité systémové volání
- = Návrat z obsluhy signálu
  - = návrat z obslužné rutiny signálu

- = Pro kód, který implementuje systémová volání je charakteristické:
  - = Společný
  - = Nízkoúrovňový
  - = Napsaný v assembleru
    - = Pro zajištění kontroly jaké strojové instrukce se dostávají do jádra
  - = Závislý na hardwarové architektuře



- = Z pohledu jádra volání začíná ve vstupním bodě, kam se dostává přes instrukci `sysenter`.
- = Používá se pro něj označení `sysenter_entry`.
- = Provedením instrukce `sysenter` se:
  - = adresa bodu `sysenter_entry` objeví v registru EIP (ukazatel instrukce)
  - = adresa zásobníku jádra v registru EXP (ukazatel zásobníku)

- = Od vstupního bodu začíná obecný kód
- = Ukládání registrů a návratové hodnoty na zásobník
  - = Registr EBP obsahuje původní obsah registru ESP (zkopírován před vstupem do volání)
  - = Přenos šestého parametru volání z uživatelského prostoru
    - = Přenos z uživatelského zásobníku prostřednictvím registru EBP (použijeme adresu uloženou v registru a následně je vložena samotná hodnota parametru)
- = Kontrola čísla volání (registr FAX) s ohledem na podporovaný rozsah
- = Vyhledání a zavolání funkce odpovídající číslu volání z tabulky `sys_call_table`





## Vykonávání systémového volání - výběr výkonné funkce

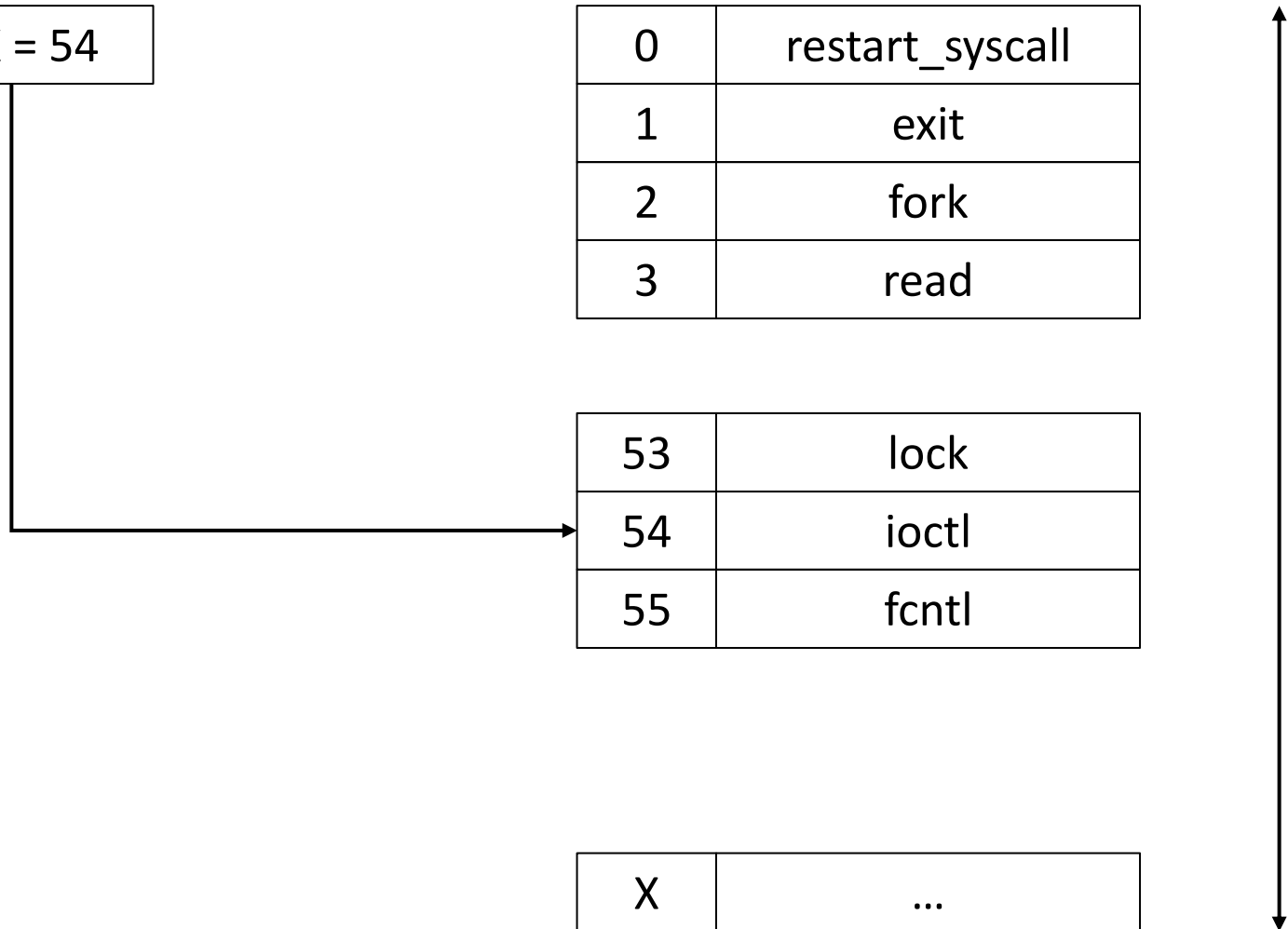
EAX = 54

0	restart_syscall
1	exit
2	fork
3	read

53	lock
54	ioctl
55	fcntl

X	...
---	-----

NR\_syscall



- = Po návratu z funkce se:
  - = Obsah registru EAX zkopíruje do zásobníku
  - = Do registru EDX je umístěna návratová adresa
    - = Při návratu do uživatelského prostoru se přesune do EIP
  - = Do registru ECX je umístěn ukazatel na uživatelský zásobník ESP
- = Po celou dobu manipulaci s registry jsou zakázána přerušování
- = Kopírování do cílových registrů a přepnutí do uživatelského režimu řeší `sys_exit`.
- = Konec systémového volání

- = Těsně po návratu z výkonné funkce volání se kontroluje, zde před návratem do uživatelského prostoru není nutné realizovat další činnosti, které typicky jsou:
  - = Přepnutí kontextu – naplánování jiné úlohy
  - = Obsluha čekajících signálů
    - = Po jejich zpracování je realizována test na potřebu realizace dalších činností.
- = Není-li již nutné řešit další úkoly dochází k obnovení registrů ze zásobníku a za využití `sys_exit` se vlákno vrací do uživatelského prostoru.



Univerzita Hradec Králové  
Fakulta informatiky a managementu

**Děkuji za pozornost**

Další téma: Signály

