

Tema 3 - Software RAID

Termen de predare: **joi, 21 Aprilie 2011, ora 23:00**

Scopul temei

Implementarea unui modul de Software RAID ce folosește un dispozitiv logic de tip bloc ce va citi și scrie date de pe două dispozitive fizice, asigurând consistența și sincronizarea datelor de pe cele două dispozitive fizice. Tipul de RAID implementat va fi asemănător unui RAID1.

Obiective didactice

- Înțelegerea aprofundată a modului de funcționare a subsistemului de I/O.
- Obținerea de deprinderi avansate de lucru cu structuri bio și IRP.
- Acomodarea cu modul de lucru cu dispozitivele de tip bloc/disc în nucleu.
- Dobândirea de abilități de parcurgere și înțelegere a codului și API-ului dedicat subsistemului de I/O în Linux.
- Acomodarea cu API-ul de lucru cu dispozitive și IRP-uri în Windows.

Enunț

Să se scrie un modul de kernel care să implementeze funcționalitatea de RAID software. [Software RAID](#) oferă o abstractizare între dispozitivul logic și dispozitivele fizice. Implementarea va utiliza schema [RAID1](#).

Mașina virtuală dispune de două hard disk-uri suplimentare ("Hard Disk 2(IDE)" și "Hard Disk 3(IDE)") care vor reprezenta dispozitivele fizice. Sistemul de operare va oferi un dispozitiv (de tip bloc) logic care va interfața accesul din user space. Al doilea hard disk ("Hard Disk 3(IDE)") va fi folosit ca backup pentru primul "Hard Disk 2(IDE)", pentru recuperarea din eroare. Cererile de scriere în dispozitivul logic vor rezulta în două scrieri, câte una pentru fiecare hard disk. Hard disk-urile **nu** sunt partiționate. Se va considera că fiecare hard disk are o singură partiție ce acoperă întregul disk.

Fiecare partiție va stoca un sector împreună cu o sumă de control asociată ([CRC32](#)) pentru a asigura recuperarea din eroare. Dacă în cazul unei citiri un sector al primei partiții deține date corupte (valoarea CRC este greșită), se va citi sectorul de pe cea de-a doua partiție; în același timp se va corecta sectorul primei partiții. Similar în cazul unei citiri a unui sector corupt de pe a doua partiție. În cazul în care ambele sectoare au valori CRC greșite, se va returna un cod de eroare corespunzător.

Precizări generale

Pentru asigurarea recuperării din eroare se asociază fiecărui sector un cod CRC. Codurile CRC sunt stocate după LOGICAL_DISK_SIZE octeți ai partiției (macro definit în fișierele header pentru [Linux](#) și [Windows](#)). Structura pe disc va avea următoarea formă:

```
+-----+-----+-----+-----+-----+
| sector1 | sector2 | sector3 | .....|C1 |C2 |C3 |.....
+-----+-----+-----+-----+-----+
```

C1, C2, C3 sunt valorile CRC pentru sectoarele sector1, sector2, sector3. Zona CRC-urilor se regăsește imediat după LOGICAL_DISK_SIZE octeți ai partiției.

Precizări Linux

- dispozitivul logic va fi accesat ca un dispozitiv de tip bloc cu majorul SSR_MAJOR și minorul SSR_FIRST_MINOR sub numele /dev/ssr (prin macro-ul LOGICAL_DISK_NAME); macro-urile sunt definite în [ssr.h](#);
- cele două discuri sunt reprezentate de dispozitivele /dev/sdb, respectiv /dev/sdc, definite prin intermediul macro-urilor PHYSICAL_DISK1_NAME, respectiv PHYSICAL_DISK2_NAME, definite în [ssr.h](#);
- pentru lucrul cu structura [struct block device](#) asociată unui dispozitiv fizic, puteți utiliza funcțiile [open_bdev_exclusive](#) și [close_bdev_exclusive](#);
- pentru tratarea cererilor din user-space, se recomandă să nu vă folosiți de o coadă de cereri, ci să faceți prelucrare la nivel de bio; puteți înregistra o rutină corespunzătoare folosind apelul [blk_queue_make_request](#);

- pentru a transmite o cerere către un dispozitiv de tip bloc, puteți utiliza funcția [submit_bio](#);
- o singură funcție de prelucrare a cererilor pentru dispozitive de tip bloc poate fi activă la un moment dat în cadrul unei stive de apeluri (mai multe detalii [aici](#)); va trebui să submitați cererile pentru dispozitivele fizice dintr-un kernel thread; se recomandă folosirea [workqueues](#);
- pentru a transmite informațiile dorite în user-space va trebui să așteptați terminarea lucrului cu structura [struct bio](#) transmisă dispozitivului fizic; utilizați câmpul [bi_end_io](#) al structurii;
- funcția apelată în momentul încheierii unui bio ([bi_end_io](#)) rulează în context întrerupere (bottom half);
- pentru sincronizarea informațiilor scrise/citite în cadrul discului virtual cu discurile fizice (flushing), modulul va trebui să exporte o operație tip ioctl (SSR_IOCTL_SYNC); acest lucru îl realizați cu ajutorul operațiilor [sync_blockdev](#) și [invalidate_bdev](#) pe dispozitivul virtual și pe dispozitivele fizice;
- pentru calculul CRC32 puteți utiliza macro-ul [crc32](#) pus la dispoziție de kernel;
- o bună sursă de documentare o reprezintă implementarea de [software RAID](#) din nucleul Linux;
- când generați structuri bio, luați în considerare că aceasta trebuie să aibă dimensiunea multiplu al sectorului de disc (`KERNEL_SECTOR_SIZE`);
- întrucât sectoarele de date sunt separate de sectoarele de CRC va trebui să construiți structuri [struct bio](#) separate pentru date și pentru valorile CRC;
- mașina virtuală dispune de două discuri suplimentare de 100 MB folosite pentru testare (`/dev/sdb` și `/dev/sdc`);
- dispozitivul virtual (`LOGICAL_DISK_NAME - /dev/ssr`) va avea capacitatea de `LOGICAL_DISK_SECTORS` (câmpul `capacity` al structurii [struct gendisk](#));
- macrodefinițiile utile se găsesc în [ssr.h](#);
- **se recomandă != este obligatoriu**; orice rezolvare care respectă cerințele temei este acceptată.

Precizări Windows

- dispozitivul va fi accesat din user space ca `LOGICAL_DISK_USER_NAME` (definit în [ssr.h](#)); va trebui să creați un dispozitiv și o legătură simbolică;
- pentru crearea dispozitivului logic folosiți apelul [IoCreateDevice](#); tipul dispozitivului trebuie să fie `FILE_DEVICE_DISK`;
- pentru crearea legăturii simbolice folosiți [IoCreateSymbolicLink](#);
- denumirile din kernel mode ale celor două dispozitive fizice sunt date de macro-urile `PHYSICAL_DISK1_DEVICE_NAME`, respectiv `PHYSICAL_DISK2_DEVICE_NAME` definite în [ssr.h](#); link-urile simbolice sunt date de macro-urile `PHYSICAL_DISK1_LINK_NAME`, respectiv `PHYSICAL_DISK2_LINK_NAME`;
- cele două dispozitive fizice vor fi accesate din user mode ca `PHYSICAL_DISK1_USER_NAME`, respectiv `PHYSICAL_DISK2_USER_NAME`;
- pentru a deschide dispozitivele fizice folosiți [IoGetDeviceObjectPointer](#); nu uitați să eliberați, după utilizare, referința la dispozitiv folosind [ObDereferenceObject](#);
- pentru a transmite o cerere read/write către dispozitivele fizice puteți folosi [IoBuildSynchronousFsdRequest](#) și [IoCallDriver](#);
- pentru calculul CRC32 nu există API, dar puteți utiliza [implementarea crc32 de aici](#) fără multe modificări;
- în cazul în care ambele valori CRC sunt incorecte, se va returna `STATUS_DEVICE_DATA_ERROR`;
- o bună sursă de documentare o reprezintă implementarea de [RAM disk](#) și [exemplele de procesare a IRP-urilor](#);
- pentru testare se vor folosi cele două discuri IDE de dimensiune 100 MB atașate mașinii virtuale;
- pentru testare proprie trebuie să aveți în vedere ca cererile trebuie să fie aliniate la 512 octeți (atât offset-ul, cât și dimensiunea datelor de citit/scriș);
- nu este nevoie de operații de tip [flush](#).

Testare

Pentru simplificarea procesului de corectare a temelor, dar și pentru a reduce greșelile temelor trimise, corectarea temelor se va face automat cu ajutorul unor teste publice ([Teste Linux](#), [Teste Windows](#)).

Testele presupun ca numele modulului de kernel este `ssr.ko` pentru Linux, respectiv `ssr.sys` pentru Windows.

Dacă, în urma procesului de testare, sectoarele de pe ambele discuri conțin date nevalide, rezultând în erori de citire ce fac imposibilă funcționarea modulului, va trebui să refaceți cele două discuri, folosind comenzile:

```
dd if=/dev/zero of=/dev/sdb bs=8M
dd if=/dev/zero of=/dev/sdc bs=8M
```

Depuneri

În corectarea temei se va ține cont de următoarele criterii de depunere:

- -1.0 implementare incorectă a cerințelor temei
- -0.2 omiterea eliberării resurselor (discuri, dispozitive etc.)
- -0.1 pentru fiecare test important picat (check corrupt, crc check, simple compare)
- -0.02 pentru alte fiecare alt test picate

De asemenea, consultați [sfaturile și depunctările generale](#).

Întrebări

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

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

Permanent link:
<http://elf.cs.pub.ro/so2/wiki/teme/tema3>

Last update: **2011/04/19 14:44**