

Tema:

- 1) Kapitola 4 knihy o xv6: "Traps and device drivers"
hw si vynucuje pozornost,
vykonavanie sw sa musi prerusit!
- 2) Kapitola 5 knihy o xv6: "Locking"
zamky

Preci si hw vyzaduje pozornost?

- 1) MMU nedokaze prelozit VA --> PA (napr. PTE_V bit nie je nastaveny)
- 2) uzivatelsky program deli nulou
- 3) uzivatelsky program sa pokusa vykonavat kod jadra
- 4) sietova karta prijala udaje, ktore treba spracovat
- 5) casovac chce zvyisit pocitadlo hodin
- 6) je potrebne vyriesit komunikaciu medzi procesormi (napr. pri vyprazdneni TLB)

Rozlisujeme 3 druhy vynimocnych stavov:

- 1) vynimky (vypadok stranky, delenie nulou...)
- 2) systemove volania (ecall)
- 3) prerusenania (hw zariadenia)

Terminologia v literature nie je ohladom tychto stavov jednotna. Napriklad co sa tyka RISC-V, vsetky 3 druhy sa obsluhuju uplne rovnakou technikou, a preto sa aj casto rovnako oznacuju, napr. tzv. traps.

Pre blizsie ozrejmenie sa zameriame na vstup z klavesnice. Ako funguje take napisanie
\$ ls

- 1) klavesnica posle udaje hardveru UART (universal asynchronous receiver/transmitter)
- 2) UART vygeneruje prerusenie
- 3) na prerusenie reaguje kod v jadre (ovladac), ktory precita udaje z UART (1 znak)
- 4) pomocou dalsieho ovladaca sa tento znak vypise na konzole

Ako kernel rozpozna, ktore zariadenie vyslalo signal prerusenania?

- * kazde zariadenie ma vlastne jedinecne IRQ (interrupt request) cislo
- * toto cislo je definovane vyrobcom hardveru
- * napr. UART0 ma cislo 10 v qemu (vid kernel/memlayout.h)

Registre RISC-V suvisiace s preruseniami>

- * sie -- supervisor interrupt enabled register
-- 1 bit pre zapnutie prerusenania
- * sstatus -- supervisor status register
-- bity rozlisujuce zdroj prerusenania: softverove, externe, od casovaca
- * scause -- supervisor cause register
-- dovod (pricina) prerusenania
- * stvec -- supervisor trap vector register
-- adresa obsluznej funkcie vynimocneho stavu

Ako funguje proces prerusenani v xv6

```
main():
  consoleinit():
    uartinit()
  plicinit()
  scheduler():
    intr_on():
      w_sie(r_sie() | SIE_SEIE | SIE_STIE | SIE_SSIE)
      w_sstatus(r_sstatus() | SSTATUS_SIE)
```

```
$
shell je v systemovom volani read; caka na vstup z konzoly; sled volani:
usertrap() pre systemove volanie:
  w_stvec((uint64)kernelvec)
  consoleread():
    sleep():
      scheduler():
        intr_on()
pozor! prerusenania su zapnute!
```

```
$ 1
pouzivatel stlacil klavesu 'l', ktora sposobi prerusenie na zbernici UART
subsystem PLIC doruci prerusenie na procesor
procesor vykona nasledovne kroky:
1) ak ide o prerusenie zo zariadenia a SIE bit je nulovy, nasledovne kroky sa NEkonaju
2) !!! vypne prerusenie vynulovanim SIE bitu !!! (preco?)
3) skopiruje pc do sepc
4) ulozi aktualny rezim cpu (user alebo supervisor) do SPP bitu v registri sstatus
5) ulozi pricinu prerusenania do scause registra
6) zmeni mod cpu na supervisor
7) skopiruje stvec do pc (v pripade xv6 je v stvec ulozena adresa funkcie kernelvec)
8) pokracuje vykonavanim instrukcii podla pc registra
```

co robi kernelvec()?

- 1) alokuje miesto na aktualnom zasobniku; o ktory zasobnik ide?
ulozi registre cpu do tohto miesta na zasobniku
v pripade xv6 sa jedna o zasobnik planovaca
- 2) kerneltrap()
 devintr()
 uartintr()
 c = uartgetc()
 consoleintr(c)
 handle ctrl characters
 echo c
 put c in buffer
 wakeup reader
 return from devintr
 return from kerneltrap
- 3) nahraj hodnoty registrov zo zasobnika
- 4) sret

ako sa spracova prerusenie z uzivatelskeho priestoru? napr. sekvencia: `usertests&\nls\n`
`usertrap()`

```
...
w_stvec((uint64)kernelvec)
...
if(r_scause() == 8){
  ...
} else if ((which_dev = devintr()) != 0){
  // ok
}
...
usertrapret()
```


hash tabulka, viacvlakovy pristup
paralelne operacie put() a get()
kolizie riesene obycajnym retazenim

```
struct entry{  
    int key, value;  
    struct entry *next;  
}
```

predpokladajme put(5) a put(10) sucasne
obe vlakna sa pokusia o zapis napr. do prvku table[0], ale v akom poradii?
race condition (subeh):
- ked iste poradie operacii moze sposobit nekorektne spravanie
- sucasny pristup (aspon jeden z pristupov je zapis)

kam umiestnit zamok?

- 1) na celu tabulku? --- tzv. big lock
- 2) na polozku tabulky?
- 3) v kazdej strukture 'entry' na ochranu polozky 'next'?

co je to vobec zamok (lock)?

je to objekt (abstraktny datovy typ - ADT)
pomocou roznych objektov je mozne riadit pristup k roznyim udajom
pouzitie:

```
lock l;  
acquire(&l);  
// critical section  
// i.e. x = x + 1;  
release(&l);
```

viacero vlakien moze sucasne vyvolat funkciu acquire(&l):
iba jedinemu vlaknu sa podari okamzite funkciu acquire dokoncit
ostatne vlakna musia cakat, pokym nebude vyvolana funkcia release(&l)

na co je dobry zamok:

- 1) aby vlakna/procesy po jednom pristupovali k zdielanej pamati
- 2) aby sa z neatomickych operacii stala atomicka!

1) kam umiestnit zamok (z pohladu riadenia pristupu k udajom)? podla nasledovneho pravidla:

- 1) mozu sucasne dve alebo viac vlakien pristupovat k nejakej pamati?
- 2) je aspon jeden z tychto pristupov zapis?
--> ak ano, NUTNE POUZIT ZAMOK

2) kam umiestnit zamok (z pohladu riadenia vykonavania kodu)?

uviaznutie (deadlock):

```
proces A:                proces B:  
    lock(&lock1)          lock(&lock2)  
    ...  
    lock(&lock2)          ...    lock(&lock1)  
    ...  
    ...
```

ani jeden z procesov nemoze pokracovat!

riesenie:

napr. zaviesť vzdy rovnake poradie ziskavania zamkov
to vyzaduje nasledovne operacie: predikcia, utriedenie, ziskanie zamkov
prilis komplexne!

kam umiestniť zámok (z pohľadu riadenia vykonávania kódu)? príklad dobrej stratégie:

- 1) napísať modul (nezavislú časť programu) korektné vzhľadom na seriové vykonávanie t.j. vzhľadom na vykonávanie jedným jadrom
- 2) POTOM pridať zámky na VYNUTENIE serializácie tam, kde môže prichádzať k subbehom každú takto uzamknutú časť kódu môže v jednom čase vykonávať iba jedno CPU, takže túto časť môžeme považovať za atomicky vykonávanú (neprerušiteľná operácia)

prečo používať túto stratégiu?

ľahšie sa nám uvažuje nad serióvo vykonávaným kódom...

aký je vzťah medzi zámkami a paralelizmom?

- zámky ZABRAŇUJU paralelnému vykonávaniu!!!
- keď chceme zvýšiť paralelizmus, zväčša musíme rozdeliť údaje, zaviesť viac zámkov (pre každú sekciu údajov), t.j. zjemniť granularitu uzamykania
- zvolit správny pomer medzi rozdelením údajov a zámkami patrí medzi najväčšie výzvy tvorby a návrhu programov:
 - a) celá tabuľka; každý riadok zvlášť; každá položka
 - b) celý suborový systém; každý adresár/subor zvlášť; každý blok disku zvlášť
 - c) celé jadro OS; každý subsystem zvlášť; každý objekt zvlášť
- niekedy je nutné pozmeniť značnú časť návrhu, aby bolo možné zaviesť paralelizmus! toto už môže byť veeeelmi náročná úloha!
príklad: ak vlákna čakajú na získanie voľného bloku pamäte z jedného zoznamu, rozdeliť jeden zoznam na zoznamy voľných blokov pamäte zvlášť pre každé CPU

voľba granularity uzamykania:

treba začať s tzv. veľkým uzamykaním (big locks): jeden zámok na jeden modul:

mensia sanca uviaznutia, pretože bude mensia sanca, že proces potrebuje dva zámky súčasne, keď rieši nejakú samostatnú úlohu (a celé riešenie úlohy sa nachádza v jednom module)

potom je potrebné merať efektívnosť riešenia a pozorovať správanie systému

v mnohých prípadoch staci granularita 'big lock', zväčša sa zníži efektívnosť kvôli čakaniu na jemnejšiu granularitu sa prechádza iba vtedy, keď efektívnosť riešenia s veľkým zámkom nestaci!

implementácia zámku:

```
struct lock { int locked; }
acquire(lock *l){
    while(1){
        if(l->locked == 0){ // riadok A
            l->locked = 1;    // riadok B
            return;
        }
    }
}
```

chybaaaa!!! subeh (race condition) medzi riadkami A a B!

je možné nejakým spôsobom vykonať oba riadky atomicky?

softverovo veeeelmi ťažko, vyžaduje sa nutná podpora hw; instrukcia exchange:

```
mov $1, %eax    // daj do registra eax hodnotu 1
xchg %eax, addr // atomicky vymen obsah registra eax s hodnotou na adrese addr
co vykona hardver pri instrukcii xchg?
```

- 1) lock addr globalne (t.j. pre všetky CPU a jadra uzamkne prístup k addr)
- 2) temp <-- *addr
- 3) *addr <-- %eax
- 4) %eax <-- temp
- 5) unlock addr

takto je to riešene na x86 (hw uzamknutím pamätovej oblasti)

rozny hw to rieši roznyim sposobom

!!! implementacia zamku s instrukciou xchg: !!!

```
struct lock { int locked; }
acquire(lock *l){
    while(1){
        if(xchg(&l->locked, 1) == 0)
            break;
    }
}
```

analiza kodu:

- 1) ak je uz hodnota l->locked 1, xchg ju (znovu) nastavi na 1, vrati 1, a pokracuje dalsim cyklom (vymeni hodnotu, otestuje navrat atd)
- 2) ak je pred vyvolanim instrukcie xchg hodnota l->locked 0, aspon jedno CPU, ktore vyvolava tuto instrukciu nad zamkom 'uvidi' hodnotu 0; nastavi hodnotu 1, instrukcia xchg vrati hodnotu 0, takze test bude splneny a cyklus while(1) sa prerusi

spin lock:

!!! typ zamku, ktory neustale (v cykle) vytazuje testom hodnoty CPU !!!
zbytočne vytazuje CPU? preco neprepnut vykonavanie ineho procesu/vlakna, ktore nemusí čakať?
implementacia tohto typu zamkov je jednoduchšia ako zamkov s blokovanim, menej rezie
pravidla pre spravne pouzivanie zamkov typu spin:

- 1) drzanie zamku iba na veeeeelmi kratky cas!
- 2) v ziadnom pripade sa nikdy pri drzani zamku nevzdat CPU volanim yield!

'klasicky' blokujuci zamok:

cakajuce vlakna sa vzdavaju procesora a zaraduju sa do fronty
rezie je oveľa vyššia ako pri type 'spin' (zaradovanie do fronty, planovac, prepnutie ctx)

zaverecne odporucania:

- 1) nikdy nezdielajte medzi vlakna/procesy nic, čo nemusíte
- 2) odladte aplikáciu pre beh na jednom jadre
- 3) identifikujte kritické oblasti: neatomicke operacie, subehy (operácii a/alebo pamate)
- 4) vyvoj zacnite s veľkými zamkami; avšak iba na miestach, ktoré vyžadujú vynutenie

serializacie

- 5) jemnejšiu granularitu zamkov použite iba v krajnej núdzi (kvôli zvýšeniu paralelného výkonu)
- 6) používajte automatizované nastroje na hľadanie subehov a uviaznutí