

## 4 Sistemul de fişiere Unix

### 4.1 Structura arborescentă şi legături suplimentare

#### 4.1.1 Tipuri de fişiere şi sisteme de fişiere

În cadrul unui sistem de fişiere, apelurile sistem Unix gestionează opt tipuri de fişiere şi anume:

1. Normale (obişnuite)
2. Directori
3. Legături hard (hard links)
4. Legături simbolice (symbolic links)
5. Socketuri (sockets)
6. FIFO - pipe cu nume (named pipes)
7. Periferice caracter
8. Periferice bloc

Pe lângă aceste opt tipuri, mai există încă patru entităţi, pe care apelurile sistem le văd, din punct de vedere sintactic, tot ca şi fişiere. Aceste entităţi sunt gestionate de nucleul Unix, au suportul fizic tot în nucleu şi folosite la comunicări între procese. Aceste entităţi sunt:

9. Pipe (anonymous pipes)
10. Segmente de memorie partajată
11. Cozi de mesaje
12. Semafoare

În acest capitol, sau în cele care urmează, vom trata aceste entităţi, exceptând 5, 10, 11, şi 12 care vor face obiectul unei lucrări viitoare, destinate special sistemelor de operare distribuite.

*Fişierele obişnuite* sunt privite ca şiruri de octeţi, accesul la un octet putându-se face fie secvenţial, fie direct prin numărul de ordine al octetului.

*Fişierele directori.* Un fişier director se deosebeşte de un fişier obişnuit numai prin informaţia conţinută în el. Un director conţine lista de nume şi adrese pentru fişierele subordonate lui. Uzual, fiecare utilizator are un director propriu care punctează la fişierele lui obişnuite, sau la alţi subdirectori definiţi de el.

*Fişierele speciale.* În această categorie putem include, pentru moment, ultimele 6 tipuri de fişiere. În particular, Unix priveşte fiecare dispozitiv de I/O ca şi un fişier de tip special. Din punct de vedere al utilizatorului, nu există nici o deosebire între lucrul cu un fişier disc normal şi lucrul cu un fişier special.

Fiecare director are două intrări cu nume speciale şi anume:

- " ." (punct) denumeşte generic (punctează spre) însuşi directorul respectiv;
- " . ." (două puncte succesive), denumeşte generic (punctează spre) directorul părinte.

Fiecare sistem de fişiere conţine un director principal numit **root** sau **/**.

În mod obișnuit, fiecare utilizator folosește un *director curent*, atașat utilizatorului la intrarea în sistem. Utilizatorul poate să-și schimbe acest director (`cd`), poate crea un nou director subordonat celui curent, (`mkdir`), să șteargă un director (`rmdir`), să afișeze *calea* de acces de la `root` la un director sau fișier (`pwd`) etc.

Apariția unui mare număr de distribuitori de Unix a condus, inevitabil, la proliferarea unui număr oarecare de "*sisteme de fișiere extinse*" proprii acestor distribuitori. De exemplu:

- Solaris utilizează sistemul de fișiere `ufs`;
- Linux utilizează cu precădere sistemul de fișiere `ext2` și mai nou, `ext3`;
- IRIX utilizează `xfs`
- etc.

Actualele distribuții de Unix permit utilizarea unor sisteme de fișiere proprii altor sisteme de operare. Printre cele mai importante amintim:

- Sistemele FAT și FAT32 de sub MS-DOS și Windows 9x;
- Sistemul NTFS propriu Windows NT și 2000.

Din fericire, aceste extinderi sunt transparente pentru utilizatorii obișnuiți. Totuși, se recomandă prudență atunci când se efectuează altfel de operații decât citirea din fișierele create sub alte sisteme de operare decât sistemul curent. De exemplu, modificarea sub Unix a unui octet într-un fișier de tip `doc` creat cu Word sub Windows poate ușor să compromită fișierul așa încât el să nu mai poată fi exploatat sub Windows!

Administratorii sistemelor Unix trebuie să țină cont de sistemele de fișiere pe care le instalează și de drepturile pe care le conferă acestora vis-a-vis de userii obișnuiți.

Principiul structurii arborescente de fișiere este acela că orice fișier sau director are un singur părinte. Automat, pentru fiecare director sau fișier există o singură cale (path) de la rădăcină la directorul curent. Legătura între un director sau fișier și părinte o vom numi *legătură naturală*. Evident ea se creează odată cu crearea directorului sau fișierului respectiv.

#### 4.1.2 Legături hard și legături simbolice

În anumite situații este utilă partajarea unei porțiuni a structurii de fișiere între mai mulți utilizatori. De exemplu, o bază de date dintr-o parte a structurii de fișiere trebuie să fie accesibilă mai multor utilizatori. Unix permite o astfel de partajare prin intermediul *legăturilor suplimentare*. O legătură suplimentară permite referirea la un fișier pe alte căi decât pe cea naturală. Legăturile suplimentare sunt de două feluri: *legături hard* și *legături simbolice (soft)*.

*Legăturile hard* sunt identice cu legăturile naturale și ele pot fi create numai de către administratorul sistemului. O astfel de legătură este o intrare într-un director care punctează spre o substructură din sistemul de fișiere spre care punctează deja legătura lui naturală. Prin aceasta, substructura este văzută ca fiind descendentă din două directoare diferite! Deci, printr-o astfel de legătură un fișier primește efectiv două nume. Din această cauză, la parcurgerea unei structuri arborescente, fișierele punctate prin legături hard apar duplicate. Fiecare duplicat apare cu numărul de legături către el.

De exemplu, dacă există un fișier cu numele `vechi`, iar administratorul dă comanda:

```
#ln vechi linknou
```

atunci în sistemul de fișiere se vor vedea două fișiere identice: `vechi` și `linknou`, fiecare dintre ele având marcat faptul că sunt două legături spre el.

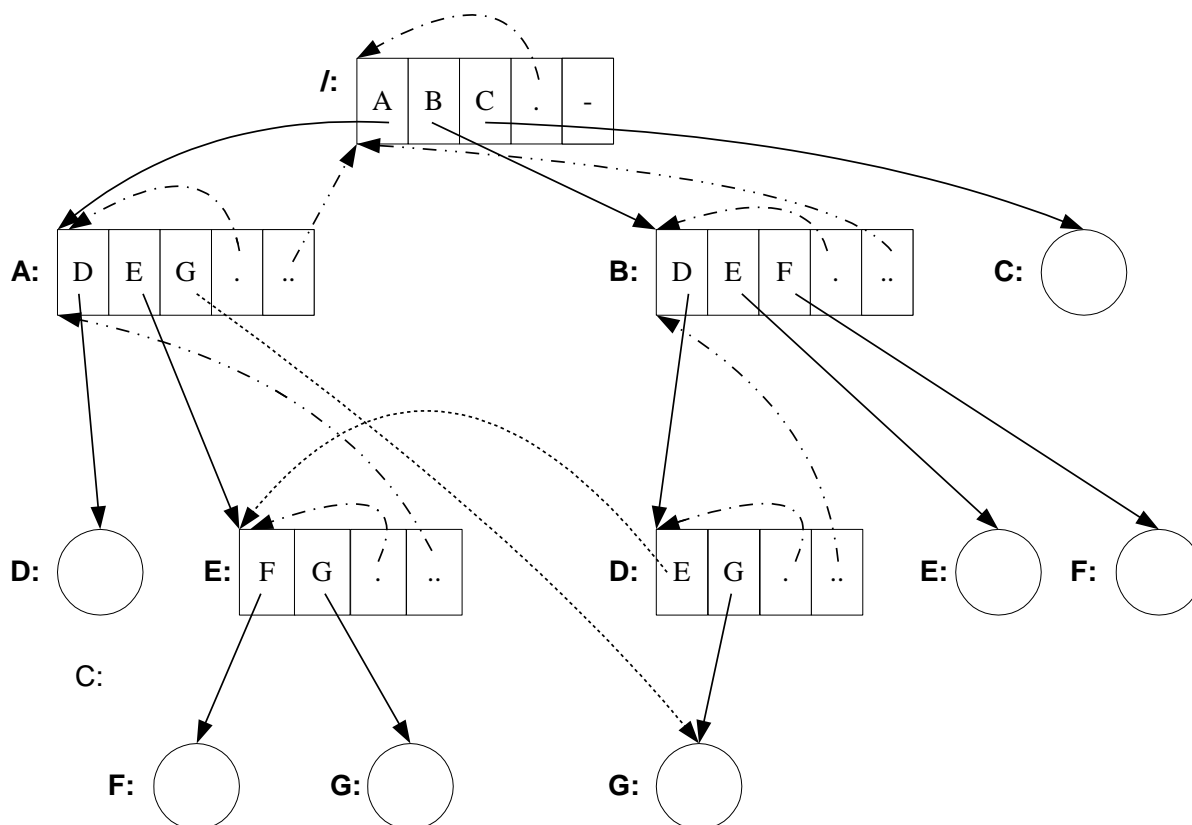
Legăturile hard pot fi făcute numai în interiorul aceluiași sistem de fișiere (detalii puțin mai târziu).

*Legăturile simbolice* sunt intrări speciale într-un director, care punctează (referă) un fișier (sau director) oarecare în structura de directori. Această intrare se comportă ca și un subdirector al directorului în care s-a creat intrarea.

În forma cea mai simplă, o legătură simbolică se creează prin comanda:

```
ln -s caleInStructuraDeDirectori numeSimbolic
```

După această comandă, `caleInStructuraDeDirectori` va avea marcată o legătură în plus, iar `numeSimbolic` va indica (numai) către această cale. Legăturile simbolice pot fi utilizate și de către userii obișnuiți. De asemenea, ele pot puncta și înafara sistemului de fișiere (detalii puțin mai târziu).



**Figura 4.1 O structură arborescentă cu legături**

Structura arborescentă împreună cu legăturile simbolice sau hard conferă sistemului de fișiere Unix o structură de graf aciclic. În fig. 4.1 este prezentat un exemplu simplu de structură de fișiere. Prin literele mari A, B, C, D, E, F, G am indicat nume de fișiere obișnuite, nume de directori și nume de legături. Este evident posibil ca același nume să apară de mai multe ori în structura de directori, grație structurii de directori care elimină ambiguitățile. Fișierele obișnuite sunt reprezentate prin cercuri, iar fișierele directori prin dreptunghiuri.

Legăturile sunt reprezentate prin săgeţi de trei tipuri:

- linie continuă – legăturile naturale;
- linie întreruptă – spre propriul director şi spre părinte;
- linie punctată – legături simbolice sau hard.

În fig. 4.1 există 12 noduri - fişiere obişnuite sau directori. Privit ca un arbore, deci considerând numai legăturile naturale, el are 7 ramuri şi 4 nivele.

Să presupunem că cele două legături (desenate cu linie punctată) sunt simbolice. Pentru comoditate, vom nota legătura simbolică cu ultima literă din specificarea căii. Crearea celor două legături se poate face, de exemplu, prin succesiunea de comenzi:

```
cd /A
ln -s /A/B/D/G G      Prima legătură
cd /A/B/D
ln -s /A/E E          A doua legătură
```

Să presupunem acum că directorul curent este B. Vom parcurge arborele în ordinea director urmat de subordonaţii lui de la stânga spre dreapta. Următoarele 12 linii indică toate cele 12 noduri din structură. Pe aceeaşi linie apar, atunci când este posibil, mai multe specificări ale acelui nod. Specificările care fac uz de legături simbolice sunt subliniate. Cele mai lungi 7 ramuri vor fi marcate cu un simbol # în partea dreaptă.

/	..				
/A	../A				
/A/D	../A/D				#
/A/E	../A/E	<u>D/E</u>	<u>../D/E</u>		
/A/E/F	../A/E/F	<u>D/E/F</u>	<u>../D/E/F</u>		#
/A/E/G	../A/E/G	<u>D/E/G</u>	<u>../D/E/G</u>		#
/B	.				
/B/D	D	<u>../D</u>			
/B/D/G	D/G	<u>../D/G</u>	<u>../A/G</u>	<u>../A/G</u>	#
/B/E	E	<u>../E</u>			#
/B/F	F	<u>../F</u>			#
/C	../C				#

Ce se întâmplă cu ştergerea în cazul legăturilor multiple? De exemplu, ce se întâmplă când se execută una dintre următoarele două comenzi?

```
rm D/G
rm /A/G
```

Este clar că fişierul trebuie să rămână activ dacă este şters numai de către una dintre specificări.

Pentru aceasta, în descriptorul fişierului respectiv există un câmp numit *contor de legare*. Acesta are valoarea 1 la crearea fişierului şi creşte cu 1 la fiecare nouă legătură. La ştergere, se radiază legătura din directorul părinte care a cerut ştergerea, iar contorul de legare scade cu 1. Abia dacă acest contor a ajuns la zero, fişierul va fi efectiv şters de pe disc şi blocurile ocupate de el vor fi eliberate.

### 4.1.3 Conceptul de montare

Spre deosebire de alte sisteme de operare ca DOS, Windows, etc. în specificarea fișierelor Unix *nu apare zona de periferic*. Acest fapt nu este întâmplător, ci este cauzat de filozofia generală de acces la fișierele Unix. Conceptul esențial prin care se rezolvă această problemă este cunoscut în Unix prin termenii de *montare* și *demontare* a unui sistem de fișiere.

Operația de *montare* constă în conectarea unui sistem de fișiere, de pe un anumit disc, la un director existent pe sistemul de fișiere implicit. Administratorul lansează comanda de montare sub forma:

```
# mount [ optiuni ] sistemDeFișiere directorDeMontare
```

Efectul este conectarea indicată prin *sistemDeFișiere* la *directorDeMontare* existent pe sistemul implicit de fișiere. Opțiunile pot să indice caracteristicile montării. De exemplu opțiunea *rw* permite atât citirea, cât și scrierea în subsistemul montat, în timp ce opțiunea *ro* permite numai citirea din subsistemul montat. Opțiunea *-t* indică tipul sistemului de fișiere care se montează, și în funcție de tip, argumentul *sistemDeFișiere* poate fi */dev/periferic* sau *dev/dsk/periferic* sau */root/periferic* ș.a.m.d. Pentru detalii se va putea consulta manualul *mount* al sistemului de operare curent.

Operația de *demontare* are efectul invers și ea se face cu comanda:

```
#/etc/unmount directorDeMontare
```

Să urmărim cele ilustrate în fig. 4.2. În fig. 4.2a este dată structura sistemului de fișiere activ. Directorul B este vid (nu are descendenți). În fig. 4.2b este dată structura de fișiere de pe un disc aflat (să zicem) pe unitatea de disc nr. 3, cu care intenționăm să lucrăm. Pentru aceasta, se va da comanda *privilegiată* Unix:

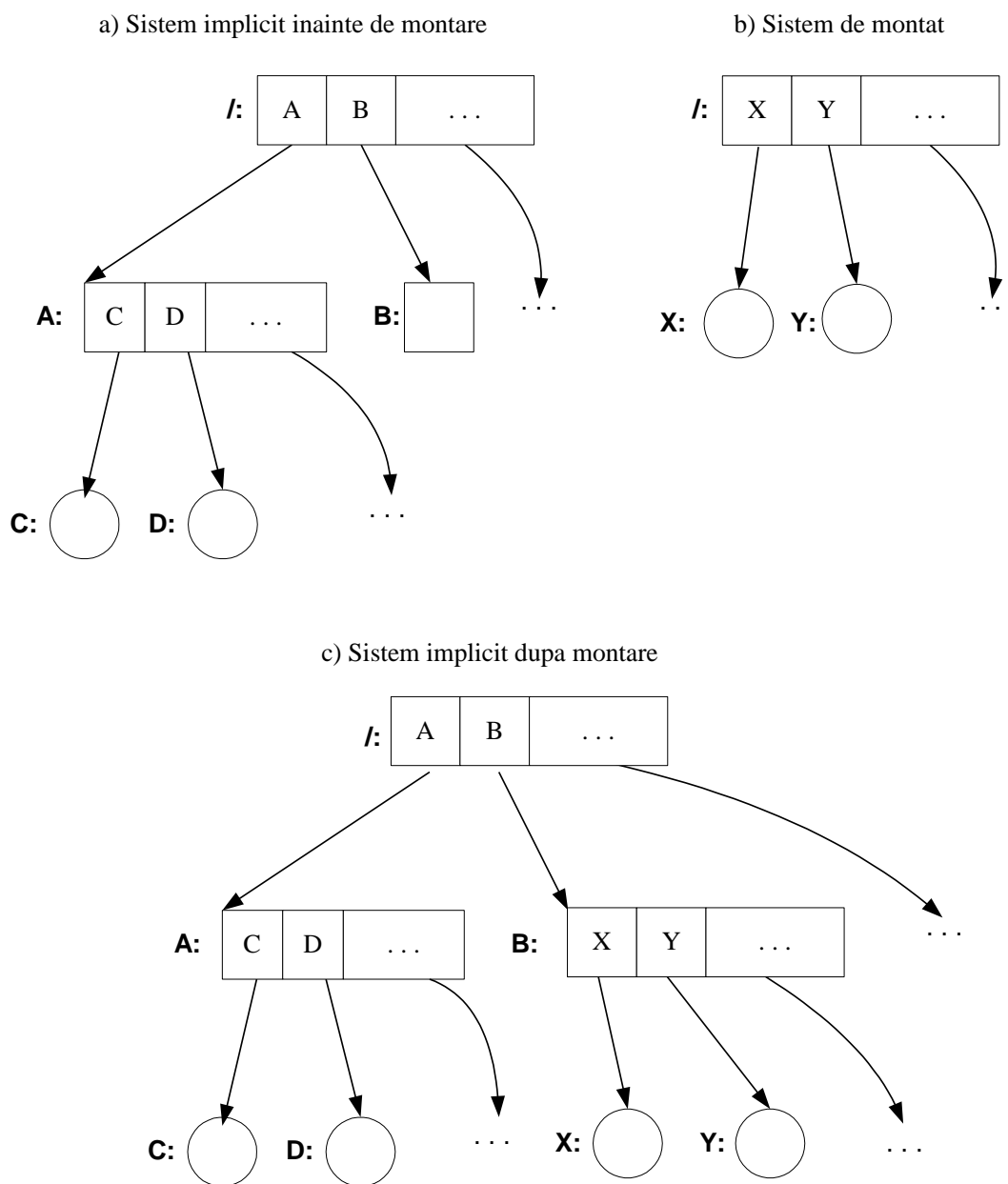
```
#/etc/mount /dev/fd3 /B
```

Prin ea se indică legarea discului al cărui fișier special (driver) poartă numele */dev/fd3* (și care se referă la discul nr. 3), la directorul vid */B*. Urmare a legării, se obține structura de fișiere activă din fig. 4.2c.

Cu scuzele de rigoare față de cititorul pe care-l plictisim, vom face câteva precizări. Deși ele sunt firești, practica arată că sunt de multe ori încălcate, ceea ce provoacă neplăceri (și nu numai atât).

- Nu se va cere accesul la un fișier de pe un disc decât dacă acesta este montat în structura de fișiere implicită.
- Nu se va cere demontarea unei substructuri decât dacă a fost montată în prealabil.
- *Scoaterea unui suport montat de pe o unitate și înlocuirea lui cu un alt suport poate provoca pagube mari*. Să presupunem, spre exemplu, că s-a montat o structură existentă pe o anumită dischetă. Dacă înaintea demontării se scoate discheta din unitate și se introduce alta în loc, atunci este posibil să se piardă informații de pe noua dischetă, și, în unele cazuri, este posibilă blocarea întregului sistem! De altfel, sistemele Unix, mai nou, nici nu permit scoaterea din unitate a unui suport până când nu se efectuează operația *unmount*.

- Nu este posibilă efectuarea de legături simbolice decât dacă directorul de unde se leagă și fișierul / directorul care se leagă se află pe același suport fizic. (Deși pe unele sisteme s-ar putea ca o astfel de legare să fie posibilă, noi nu o recomandăm :)



**Figura 4.2 Operația de montare**

În practica implementărilor Unix, la încărcarea **SO** se fac automat o serie de operații de montare, în conformitate cu configurarea sistemului. Indicațiile de montare automată sunt trecute în fișierul `/etc/fstab`. Acesta este un fișier text, având pe fiecare linie câte o montare, în care se specifică:

- partiția Unix (periferic, sistemDeFișiere) care se montează;
- directorDeMontare sub care este montată partiția;
- tipul sistemului de fișiere conținut;
- diverse opțiuni.

Iată, spre exemplu, o porțiune din acest fișier:

/dev/hda2	/	ext3	defaults	1	1
/dev/hda7	/home	ext3	defaults,nosuid,nodev,noexec		
/dev/cdrom	/mnt/cdrom	iso9660	noauto,owner,ro	0	0
/dev/fd0	/mnt/floppy	auto	noauto,owner	0	0
/dev/hda3	/usr	ext3	defaults,nodev	1	2
/dev/hda5	/usr/local	ext3	defaults,nodev	1	2
/dev/hda6	/var	ext3	defaults,nosuid,nodev,noexec		
/dev/hda1	swap	swap	defaults	0	0

De asemenea, operația de montare este practică pentru *înglobarea SO Unix într-o rețea de calculatoare*. Cele mai cunoscute astfel de sisteme care dirijează montările de fișiere în rețele Unix (și nu numai) sunt *NFS* (*Network File System* al firmei Sun Microsystems), *RFS* (*Remote File System* al firmei AT&T), *SAMBA* (produs open source). Sistemul care montează o structură de directoare de pe o altă mașină poartă numele de *client* (NFS, RFS, SAMBA, etc.). Mașina care oferă structura spre montare (exportă) se numește *server* (NFS, RFS, SAMBA, etc.).

Montările pentru NFS sunt specificate tot în fișierul `/etc/fstab`, prin linii de forma:

server1:/export/home	/home	nfs	rw,bg,intr	0	0
server1:/export/usr/local	/usr/local	nfs	rw,bg,intr	0	0
server2:/export/var/spool/mail	/var/spool/mail	nfs	rw,bg,intr	0	0

În acest context se poate da o nouă caracterizare a diferențelor dintre legăturile hard și cele simbolice: legăturile hard funcționează numai în interiorul aceluiași sistem de fișiere, în timp ce legăturile simbolice pot puncta și spre noduri ale altui sistem de fișiere montat împreună cu sistemul de fișiere ce conține legătura limboică.

#### 4.1.4 Protecția fișierelor Unix

##### 4.1.4.1 Drepturi de acces

Reamintim iarăși principalele entități participante la Unix: utilizatori, fișiere, procese și mai ales reamintim faptul că între ele există multe interdependențe. În această secțiune vom trata una dintre aceste interdependențe.

Vis a vis de un fișier sau de un director, utilizatorii se împart în trei categorii:

- proprietarul fișierului (*u* - *user*).
- grupul de utilizatori (*g* - *group*), de exemplu o grupă de studenți participă la un același proiect, motiv pentru care administratorul poate constitui un astfel de grup, cu drepturi specifice – de regulă mai slabe decât ale proprietarului, dar mai puternice decât ale restului utilizatorilor.
- restul utilizatorilor (*o* – *others*) cei care nu sunt în primele două categorii.

Nucleul SO Unix identifică utilizatorii prin numere naturale asociate unic, numite UID-uri (User IDentifications). De asemenea, identifică grupurile de utilizatori prin numere numite GID-uri (Group IDentifications).

Un utilizator aparte, cu drepturi depline asupra tuturor fișierelor este *root* sau *superuserul*.

Pentru fiecare categorie de utilizatori, fișierul permite maximum trei drepturi:

- dreptul de citire (*r* – *read*)
- dreptul de scriere (*w* – *write*) care include crearea de subdirectori, stergerea de subdirectori, adăugarea sau ștergerea de intrări în director, modificarea fișierului, etc.
- dreptul de execuție (*x* – *execution*) care permite lansarea în execuție a unui fișier. Acest drept, conferit unui director, permite accesul în directorul respectiv (*cd*).

În consecință, pentru specificarea drepturilor de mai sus asupra unui fișier sau director sunt necesari 9 (nouă) biți. Reprezentarea externă a acestei configurații se face printr-un grup de 9 (nouă) caractere: *rwxrwxrwx* în care absența unuia dintre drepturi la o categorie de useri este indicată prin – (minus).

Modul de atribuire a acestor drepturi se poate face cu ajutorul comenzii Shell *chmod*. Cea mai simplă formă a ei este:

```
chmod o1o2o3 fișier . . .
```

unde o<sub>1</sub>o<sub>2</sub>o<sub>3</sub> sunt trei cifre octale. Biții 1 ai primeia indică drepturile userului, biții 1 ai celei de-a doua indică drepturile grupului, iar biții 1 ai celei de-a treia cifră indică drepturile restului userilor. De exemplu, comanda:

```
chmod 754 A B
```

fixează la fișierele A și B toate drepturile pentru user, drepturile de citire și de execuție pentru grup și doar dreptul de citire pentru ceilalți useri. În urma acestei comenzi, drepturile fișierelor A și B vor deveni: *rwxr-xr--*

Invităm cititorul să consulte manualul comenzii *chmod* pentru prezentarea completă a acestei comenzi.

#### 4.1.4.2 Drepturi implicite: *umask*

În momentul intrării unui user în sistem, acestuia i se vor acorda niște drepturi implicite pentru fișiere nou create. Toate fișierele și directoarele create de user pe durata sesiunii de lucru îl vor avea ca proprietar, iar drepturile vor fi cele primite implicit. Fixarea drepturilor implicite, sau aflarea valorii acestora se poate face folosind comanda *umask*. Drepturile implicite sunt stabilite scăzând (octal) masca definită prin *umask* din 777.

Pentru a se afla valoarea măștii se lansează:

```
umask
```

Rezultatul este afișarea măștii, de regulă 022. Deci drepturile implicite vor fi: 777-022=755, adică userul are toate drepturile, iar grupul și restul vor avea doar drepturi de citire și de execuție.



Dacă se doreşte schimbarea acestei măşti pentru a da userului toate drepturile, grupului de citire şi execuţie iar restului nici un drept, adică drepturile implicite 750, fiindcă  $777-027=750$ . Deci se va da comanda:

```
umask 027
```

Efectul ei va rămâne valabil până la un nou `umask` sau până la încheierea sesiunii.

#### 4.1.4.3 Drepturi de lansare, drepturi program executabil, biţii `setuid` şi `setgid`

În această secţiune vom analiza, din punct de vedere al drepturilor de acces, relaţia dintre un utilizator, programul executabil pe care îl lansează şi fişierele asupra cărora acţionează programul în cursul execuţiei lui.

Pentru fixarea ideilor, să considerăm un exemplu. Presupunem că:

1. un utilizator `U`, care face parte dintr-un grup `G`, lansează în execuţie un program `P`. Pe durata execuţiei, programul `P` acţionează asupra unui fişier `F`.
2. programul `P` are ca proprietar utilizatorul `UP` care face parte din grupul `GP`, iar drepturile fişierului executabil `P` sunt `rwxr-xr-x`.
3. fişierul `F` asupra căruia se acţionează are proprietarul `UF` care face parte din grupul `GF`, iar drepturile fişierului `F` sunt `rwxr-xr--`.

În aceste ipoteze, userul `U` poate să lanseze în execuţie programul `P`, deoarece acesta conferă drepturi de execuţie tuturor utilizatorilor. (Dacă drepturile lui `P` ar fi fost `rwxr--r--`, atunci lansarea în execuţie ar fi fost posibilă numai dacă `U = UP`. Dacă drepturile lui `P` ar fi fost `rwxr-xr--`, atunci lansarea ar fi fost posibilă numai dacă `U = UP` sau `G = GP`.)

După lansarea în execuţie de către `U` a lui `P`, în mod implicit acţiunile pe care programul `P` le poate efectua asupra fişierului `F` sunt cele permise de drepturile pe care `U` le are asupra lui `F`. În exemplul nostru, dacă `U = UF` atunci `P` poate citi, scrie sau executa `F`. Dacă `G = GF` atunci `P` poate citi sau executa `F`, iar dacă `G` şi `GF` sunt diferite atunci `P` poate numai să citească din `F`. (Într-o secţiune următoare vom arăta cum un program poate să lanseze în execuţie un alt program).

Această regulă, prin care drepturile lansatorului de program permit sau nu unele operaţii pe care programul lansat le aplică unui fişier, este regula cvasigenerală de acţiune asupra fişierelor.

Există însă situaţii, nu foarte frecvente, în care această regulă se poate schimba. Schimbarea este cunoscută sub numele *setuid* / *setgid* şi se referă, cu notaţiile de mai sus, la faptul că: După lansarea în execuţie de către `U` a lui `P`, în mod setuid / setgid acţiunile pe care programul `P` le poate efectua asupra fişierului `F` sunt cele permise de drepturile proprietarului programului `P` şi ale grupului din care face parte acesta.

Este vorba de doi biţi importanţi relativ la drepturile de acces, aşa numiţii *bit setuid* (*set-user-id*) care pus pe 1 schimbă drepturile lui `U` cu drepturile lui `UP` şi bitul *setgid*, care pus pe 1 schimbă drepturile lui `G` cu drepturile lui `GP`. (Aceste schimbări au loc numai la execuţia programului `P`.)

Cu alte cuvinte, dacă pentru un fișier executabil bitul `setuid` este 1, atunci **un utilizator care lansează în execuție acest fișier** (evident, dacă are dreptul să-l lanseze) **primește, pe timpul execuției, aceleași drepturi de acces la resurse** (fișiere, semafoare, zone de memorie etc.) **ca și proprietarul fișierului executabil.**

Să vedem o situație concretă în care este utilă folosirea bitului `setuid`. Un utilizator cu numele **profesor** întreține un fișier **note** al cărui proprietar este. Din rațiuni lesne de înțeles, drepturile fișierului **note** sunt fixate la **`rw-----`**. Utilizatorul **profesor** dorește să permită utilizatorilor din grupul **studenti** să vadă unele informații din fișierul **note**.

Pentru aceasta, **profesor** creează un fișier executabil **examen** (proprietarul lui **examen** este **profesor**) care permite citirea (eventual selectivă) de informații din fișierul **note**. Proprietarul atribuie pentru **examen** drepturile **`rw-x--x--x`** și pune bitul `setuid` al lui **examen** pe 1. Această atribuire se face cu comanda **`chmod +s examen`**. Noile drepturi afișate ale lui **examen** sunt: **`rws--x--x`**

Utilizatorii **studenti**, în momentul lansării programului **examen**, primesc aceleași drepturi de acces la fișiere ca și **profesor**. În particular programul **examen** poate accesa fișierul **note** (vezi drepturile acestui fișier) chiar dacă el nu a fost lansat în execuție de către **profesor**. În absența lui `setuid` pentru **examen**, acesta poate fi, totuși lansat în execuție, însă nu poate să acceseze fișierul **note**.

După cum spuneam mai sus, nucleul SO Unix identifică utilizatorii prin numere naturale asociate unic, numite UID-uri (User IDentifications). De asemenea, identifică grupurile de utilizatori prin numere numite GID-uri (Group IDentifications).

Pe parcursul execuției programului **examen** acestuia i se mai asociază în plus identificatorul EUID (effective UID), care coincide cu UID-ul lui **profesor**, prin care asigură accesul la resurse.

Mecanismul `setuid` permite o foarte elastică manevrare a fișierelor. În schimb, dacă superuserul gestionează prost acest mecanism, atunci potențialii infractori au un câmp larg de acțiune. Să considerăm, din rațiuni evidente, doar un scenariu simplu: Presupunem ca **root** este proprietar al unui fișier executabil cu bitul `setuid` setat și cu drept de scriere pentru alții. În această situație, un răuvoitor poate să modifice acest fișier executabil așa încât să aibă o acțiune malefică ce presupune acces la resurse ale superuserului!. Acțiunea se va putea executa deoarece EUID-ul este UID-ul lui **root**!

Modificarea drepturilor de acces la fișiere se poate face numai de către proprietarul fișierului (sau de către superuser), folosind comanda `chmod`. Modificarea proprietarului sau a grupului se poate face, în aceleași condiții folosind comanda `chown`.

Un exemplu tipic în care se folosește `setuid` este comanda `/usr/bin/passwd`. Aceasta este lansată de către fiecare utilizator atunci când dorește să-și schimbe parola. Efectul ei se răsfrânge asupra fișierului `/etc/shadow`. În acest scop, se stabilesc drepturile:

```
-r-s--x--x  1 root  root  22312 Sep 25 18:52 /usr/bin/passwd
-r-----  1 root  root  10256 Mar  2 14:40 /etc/shadow
```

Deci programul `passwd` are `setuid`, ceea ce permite accesul la `/etc/shadow` numai prin programul `/usr/bin/passwd`.

#### 4.1.5 Principalele directoare ale unui sistem de fişiere Unix

De-a lungul evoluţiei sistemelor din familia Unix, partea superioară a structurii sistemului de fişiere a avut mai mult sau mai puţin o formă standard. De fapt, fiecare versiune Unix şi-a fixat o structură specifică a părţii superioare din sistemul de fişiere. Diferenţele între aceste structuri nu sunt prea mari. Mai mult, din raţiuni de compatibilitate, versiunile mai noi definesc legături suplimentare hard (nu legături simbolice), pentru a asigura compatibilitatea cu sistemele de fişiere mai vechi. Din această cauză este cel mai nimerit să se studieze o reuniune a celor mai răspândite structuri. O astfel de structură este prezentată de noi în fig. 4.3.

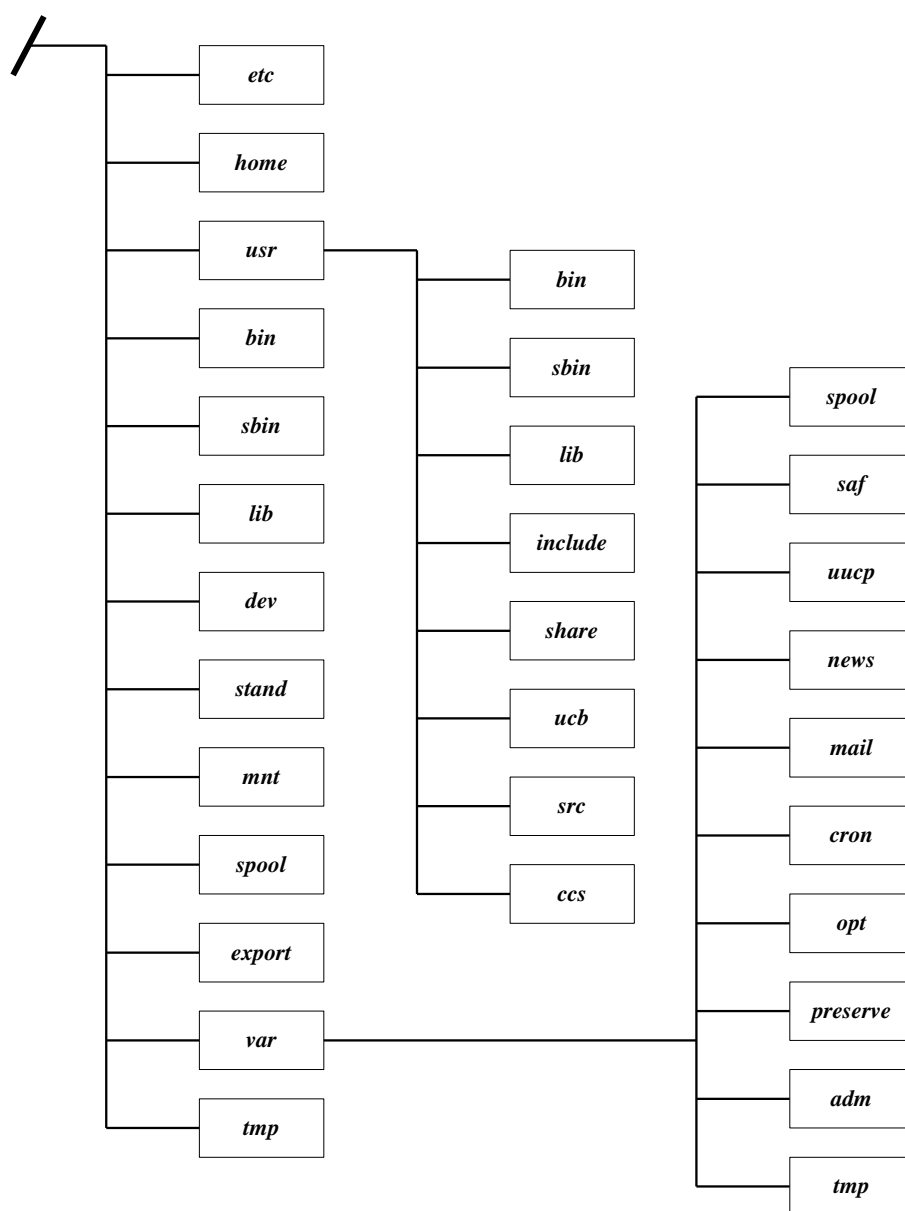


Figura 4.3 Structura superioară a unui sistem de fişiere Unix

Directorul `/etc` conține informații specifice mașinii necesare întreținerii sistemului. Acestea sunt de fapt datele restrictive și periculoase din sistem. Un exemplu de fișier ce conține informații specifice mașinii este, spre exemplu, `/etc/rc2.d` care este un director ce conține programe shell executate de procesul *init* la schimbarea stării. Tot aici sunt plasate fișierele `/etc/passwd`, `/etc/group`, `/etc/shadow` care sunt folosite pentru administrarea utilizatorilor.

Directorul `/home` este folosit pentru directorii gazdă ai utilizatorilor. În momentul intrării în sistem, fiecare utilizator indică directorul lui, care este fixat ca director curent (*home directory*) în `/etc/passwd`. Deși se poate trece ușor de la un director la altul, în mod normal fiecare utilizator rămâne în propriul lui director, unde își dezvoltă propria lui structură arborescentă de directori.

Directorul `/usr` este folosit în mod tradițional pentru a stoca fișiere ce pot fi modificate. La primele versiuni de Unix conținea și fișierele utilizatorilor. În prezent este punctul de montare pentru partiția ce conține `usr`. El conține programe executabile dependente și independente de sistemul de fișiere, precum și fișiere necesare acestora, dar care nu cresc dinamic. Asupra conținutului de subdirectoare ale lui `/usr` vom reveni puțin mai târziu.

Directorul `/bin` conține programele principalelor comenzi standard Unix: compilatoare, asamblare, editoare de texte, utilitare etc. Versiunile mai noi de Unix plasează acest director în `/usr/bin`.

Directorul `/sbin` (super-utilizator bin) conține comenzi critice pentru procedura de încărcare a sistemului. Orice comenzi administrative care necesită lucrul mono-utilizator sunt în acest director. Copii ale acestor comenzi se află în directoarele `/usr/bin` și `/usr/sbin`, astfel încât el (`/sbin`) poate fi refăcut dacă este necesar.

Directorul `/lib` conține diverse biblioteci și baze de date necesare apelurilor sistem. Doar versiunile mai vechi de Unix plasează acest director în `/`, cele actuale îl plasează în `/usr/lib`.

Directorul `/dev` este folosit pentru memorarea fișierelor speciale (fișierele devices). Practic, fiecare tip de terminal și fiecare tip de unitate de disc trebuie să aibă asociat un astfel de fișier special. Incepând cu *SVR4* se permite ca în `dev` să existe și subdirectoare care să grupeze astfel de device-uri.

Directorul `/stand` conține informațiile necesare încărcării sistemului.

Directorul `/mnt` este folosit pentru a monta un sistem de fișiere temporar. De exemplu, un sistem de fișiere de pe un disc flexibil poate fi montat în `/mnt` pentru a verifica fișierele de pe această dischetă.

Directorul `/spool` este plasat în `/` doar în versiunile mai vechi de Unix. În el sunt memorate fișierele tampon temporare destinate prelucrărilor asincrone: listări asincrone de fișiere (`/spool/lpd`) și execuția la termen a unor comenzi (`/spool/at`).

Directorul `/export` este folosit ca punct implicit de montare pentru un arbore de sistem de fișiere exportat, pentru fișierele în rețea gestionate prin pachetul NFS (Network File System).

Directorul `/var` este folosit pentru memorarea fișierelor care cresc dinamic. În particular, multe dintre versiunile de Unix țin în acest director fișierele `INBOX` cu căsuțele poștale ale utilizatorilor. Structura de subdirectoare a lui `/var` o vom descrie ceva mai încolo.

Directorul `/tmp` este folosit pentru a memora fișiere temporare pentru aplicații. În mod normal, aceste fișiere nu sunt salvate și sunt șterse după o perioadă de timp.

În continuare vom descrie pe scurt conținutul directorului `/usr`. După cum se poate vedea, unele dintre subdirectoare sunt plasate (și au fost descrise) la nivel de rădăcină `/`. Versiunile mai noi de Unix, începând cu *SVR4*, le-au coborât din rădăcină ca subdirectoare ale lui `/usr`. Este cazul directoarelor `bin`, `sbin`, `lib`.

Directorul `/usr/include` conține fișierele header (`*.h`) standard ale limbajului C de sub Unix.

Directorul `/usr/share` conține o serie de directoare partajabile în rețea. În multe dintre sistemele noi, în el se află directoarele: `man` cu manualele Unix, `src` cu sursele C ale nucleului Unix și `lib` mai sus prezentat.

Directorul `/usr/ucb` conține programele executabile compatibile Unix BSD.

Directorul `/usr/src` conține textele sursă C ale nucleului Unix de pe mașina respectivă.

Directorul `/usr/ccs` conține instrumentele de dezvoltare a programelor C oferite de Unix: `cc`, `gcc`, `dbx`, `cb`, `indent`, etc.

În continuare vom descrie conținutul directorului `/var`. Ca și mai sus, unele subdirectoare de la nivelele superioare au fost mutate aici de către versiunile mai noi de Unix. Este vorba de directoarele `spool` și `tmp`.

Directorul `/var/saf` conține fișiere jurnal și de contabilizare a serviciilor oferite.

Directorul `/var/uucp` conține programele necesare efectuării de copii de fișiere între sisteme Unix (Unix to Unix CoPy). Acest gen de servicii este primul pachet de comunicații instalat pe sisteme Unix, este operațional încă din 1978 și este utilizat și astăzi atunci când nu există alt mijloc mai modern de comunicații. Un astfel de sistem permite, de exemplu, apelul telefonic între două sisteme Unix, iar după luarea contactului cele două sisteme își schimbă între ele o serie de fișiere, ca de exemplu mesajele de poștă electronică ce le sunt destinate.

Directorul `/var/news` conține fișierele necesare serviciului de (știri) noutăți (`news`) care poate fi instalat pe mașinile Unix.

Directorul `/var/mail` conține căsuțele poștale implicite ale utilizatorilor (`INBOX`). Pe unele sisteme, ca de exemplu pe Linux, acestea se află în `/var/spool/mail`.

Directorul `/var/cron` conține fișierele jurnal necesare serviciilor executate la termen.

Directorul `/var/opt` constituie un punct de montare pentru diferite pachete de aplicații.

Directorul `/var/preserve` conține, la unele implementări Unix (SVR4) fișiere jurnal destinate refacerii stării editoarelor de texte “picate” ca urmare a unor incidente.

Directorul `/var/adm` conține fișiere jurnal (log-uri) de contabilizare și administrare a sistemului. La versiunile mai noi acestea sunt în `/var/log`.

După cum se poate vedea ușor, structura de directori Unix începând de la rădăcină este relativ dependentă de tipul și versiunea de Unix. De fapt, este vorba de “Unix vechi” și “Unix noi”. De asemenea, multe dintre directoare au fost înlocuite sau li s-a schimbat poziția în structura de directori. Tabelul de mai jos prezintă câteva corespondențe între vechile și noile plasări de fișiere.

Nume vechi	Nume nou
<code>/bin</code>	<code>/usr/bin</code>
<code>/lib</code>	<code>/usr/lib</code>
<code>/usr/adm</code>	<code>/var/adm</code>
<code>/usr/spool</code>	<code>/usr/spool</code>
<code>/usr/tmp</code>	<code>/var/tmp</code>
<code>/etc/termcap</code>	<code>/usr/share/lib/termcap</code>
<code>/usr/lib/terminfo</code>	<code>/usr/share/lib/terminfo</code>
<code>/usr/lib/cron</code>	<code>/etc/cron.d</code>
<code>/usr/man</code>	<code>/usr/share/man</code>
<code>/etc/&lt;programe&gt;</code>	<code>/usr/bin/&lt;programe&gt;</code>
<code>/etc/&lt;programe&gt;</code>	<code>/sbin/&lt;programe&gt;</code>

## 4.2 Structura internă a discului Unix

### 4.2.1 Partiții și blocuri

Un sistem de fișiere Unix este găzduit fie pe un periferic oarecare (hard-disc, CD, dischetă etc.), fie pe o *partiție* a unui hard-disc. Partiționarea unui hard-disc este o operație (relativ) independentă de sistemul de operare ce va fi găzduit în partiția respectivă. De aceea, atât partițiilor, cât și suporturilor fizice reale le vom spune generic, *discuri Unix*.

Blocul 0	- bloc de boot
Blocul 1	- Superbloc
Blocul 2	- inod
- - - - -	- - - - -
Blocul n	- inod
Blocul n+1	zona fișier
- - - - -	- - - - -
Blocul n+m	zona fișier

**Figura 4.4 Structura unui disc Unix**

Un fișier Unix este o succesiune de octeți, fiecare octet putând fi adresat în mod individual. Este permis atât accesul secvențial, cât și cel direct.

Unitatea de schimb dintre disc şi memorie este *blocul*. La sistemele mai vechi acesta are 512 octeţi, iar la cele mai noi poate avea şi 1Ko sau 4Ko, pentru o mai eficientă gestiune a spaţiului.

Un sistem de fişiere Unix este o structură de date rezidentă pe disc. Aşa după cum se vede din fig. 4.4, un disc este compus din patru categorii de blocuri.

Blocul 0 conţine programul de încărcare al **SO**. Acest program este dependent de maşina sub care se lucrează. Acţiunea lui se declanşează la pornirea sistemului - apăsarea butonului **<Reset>** sau **<Power>**. La acest moment, intră în lucru un mic program din memoria ROM, aşa-numitul *cod startup din BIOS*. Acesta ştie să citească primii 512 octeţi de pe disc, să îi depună undeva în memoria RAM şi să lanseze în execuţie secvenţa de octeţi citită. Evident, în primii 512 octeţi de pe disc va fi un program, *bootprogramul* sau *programul de pornire a încărcării sistemului de operare*. Principala sarcină a *bootprogramului* este aceea de a încărca în RAM partea din *kernel*-ul Unix (sau a altui sistem de operare) special destinată încărcării comple a sistemului de operare.

Nu este obligatoriu ca întregul program de boot să se găsească în blocul 0, ci doar partea care odată executată să permită sistemului lucrul cu alte tipuri de memorie. Odată acest punct atins, execuţia programului poate continua cu părţi care se găsesc pe alte medii de stocare nevolatile, cum ar fi un harddisk, o dischetă, un CR-ROM sau chiar pe un server de boot în cazul staţiilor fără harddisk.

În majoritatea cazurilor în care nu este vorba de o primă instalare a sistemului, restul programului de boot se află stocat pe harddisk într-o zonă specială a acestuia numită *partiţie* sau *segment de boot*. Această partiţie are o structură de date extrem de simplă (mult mai simplă decât a sistemului de fişiere) şi poate fi accesată foarte uşor.

Blocul 1 este numit şi *superbloc*. În el sunt trecute o serie de informaţii prin care se defineşte sistemul de fişiere de pe disc. Printre aceste informaţii amintim:

- numărul  $n$  de *inoduri* (detaliem imediat);
- numărul de zone definite pe disc;
- pointeri spre harta de biţi a alocării inodurilor;
- pointeri spre harta de biţi a spaţiului liber disc;
- dimensiunile zonelor disc, etc.

Blocurile 2 la  $n$ , unde  $n$  este o constantă a formatării discului, se constituie în zona de *inoduri*. Un *inod* (sau *i-nod*) este numele, în terminologia Unix, a *descriptorului* unui fişier. Inodurile sunt memorate pe disc sub forma unei liste (numită *i-listă*). Numărul de ordine al unui inod în cadrul *i-listei* se reprezintă pe doi octeţi şi se numeşte *i-număr*. Acest *i-număr* constituie legătura dintre fişier şi programele utilizator.

Blocurile  $n+1$  la  $n+m$  reprezintă *zona fişierelor*. Este partea cea mai mare din cadrul discului. Alocarea spaţiului pentru fişiere se face printr-o variantă elegantă de indexare. Informaţiile de plecare pentru alocare sunt fixate în inoduri. La discurile Unix actuale există, de regulă, mai multe zone de inoduri intercalate cu mai multe zone de fişiere.

#### 4.2.2 Directori şi inoduri

Structura unei intrări într-un fişier director este ilustrată în fig. 4.5.

Numele fişierului (practic oricât de lung)	inumăr
--	--------

**Figura 4.5 Structura unei intrări în director**

Deci, în director se află numele fişierului şi referinţa spre inodul descriptor al fişierului.

Un *inod* are, de regulă, între 64 şi 128 de octeţi şi el conţine informaţiile din tabelul următor:

mode	Drepturile de acces şi tipul fişierului.
link count	Numărul de directoare care conţin referiri la acest inumăr, adică numărul de legături spre acest fişier.
user ID	Numărul (UID) de identificare a proprietarului.
group ID	Numărul (GID) de identificare a grupului.
Size	Numărul de octeţi (lungimea) fişierului.
access time	Momentul ultimului acces la fişier.
mod time	Momentul ultimei modificări a fişierului.
inode time	Momentul ultimei modificări a structurii inodului.
block list	Lista adreselor disc pentru primele blocuri care aparţin fişierului.
indirect list	Referinţe către celelalte blocuri care aparţin fişierului.

#### 4.2.3 Schema de alocare a blocurilor disc pentru un fişier

Fiecare sistem de fişiere Unix are câteva constante proprii, printre care amintim:

- lungimea unui inod,
- lungimea unui bloc,
- lungimea unei adrese disc (implicit câte adrese disc încap într-un bloc),
- câte adrese de prime blocuri se înregistrează direct în inod,
- câte referinţe se trec în lista de referinţe indirecte.

Indiferent de valorile acestor constante, principiile de înregistrare / regăsire sunt aceleaşi şi le vom prezenta în cele ce urmează. Pentru fixarea ideilor, vom alege aceste constante cu valorile întâlnite mai frecvent la sistemele de fişiere deja consacrate. Cu aceste constante, în fig. 4.6 este prezentată structura pointerilor spre blocurile ataşate unui fişier Unix. Aceste constante sunt:

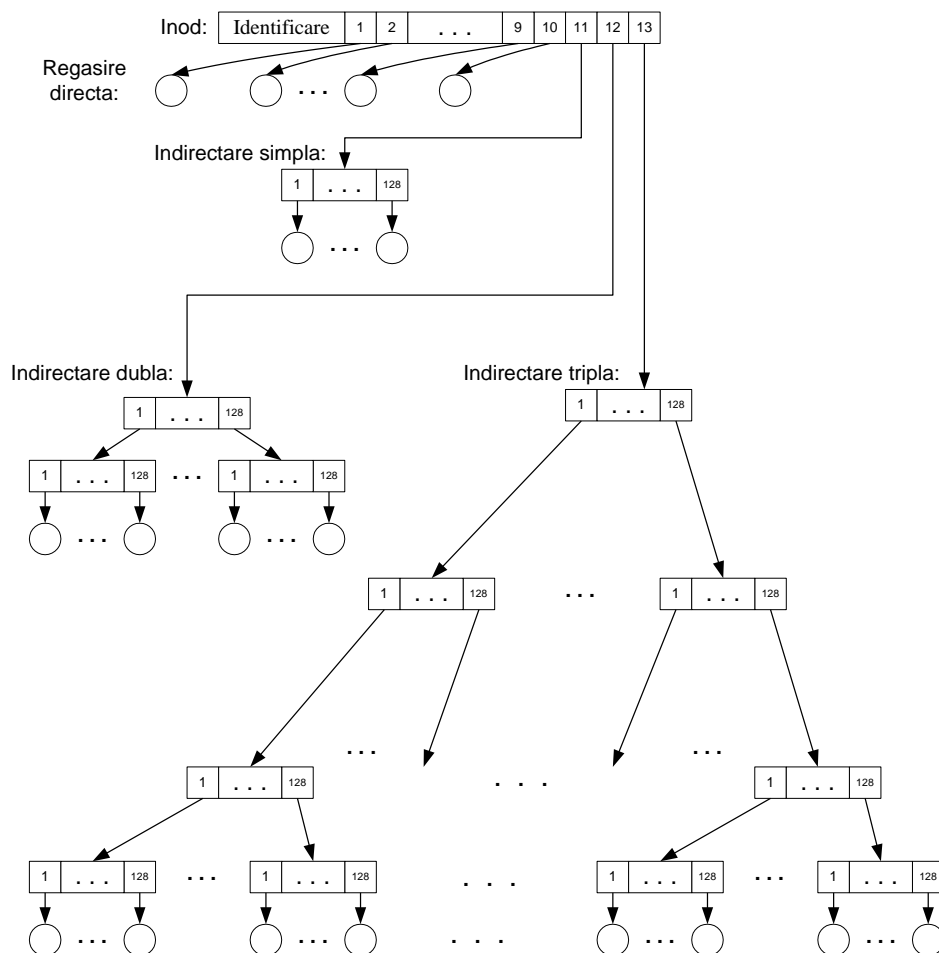
- un inod se reprezintă pe 64 octeţi,
- un bloc are lungimea de 512 octeţi,
- adresă disc se reprezintă pe 4 octeţi, deci încap 128 adrese disc într-un bloc,
- în inod trec direct primele 10 adrese de blocuri,
- lista de adrese indirecte are 3 elemente.

În inodul fişierului se află o listă cu 13 intrări, care desemnează blocurile fizice aparţinând fişierului.

- Primele 10 intrări conţin *adresele primelor* 10 blocuri de câte 512 octeţi care aparţin fişierului.
- Intrarea nr. 11 conţine adresa unui bloc, numit *bloc de indirectare simplă*. El conţine adresele următoarelor 128 blocuri de câte 512 octeţi, care aparţin fişierului.
- Intrarea nr. 12 conţine adresa unui bloc, numit *bloc de indirectare dublă*. El conţine adresele a 128 blocuri de indirectare simplă, care la rândul lor conţin, fiecare, adresele a câte 128 blocuri, de 512 octeţi fiecare, cu informaţii aparţinând fişierului.



- Intrarea nr. 13 conține adresa unui bloc, numit *bloc de indirectare triplă*. În acest bloc sunt conținute adresele a 128 blocuri de indirectare dublă, fiecare dintre acestea conținând adresele a câte 128 blocuri de indirectare simplă, iar fiecare dintre acestea conține adresele a câte 128 blocuri, de câte 512 octeți, cu informații ale fișierului.



**Figura 4.6 Structura unui inod și accesul la blocurile unui fișier**

În fig. 4.6 am ilustrat prin cercuri blocurile de informație care aparțin fișierului, iar prin dreptunghiuri blocurile de referințe, în interiorul acestora marcând referințele.

Numărul de accese necesare pentru a obține direct un octet oarecare este cel mult 4. Pentru fișiere mici acest număr este și mai mic. Atât timp cât fișierul este deschis, inodul lui este prezent în memoria internă. Tabelul următor dă numărul maxim de accese la disc pentru a obține, în acces direct orice octet dintr-un fișier, în funcție de lungimea fișierului.

Lungime maximă (blocuri)	Lungime maximă (octeți)	Accese indirecte	Accese la informație	Total accese
10	5120	–	1	1
$10+128 = 138$	70656	1	1	2
$10+128+128^2 = 16522$	8459264	2	1	3
$10+128+128^2+128^3 = 2113674$	1082201088	3	1	4

La sistemele Unix actuale lungimea unui bloc este de 4096 octeți care poate înregistra 1024 adrese, iar în inod se înregistrează direct adresele primelor 12 blocuri. În aceste condiții, tabelul de mai sus se transformă în:

Lungime maximă (blocuri)	Lungime maximă (octeți)	Accese indirecte	Accese la informație	Total accese
12	49152	–	1	1
$12+1024 = 1036$	4243456	1	1	2
$12++1024+1024^2 = 1049612$	4299210752	2	1	3
$12+1024+1024^2+1024^3 = 1073741824$	4398046511104 (peste 5000Go)	3	1	4

Practic, orice fișier actual, indiferent de mărimea lui, poate fi reprezentat printr-o astfel de schemă.

#### 4.2.4 Accesul proceselor la fișiere

Unix privește conceptul de fișier într-un sens ceva mai larg decât o fac alte sisteme de operare. Așa cum am mai arătat mai sus, există opt tipuri de fișiere:

- normale,
- directori,
- legături hard,
- legături simbolice,
- FIFO, (pipe cu nume),
- socketuri,
- periferice caracter,
- periferice bloc.

Pe lângă acestea, nucleul mai gestionează, într-o sintaxă similară fișierelor, următoarele patru tipuri de comunicații între procese:

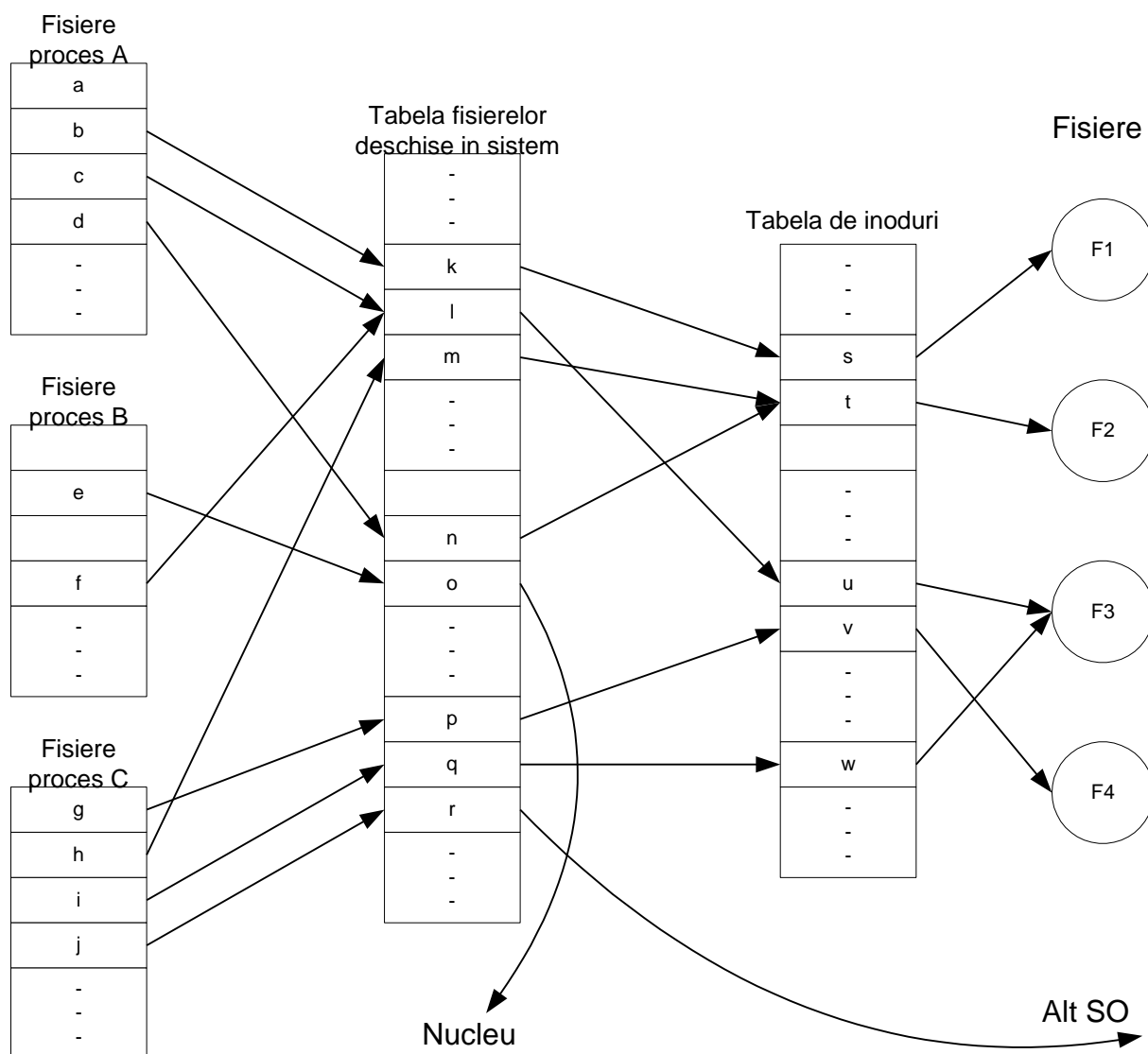
- pipe anonime (a se deosebi de FIFO - pipe cu nume);
- segmente de memorie partajată;
- cozi de mesaje;
- semafoare.

Suporturile fizice pentru aceste 12 tipuri de fișiere sunt, în ultimă instanță:

- *Zona fișierelor pe disc*, pentru fișierele normale, directori, legături hard și simbolice, FIFO și socket din familia Unix.
- *Perifericul respectiv* pentru perifericele caracter și bloc.
- *Zone rezervate de nucleu în memoria internă*, pentru pipe, memorie partajată, cozi de mesaje și semafoare.
- *Interfața de comunicație prin rețea*, pentru socket din familia Internet.

Pentru a putea asigura o tratare uniformă, traseul accesului unui proces la un fișier trece prin mai multe nivele: proces, sistem, inod, fișier, așa cum se vede în fig. 4.7. În fig. 4.7 prezentăm un exemplu în care există trei procese: **A**, **B**, **C** și patru fișiere **F1**, **F2**, **F3**, **F4**. Intrările

pentru legături la diversele nivele le-am notat cu litere mici **a – w**. În cele ce urmează descriem cele patru nivele de legătură.



**Figura 4.7 Corespondența Unix între procese și fişiere**

*Nivelul proces*. Fiecare proces își întreține o tabelă proprie în care înregistrează toate fişierele lui deschise. În fig. 4.7 am notat cu **a–j** câteva intrări de pe acest nivel.

*Nivelul sistem* întreține o tabelă unică cu toate fişierele deschise de către toate procesele din sistem. În fig. 4.7 am notat prin **k–r** câteva intrări din această tabelă.

*Nivelul inod* este de fapt zona (zonele) de inoduri de pe disc. Pentru fişierele deschise, se păstrează în memoria internă copii ale inodurilor corespunzătoare. În fig. 4.7 am notat prin **s–w** câteva astfel de intrări.

*Nivelul fişier* este reprezentat de blocurile disc ce aparțin fişierului. În fig. 4.7 am notat **F1–F4** astfel de fişiere.

Tabela de fişiere la nivel proces are intrările numerotate începând de la 0. Primele trei intrări sunt rezervate astfel:

- intrarea 0 este rezervată intrării standard a procesului (vezi în fig. 4.7 intrările **a** din procesul **A** şi **g** din procesul **C**);
- intrarea 1 este rezervată ieşirii standard (vezi intrările **b**, **e**, **h** din fig. 4.7);
- intrarea 2 este rezervată fişierului standard de erori (unde sistemul afişează mesajele de eroare pentru proces - vezi intrările **c**, **i** din fig. 4.7).

După cum se va vedea în secţiunile următoare, toate apelurile sistem de lucru cu fişiere folosesc pentru identificarea fişierului un număr întreg numit *handle* sau *descriptor de fişier*. Acest întreg este chiar indexul intrării fişierului în tabela procesului respectiv.

În gestiunea fişierelor deschise pe aceste patru nivele, marea majoritate a fişierelor au exact câte o singură intrare pe fiecare nivel, corespunzătoare fişierului respectiv. De exemplu, procesul **A** din fig. 4.7 vede fişierul **F1** prin intermediul intrărilor **b**, **k**, **s** din cele trei tabele. De asemenea, procesul **C** vede fişierul **F4** prin intermediul intrărilor **g**, **p**, **v**.

În Unix este posibil ca mai multe procese să deschidă, în acelaşi timp, un acelaşi fişier - *multiacces la fişiere*. Acest lucru este posibil prin faptul că două intrări de fişiere din două procese pot puncta spre aceeaşi intrare din tabela sistem. În fig. 4.7 intrarea **c** de la procesul **A** şi intrarea **f** de la procesul **B** folosesc în comun acelaşi fişier, cel localizat de intrarea **l** din tabela sistem.

De regulă, fiecare intrare din tabela sistem punctează spre un inod, iar intrări diferite punctează spre inoduri diferite. Aşa sunt, de exemplu, legăturile **k-s**, **l-u**, **p-v**, **q-w**.

Sunt interesante abaterile de la regula de mai sus. Astfel, spre exemplu avem legăturile **m-t** şi **n-t**. Această corespondenţă este posibilă în cazul în care cel puţin una dintre intrările **m** sau **n** se referă la o *legătură simbolică*.

Legăturile din **o** şi din **r** nu punctează spre nici un inod. Legătura **o** este un canal de tip pipe, memorie partajată, coadă de mesaje sau semafor, acestea fiind găzduite în nucleu. Legătura **r** este un socket, prin care se realizează legătura prin reţea cu un alt sistem de operare.

În fine, în marea majoritate a cazurilor există o corespondenţă biunivocă între inoduri şi fişierele corespunzătoare. În fig. 4.7 avem corespondenţele **s-F1**, **t-F2**, **v-F4**.

Abaterea de la această regulă este făcută doar de *legăturile hard*. Fiecare astfel de legătură creează un nou inod pentru acelaşi fişier. În fig. 4.7 avem corespondenţele **u-F3** şi **w-F3**, cel puţin una dintre **u** şi **w** fiind o legătură hard.

### 4.3 Apeluri sistem pentru lucrul cu fişiere

#### 4.3.1 Operaţii I/O

Există două posibilităţi de efectuare a operaţiilor I/O din programe C:

- Prin funcțiile standard C (`fopen`, `fclose`, `fgets`, `fprintf`, `fread`, `fwrite`, `fseek`, etc.) existente în bibliotecile standard C; prototipurile acestora se află în fișierul header `<stdio.h>` (*nivelul superior de prelucrare al fișierelor*).
- Prin funcții standardizate POSIX (`open`, `close`, `read`, `write`, `lseek`, `dup`, `dup2`, `fcntl`, etc.) care reprezintă puncte de intrare în nucleul Unix și ale căror prototipuri se află de regulă în fișierul header `<unistd.h>`, dar uneori se pot afla și în `<sys/types.h>`, `<sys/stat.h>` sau `<fcntl.h>` (*nivelul inferior de prelucrare al fișierelor*).

Prima categorie de funcții o presupunem cunoscută deoarece face parte din standardul C (ANSI). Funcțiile din această categorie reperează orice fișier printr-o structură `FILE *`, pe care o vom numi *descriptor de fișier*.

Funcțiile din a doua categorie constituie apeluri sistem Unix pentru lucrul cu fișiere și fac obiectul secțiunii care urmează. Ele (antetul lor) sunt cuprinse în standardul POSIX. Funcțiile din această categorie reperează orice fișier printr-un întreg nenegativ, numit *handle*, dar atunci când confuzia nu este posibilă îl vom numi tot *descriptor de fișier*. Fiecare astfel de descriptor este index în tabela de fișiere deschise a procesului. De exemplu, să considerăm procesul **C** din fig. 4.7. Aici, aceste handle sau descriptori sunt indecși cu valorile **0**, **1**, **2**, **3** ... care punctează la intrările **g**, **h**, **i**, **j**, ... din tabela de fișiere a procesului **C**.

#### 4.3.2 Apelul sistem *open*

Prototipul funcției sistem este:

```
int open (char *nume, int flag [, unsigned int drepturi ]);
```

Funcția `open` întoarce un întreg - *handle* sau *descriptor de fișier*, folosit ca prim argument de către celelalte funcții POSIX de acces la fișier. În caz de eșec `open` întoarce valoarea -1 și poziționează corespunzător variabila `errno`. În cele ce urmează vom numi `descr` acest număr.

`nume` - specifică printr-un string C, calea și numele fișierului în conformitate cu standardul Unix.

Modul de deschidere este precizat de parametrul de deschidere `flag`. Principalele lui valori posibile sunt:

- `O_RDONLY` deschide fișierul numai pentru citire
- `O_WRONLY` deschide fișierul numai pentru scriere
- `O_RDWR` deschide fișierul atât pentru citire și pentru scriere
- `O_APPEND` deschide pentru adăugarea - scrierea la sfârșitul fișierului
- `O_CREAT` creează un fișier nou dacă acesta nu există, sau nu are efect dacă fișierul deja există; următoarele două constante completează crearea unui fișier
- `O_TRUNC` asociat cu `O_CREAT` (și exclus `O_EXCL` vezi mai jos) indică crearea necondiționată, indiferent dacă fișierul există sau nu
- `O_EXCL` asociat cu `O_CREAT` (și exclus `O_TRUNC`), în cazul în care fișierul există deja, `open` eșuează și semnalează eroare

- `O_NDELAY` este valabil doar pentru fişiere de tip pipe sau FIFO şi vom reveni asupra lui când vom vorbi despre pipe şi FIFO.

În cazul în care se folosesc mai multe constante, acestea se leagă prin operatorul `|`, ca de exemplu: `O_CREAT | O_TRUNC | O_WRONLY`.

Parametrul `drepturi` este necesar doar la crearea fişierului şi indică drepturile de acces la fişier (prin cei 9 biţi de protecţie) şi acţionează în concordanţă cu specificarea `umask`.

### 4.3.3 *Apelul sistem close*

Inchiderea unui fişier se face prin apelul sistem:

```
int close (int descr);
```

Parametrul `descr` este cel întors de apelul `open` cu care s-a deschis fişierul. Funcţia întoarce 0 la succes sau -1 în caz de eşec.

### 4.3.4 *Apelurile sistem read şi write*

Acestea sunt cele mai importante funcţii sistem de acces la conţinutul fişierului. Prototipurile lor sunt:

```
int read (int descr, void *mem, unsigned int noct);  
int write (int descr, const void *mem, unsigned int noct);
```

Efectul lor este acela de a citi (scrie) din (în) fişierul indicat de `descr` un număr de `noct` octeţi, depunând (luând) informaţiile în (din) zona de memorie aflată la adresa `mem`. În majoritatea cazurilor de utilizare, `mem` referă un şir de caractere (`char* mem`).

Cele două funcţii întorc numărul de octeţi efectiv transferaţi între memorie şi suportul fizic al fişierului. De regulă, acest număr coincide cu `noct`, dar sunt situaţii în care se transferă mai puţin de `noct` octeţi.

Aceste două operaţii sunt atomice - indivizibile şi neîntreruptibile de către alte procese. Deci, dacă un proces a început un transfer între memorie şi suportul fişierului, atunci nici un alt proces nu mai are acces la fişier până la terminarea operaţiei `read` sau `write` curente.

Operaţia se consideră terminată dacă s-a transferat cel puţin un octet, dar nu mai mult de maximum dintre `noct` şi ceea ce permite suportul fişierului. Astfel, dacă în fişier există doar `n` < `noct` octeţi rămaşi necitiţi atunci se vor transfera în memorie doar `n` octeţi şi `read` va întoarce valoarea `n`. Dacă în suportul fişierului există cel mult `n` < `noct` octeţi disponibili, atunci se depun din memorie în fişier doar `n` octeţi şi `write` va întoarce valoarea `n`. În ambele situaţii, pointerul curent al fişierului avansează cu numărul de octeţi transferaţi. Dacă la lansarea unui `read` nu mai există nici un octet necitit din fişier - s-a întâlnit marcajul de sfârşit de fişier, atunci funcţia întoarce valoarea 0.

Dacă operația de citire sau de scriere nu se poate termina - a apărut o eroare în timpul transferului - funcția întoarce valoarea -1 și se poziționează corespunzător variabila `errno`.

Ca exemplu de folosire a apelurilor sistem `open`, `close`, `read` și `write` vom scrie un program care face același lucru ca și comanda shell `cp`, adică copiază conținutul unui fișier într-altul. Cele două fișiere se vor da ca parametri în linia de comandă. Dacă al doilea fișier (fișierul destinație) nu există, el se va crea, iar dacă există deja, el se va suprascrie. Sursa programului este prezentată în fig. 4.8.

```
#include <unistd.h>
#include <sys/types.h>
#include <sys/stat.h>
#include <fcntl.h>
#include <stdio.h>

main(int argc, char* argv[]) {
    int fd_sursa, fd_dest, n, i;
    char buf[100], *p;
    if (argc!=3) {
        fprintf(stderr, "Eroare: trebuie dati 2 parametrii.\n");
        exit(1);
    }
    //deschidem primul fișier în modul read-only
    fd_sursa = open(argv[1], O_RDONLY);
    if (fd_sursa<0) {
        fprintf(stderr, "Eroare: %s nu poate fi deschis.\n", argv[1]);
        exit(1);
    }
    /* deschidem al doilea fișier în modul write-only. Dacă el nu
     * există, se va crea, sau dacă există va fi trunchiat. 0755
     * specifică drepturile de acces ale fișierului nou creat
     * (citire+scriere+executie pentru proprietar, citire+executie
     * pentru grup și alții).
     */
    fd_dest = open(argv[2], O_WRONLY | O_CREAT | O_TRUNC, 0755);
    if (fd_dest<0) {
        fprintf(stderr, "Eroare: fisierul %s nu poate fi deschis.\n",
            argv[2]);
        exit(1);
    }
    //citim din fișierul fd_sursa bucăți de maxim 100 octeți și le
    //scriem în fișierul fd_dest, până când nu mai avem ce citi.
    for ( ; ; ) {
        n=read(fd_sursa, buf, sizeof(buf));
        if ((n == 0) break; // S-a terminat de citit fișierul
        p = buf;
        for( ; n > 0 ; ) { // Poate nu se poate scrie odata n octeți
            i = write(fd_dest, p, n);
            if (i == n) break;
            p += i;
            n -= i;
        }
    }
    //închidem cele două fișiere
    close(fd_sursa);
    close(fd_dest);
}
```

**Figura 4.8 Copierea unui fișier cu apelurile sistem `read` și `write`**

### 4.3.5 Apelul sistem lseek

Funcția sistem `lseek` facilitează accesul direct la orice octet din fișier. Evident, pentru aceasta trebuie ca suportul fișierului să fie unul adresabil. Prototipul acestei funcții sistem este:

```
long lseek (int descr, long noct, int deUnde);
```

Se modifică pointerul curent în fișierul indicat de `descr` cu un număr de `noct` octeți. Punctul de unde începe numărarea celor `noct` octeți este indicat de către valoarea parametrului `deUnde`, astfel:

- de la începutul fișierului, dacă are valoarea `SEEK_SET` (valoarea 0)
- de la poziția curentă dacă are valoarea `SEEK_CUR` (valoarea 1)
- de la sfârșitul fișierului dacă are valoarea `SEEK_END` (valoarea 2)

### 4.3.6 Apelurile sistem dup și dup2

Aceste două funcții sistem permit ca un același fișier să fie accesibil prin doi descriptori diferiți. Presupunem că `descrvechi` este o intrare în tabela de fișiere deschise a unui proces, care punctează, așa cum am văzut în fig. 4.7 din 4.2.4, spre o intrare în tabela de fișiere deschise a sistemului. În urma apelului sistem, prin `dup` sau `dup2` se ocupă o nouă intrare `descrnou` din tabela de fișiere a procesului, care va puncta spre același fișier în tabela de fișiere deschise a sistemului. Această duplicare are loc în următoarele condiții:

- `descrvechi` și `descrnou` referă același fișier fizic
- ambele păstrează modul de acces la fișier stabilit la deschidere
- ambii descriptori partajează același pointer curent în fișier

Prototipurile celor două apeluri sistem sunt:

```
int dup (int descrvechi);  
int dup2 (int descrvechi, int descrnou);
```

Apelul sistem `dup` face o copie a `descrvechi` în primul (cel mai mic număr) descriptor liber din tabela de fișiere a procesului.

Apelul sistem `dup2` face o copie a `descrvechi` în `descrnou`, închizând, dacă este cazul, fișierul către care puncta înainte `descrnou`.

Ambele apeluri întorc noul descriptor (`descrnou`) care duplică accesul la fișierul reperat prin `descrvechi`. În caz de eșec, ambele întorc -1 și poziționează corespunzător variabila `errno`.

De exemplu, dacă se dorește ca dintr-un program C mesajele de eroare să fie trecute într-un fișier de pe disc numit "ERORI", se poate proceda ca în programul din fig. 4.9.

```
#include <stdio.h>  
#include <unistd.h>  
#include <fcntl.h>  
  
main () {  
    int descrvechi, descrnou;
```



```

descrvechi = open("ERORI", O_CREAT|O_WRONLY, 0755);
if (descrvechi < 0) {
    fprintf(stderr, "Pe terminal prin stderr: open ERORI imposibil\n");
    fprintf(stdout, "Pe terminal prin stdout: open ERORI imposibil\n");
    exit(1);
}
descrnou = dup2(descrvechi, 2);
if (descrnou != 2) {
    fprintf(stderr, "Pe terminal prin stderr: dup2 imposibil\n");
    fprintf(stdout, "Pe terminal prin stdout: dup2 imposibil\n");
    exit(1);
}
fprintf(stderr, "Mesaj in ERORI prin stderr: dup2 REUSIT\n");
fprintf(stdout, "Mesaj pe terminal prin stdout: dup2 REUSIT\n");
} //main

```

**Figura 4.9 Sursa testdup2.c**

Presupunem că programul se lansează fără nici o redirectare pentru I/O standard. Partea esențială a programului este `descrnou = dup2(descrvechi, 2);`. Prin aceasta, se închide automat fișierul cu descriptorul 2 care referă fișierul standard `stderr` și noul fișier standard de erori va fi fișierul `ERORI`. Dacă deschiderea sau `dup2` nu pot fi executate, atunci pe terminal va apare un mesaj în dublu exemplar: trimis prin `stderr` și prin `stdout`. În caz de succes, pe terminal va apare o linie trimisă prin `stdout` iar în fișierul `ERORI` linia trimisă prin `stderr`.

### 4.3.7 Apelul sistem *fcntl*

Această funcție sistem furnizează sau schimbă proprietăți ale unui fișier deschis indicat printr-un descriptor. Această funcție poate fi apelată prin unul dintre următoarele trei prototipuri:

```

int fcntl(int fd, int cmd);
int fcntl(int fd, int cmd, long arg);
int fcntl(int fd, int cmd, struct flock *lock);

```

Primul parametru, `fd`, referă fișierul deschis asupra căruia se dorește aplicarea controlului.

Parametrul `cmd` este o constantă care specifică operația dorită și în funcție de ea mai urmează sau nu încă un argument. Vom prezenta doar o parte dintre posibilitățile oferite de `fcntl`, urmând ca în secțiunile următoare să revenim cu prezentarea de noi facilități:

- `fcntl(fd, F_DUPFD, arg)` are aproximativ același efect ca și `dup(fd)`, doar că `fcntl` copiază descriptorul `fd` în cel mai mic descriptor, mai mare sau egal cu `arg`.
- `fcntl(fd, F_SETFL, arg)` modifică flagurile de tratare a fișierului față de cum au fost fixate prin `open` sau printr-un `fcntl` precedent. Parametrul `arg` are aceleași valori ca și parametrul `flag` de la `open`: una sau mai multe constante folosite pentru a specifica modul de deschidere a fișierului.
- `fcntl(fd, F_GETFL)` întoarce o configurație de biți reprezentând reuniunea valorilor curente ale flagurilor de tratare a fișierului, așa cum au fost ele fixate prin `open` sau printr-un `fcntl( ... F_SETFL ...)`.

Ce-a de a treia formă va fi tratată în secțiunea 4.6.3.

## 4.4 Gestiunea fişierelor

Unix permite efectuarea din C a principalelor operaţii asupra sistemului de fişiere, oferind în acest sens câteva apeluri sistem. Ele sunt echivalente cu comenzile Shell care efectuează aceleaşi operaţii. Aceste operaţii sunt definite prin specificaţii POSIX corespunzătoare.

De regulă, aceste funcţii au prototipurile într-unul dintre fişierele header:

```
#include <sys/stat.h>
#include <sys/types.h>
#include <unistd.h>
#include <fcntl.h>
```

### 4.4.1 Manevrarea fişierelor în sistemul de fişiere

Iată prototipurile celor mai importante dintre aceste apeluri sistem:

```
int chdir (const char *nume);
char *getcwd(char *mem, int dimensiune);
int mkdir (const char *nume, unsigned int drepturi);
int rmdir (const char *nume);
int unlink(const char *nume);
int link(const char *numevechi, const char *numenou);
int symlink(const char *numevechi, const char *numenou);
int chmod (const char *nume, unsigned int drepturi);
int stat (const char *nume, struct stat *stare);
int mknod(const char *nume, unsigned int mod, dev_t dev);
int chown(const char *nume, unsigned int proprietar,
          unsigned int grup);
int access(const char *nume, int permisiuni);
int rename(const char *numevechi, const char *numenou);
```

- `chdir` schimbă directorul curent în cel specificat prin `nume`. Pe sistemele BSD şi Unix System V release 4 mai este prezent şi apelul `fchdir` care face acelaşi lucru ca şi `chdir`, doar că directorul este specificat printr-un descriptor de fişier deschis, nu prin `nume`.
- `getcwd` copiază în zona de memorie indicată de `mem`, de lungime `dimensiune` octeţi, calea absolută a directorului curent. Dacă calea absolută a directorului are lungimea mai mare decât `dimensiune`, se returnează `NULL` şi `errno` primeşte valoarea `ERANGE`.
- `mkdir` creează un nou director, având calea şi numele specificate prin `nume` şi drepturile indicate prin întregul `drepturi`, din care se reţin numai primii 9 biţi.
- `rmdir` şterge directorul specificat prin `nume` (acest director trebuie să fie gol).
- `unlink` şterge fişierul specificat prin `nume`.
- `link` creează o legătură hard cu numele `numenou` spre un fişier existent, `numevechi`. Numele nou se poate folosi în locul celui vechi în toate operaţiile şi nu se poate spune care dintre cele două nume este cel iniţial.
- `symlink` creează o legătură simbolică (soft) cu numele `numenou` spre un fişier existent, `numevechi`. O legătură soft este doar o referinţă la un fişier existent sau

inexistent. Dacă șterg o legătură soft se șterge doar legătura, nu și fișierul original, pe când dacă șterg o legătură hard, se șterge fișierul original.

- `chmod` atribuie drepturile de acces `drepturi` la fișierul specificat prin `nume`. Pe sistemele BSD și Unix System V release 4 mai este prezent și apelul `fchmod` care face același lucru ca și `chmod`, doar că fișierul este specificat printr-un descriptor de fișier deschis, nu prin `nume`.
- `stat` depune la adresa `stare` informații privind fișierul specificat prin `nume` (inod-ul fișierului, drepturile de acces, id-ul proprietarului, id-ul grupului, lungimea în octeți, numărul de blocuri ale fișierului, data ultimului acces la fișier, etc. – vezi exemplul de mai jos). Sistemele BSD și Unix System V release 4 au și apelul `fstat` care face același lucru ca și `stat`, doar că fișierul este specificat printr-un descriptor de fișier deschis, nu prin `nume`.
- `mknod` creează un fișier simplu sau un fișier special desemnat prin `nume`. Parametrul `mod` specifică printr-o combinație de constante simbolice legate prin simbolul `'|'` atât drepturile de acces la fișierul nou creat, cât și tipul fișierului care poate fi unul dintre următoarele:
  - `S_IFREG` (fișier normal)
  - `S_IFCHR` (fișier special de tip caracter)
  - `S_IFBLK` (fișier special de tip bloc)
  - `S_IFIFO` (*pipe* cu nume sau FIFO – vezi capitolul 5)
  - `S_IFSOCK` (Unix domain socket – folosit pentru comunicarea locală între procese)

Dacă tipul fișierului este `S_IFCHR` sau `S_IFBLK`, atunci `dev` conține numărul minor și numărul major al fișierului special nou creat; altfel, acest parametru se ignoră.

- `chown` schimbă proprietarul și grupul din care face parte proprietarul unui fișier specificat prin `nume`. Noul proprietar al fișierului va fi cel indicat de parametrul `proprietar`, iar noul grup al fișierului va fi cel indicat de parametrul `grup`. Pe sistemele BSD și Unix System V release 4 mai este prezent și apelul `fchown` care face același lucru ca și `chown`, doar că fișierul este specificat printr-un descriptor de fișier deschis, nu prin `nume`.
- `access` verifică dacă procesul curent are dreptul specificat de `permisiuni` relativ la fișierul specificat prin `nume`. `permisiuni` va conține una sau mai multe valori legate prin `'|'` dintre următoarele: `R_OK` - citire, `W_OK` - scriere, `X_OK` – execuție și `F_OK` – existență fișier. Verificarea se face cu UID-ul și GID-ul reale ale procesului, nu cele efective (vezi capitolul 5).
- `rename` redenumeste fișierul specificat de `numevechi` în `numenou`.

#### 4.4.2 *Creat, truncate, readdir*

Următoarele trei apeluri nu au corespondent direct printre comenzile Shell:

```
int creat(const char *pathname, unsigned int mod);
int truncate(const char *nume, long int lung);

#include <dirent.h>
struct dirent *readdir(DIR *dir);
```

Apelul sistem `creat` este echivalent cu următorul apel sistem:

```
open (const char *nume, O_CREAT|O_WRONLY|O_TRUNC);
```

Apelul sistem `truncate` trunchiază fişierul specificat prin nume la exact lung octeţi.

Funcţia `readdir` nu este o funcţie sistem, ci este funcţia POSIX pentru parcurgerea subdirectoarelor şi fişierelor dintr-un director. Ea întoarce într-o structură de tipul `dirent` următorul subdirector sau fişier din directorul dat de parametrul `dir`. Această funcţie încapsulează de fapt apelul sistem `getdents`. Alte funcţii în legătură cu `readdir` şi specificate de standardul POSIX sunt: `opendir`, `closedir`, `rewinddir`, `scandir`, `seekdir` şi `telldir`. Nu intrăm în detalierea acestor funcţii deoarece ele nu sunt funcţii sistem.

Majoritatea apelurilor de mai sus întorc valoarea 0 în caz de succes şi -1 (plus setarea corespunzătoare a variabilei `errno`) în caz de eroare. Excepţie de la această regulă fac apelurile: `getcwd` care returnează NULL în caz de eroare şi setează corespunzător `errno` şi `readdir` care întoarce NULL în caz de eroare.

#### 4.4.3 Un exemplu: obţinerea tipului de fişier prin apelul sistem `stat`

În această secţiune vom da un exemplu de utilizare a lui `stat`. Exemplul se referă la afişarea tipului de fişier pentru fişierele ale căror nume sunt date ca argumente la linia de comandă.

Programul `tipfis.c` din fig. 4.10 exemplifică folosirea apelului sistem `stat` pentru determinarea tipului de fişier recunoscut de către sistem. Tipurile de fişiere le-am prezentat într-o secţiune precedentă.

```
/* tipfis.c:
 * Tipăreşte tipurile fişierelor date în linia de comandă
 */
#include <stdio.h>
#include <sys/types.h>
#include <sys/stat.h>
#include <string.h>

main (int argc, char *argv[]) {
    int i;
    struct stat statbuff;
    char tip[40];
    strcpy(tip, "");
    for (i = 1; i < argc; i++) {
        printf ("%s: ", argv[i]);
        if (stat (argv[i], &statbuff) < 0)
            fprintf (stderr, "Eroare stat");
        switch (statbuff.st_mode & S_IFMT) {
            case S_IFDIR:
                strcpy(tip, "Director");
                break;
            case S_IFCHR:
                strcpy(tip, "Special de tip caracter");
                break;
            case S_IFBLK:
                strcpy(tip, "Special de tip bloc");
                break;
            case S_IFREG:
                strcpy(tip, "Obişnuit");
                break;
            case S_IFLNK:
                break;
        }
        printf ("%s\n", tip);
    }
}
```

```

/* acest test nu va fi adevărat niciodată deoarece stat
 * verifică în cazul unei legături simbolice, fişierul
 * pe care îl referă legătura şi nu legătura în sine.
Il
 * scriem aici pentru completitudine. Ca să verificăm
 * tipul unei legături simbolice putem folosi funcţia
 * lstat asemănătoare cu stat şi disponibilă pe
 * versiunile BSD şi Unix System V.
 */
strcpy(tip, "Legatura simbolica");
break;
case S_IFSOCK:
    strcpy(tip, "Socket");
    break;
case S_IFIFO:
    strcpy(tip, "FIFO");
    break;
default:
    strcpy(tip, "Tip necunoscut");
} //switch
printf ("%s\n", tip);
} //for
} //main

```

**Figura 4.10 Sursa tipfis.c**

## 4.5 Alte apeluri sistem

În acest subcapitol, vom prezenta câteva apeluri sistem care pot fi utile programatorului Unix şi care nu au mai fost prezentate până aici. Înainte de prezentarea fiecărui prototip vom scrie şi directiva `#include` necesară la începutul sursei programului C. În unele cazuri apar două astfel de directive, iar programatorul o va alege, încercând, pe cea care este potrivită pentru varianta lui de Unix.

### 4.5.1 Time

```

#include <time.h>
time_t time(time_t *t);

```

Apelul sistem `time` returnează timpul scurs de la 1 Ianuarie 1970 00:00 UTC, măsurat în număr de secunde. Dacă `t` este nenul, numărul de secunde va fi salvat la adresa referită de `t`.

### 4.5.2 Umask

```

#include <sys/types.h>
#include <sys/stat.h>
unsigned int umask(unsigned int masca);

```

Apelul sistem `umask` setează masca drepturilor pentru fişierele nou create la valoarea dată de `masca` & 0777. Această mască de biţi indică drepturile de acces care nu sunt setate implicit la crearea unui fişier. Valoarea anterioară a lui `umask` este returnată.

### 4.5.3 *Gethostname*

```
#include <unistd.h>
int gethostname(char *name, int lung);
```

Apelul sistem `gethostname` pune la adresa dată de `name`, pe lungime de maxim `lung` octeți, numele mașinii (calculatorului).

### 4.5.4 *Gettimeofday*

```
#include <sys/time.h>
int gettimeofday(struct timeval *tv, 0);
```

Apelul sistem `gettimeofday` pune în structura `tv`, de tip `timeval`, valoarea timpului curent. Structura `timeval` este definită în felul următor:

```
struct timeval {
    long        tv_sec;    /* numărul de secunde */
    long        tv_usec;   /* numărul de microsecunde */
};
```

### 4.5.5 *Mmap și munmap*

```
#include <sys/mman.h>
void* mmap(void *start, size_t lung, int prot,
           int atribut, int fd, off_t offset);
int munmap(void *start, size_t lung);
```

Apelul sistem `mmap` mapează `lung` octeți începând de la deplasamentul `offset` din fișierul (sau alt obiect) indicat de descriptorul de fișier `fd` în memorie, preferabil la adresa `start`. În general, `start` are valoarea 0 și adresa unde are loc maparea este returnată de `mmap`. Apelul sistem `munmap` șterge maparea de la adresa `start`, pe lungime `lung`. Parametrul `prot` specifică protecția de memorie dorită și poate avea valorile:

- `PROT_EXEC` - paginile de memorie pot fi executate
- `PROT_READ` - paginile de memorie pot fi citite
- `PROT_WRITE` - paginile de memorie pot fi modificate
- `PROT_NONE` - paginile de memorie nu pot fi accesate

Parametrul `atribut` specifică tipul de obiect mapat, opțiuni de mapare și dacă modificări ale conținutului memoriei unde a avut loc maparea sunt vizibile numai procesului curent sau și altor procese. Cele mai importante valori ale sale sunt: `MAP_FIXED`, `MAP_SHARED`, `MAP_PRIVATE`.

### 4.5.6 *Fsync și fdatasync*

```
#include <unistd.h>
int fsync(int fd);
int fdatasync(int fd);
```

Apelul sistem `fsync` scrie pe disc datele fișierului `fd` modificate în buffer-ele nucleului, dar nesalvate încă pe disc. Diferența dintre `fsync` și `fdatasync` este că `fdatasync` nu salvează și date de control despre fișier ca timpul ultimului acces la fișier, ci doar conținutul fișierului.

#### 4.5.7 *Uname*

```
#include <sys/utsname.h>
int uname(struct utsname *buf);
```

Apelul sistem `uname` returnează în structura `buf` de tip `utsname` date despre nucleul sistemului de operare. Tipul `utsname` are următoarele câmpuri:

```
struct utsname {
    char sysname[];
    char nodename[];
    char release[];
    char version[];
    char machine[];
#ifdef _GNU_SOURCE
    char domainname[];
#endif
};
```

Aceste câmpuri au aceleași sensuri ca și câmpurile afișare de comanda shell `uname`.

#### 4.5.8 *Select și seturi de descriptori*

```
#include <sys/select.h>
int select(int n, fd_set *readfds, fd_set *writefds,
           fd_set *exceptfds, struct timeval *timeout);
```

Apelul sistem `select` monitorizează într-un interval de timp dat de `timeout` starea descriptorilor de fișiere din seturile `readfds`, `writefds` și `exceptfds` și raportează (returnează) câte dintre ele și-au modificat starea. Descriptorii din `readfds` vor fi monitorizați pentru a vedea dacă au apărut caractere noi disponibile pentru citire, cei din `writefds` pentru a vedea dacă o nouă operație de scriere nu se va bloca, iar cei din `exceptfds` vor fi monitorizați pentru excepții. Parametrul `n` trebuie să aibă ca valoare cel mai mare număr de descriptor din cele trei seturi plus 1.

Standardul POSIX specifică și patru macrouri pentru a manipula cele trei seturi de descriptori. Prototipurile lor sunt date mai jos, unde prin **fd** am notat un descriptor de fișiere, iar prin **set** unul dintre seturile de descriptori din **select**:

```
FD_CLR(int fd, fd_set *set);    //sterge fd din set
FD_ISSET(int fd, fd_set *set);  //este fd in set?
FD_SET(int fd, fd_set *set);    //adauga fd la set
FD_ZERO(fd_set *set);          //goleste set
```

Apelul sistem **select** multiplexează de fapt mai multe canale de intrare-ieșire sincrone.

## 4.6 Blocarea fişierelor

### 4.6.1 Un (contra)exemplu

Una din problemele ce apar frecvent în medii concurente este partajarea resurselor. Să considerăm următoarea situație: se dă un fişier cu numele `secv` care conține pe prima linie un număr întreg de 5 cifre, reprezentate în ASCII. Să considerăm că fiecare proces (deocamdată prin proces vom înțelege program aflat în execuție) face următoarele acțiuni:

- Citește acest număr din `secv`.
- Mărește numărul cu o unitate.
- Rescrie numărul rezultat în `secv`.

Programul `lockfile.c` din fig. 4.11 realizează aceste acțiuni. Funcțiile `my_lock(fd)` și `my_unlock(fd)` au sarcina de a asigura accesul exclusiv la fişierul `secv` pe durata acțiunii. Într-o primă variantă, comentăm apelurile celor două funcții.

```
#include <stdio.h>
#include <unistd.h>
#include <fcntl.h>
#include <string.h>

main () {
    int  fd, i, n, nrsecv;
    char buf[6];
    fd = open ("secv", O_RDWR);
    for (i = 0; i < 10000; i++) {

        // my_lock(fd);          Blocheaza fisierul

        lseek (fd, 0L, 0);      //pozitionare la inceput
        n = read (fd, buf, 5);  //citim numarul
        buf[n] = '\0';          //zeroul terminal
        sscanf (buf, "%d\n", &nrsecv);
                                //convertim stringul citit in numar
        printf ("pid = %d, secventa = %d\n", getpid(), nrsecv);
        nrsecv++;               //incrementeaza secventa
        sprintf (buf, "%5d\n", nrsecv);
                                //convertim numarul in string
        lseek (fd, 0L, 0);      //revenire inainte de scriere
        write (fd, buf, strlen(buf)); //scriem numarul inapoi

        // my_unlock(fd);       Deblocheaza fisierul
    }
}
```

**Figura 4.11 Sursa `lockfile.c`**

Înainte de rularea programului, se va crea cu un editor de texte fişierul `secv` în care se va introduce o singură linie conținând numărul 1. Se compilează programul și să presupunem că fişierul executabil are numele `a.out`. Se rulează apoi programul, simultan în două procese, folosind comanda:

```
$ a.out & a.out &
```



Cele două procese vor afişa, intercalat, linii pe ieşirea standard. Ultima linie afişată va fi ceva de forma:

```
pid = 1265 secventa = 11953
```

Rezultatul pare ciudat, cel puţin la prima vedere. În mod normal, valoarea numărului de secvenţă ar trebui să fie 20000! De ce nu se întâmplă aşa? Deoarece prin comentarea apelului funcţiilor `my_lock` şi `my_unlock` secvenţa dintre ele poate fi executată, total sau parţial, în acelaşi moment de către cele două procese! Frecvent se va întâmpla că ambele vor citi din fişier aceeaşi valoare, o vor incrementa şi o vor scrie. În consecinţă valoarea va fi mărită doar cu o unitate, deşi au acţionat asupra ei două procese!

#### 4.6.2 Tipuri de blocare

Exemplul din secţiunea precedentă ilustrează faptul că în anumite condiţii se impune ca asupra unui fişier sau asupra unei porţiuni din fişier, să aibă dreptul să acţioneze doar un singur proces. Această restricţie este cunoscută sub numele de *blocarea* unui fişier. Sub Unix sunt utilizate două feluri de blocare: conciliantă şi obligatorie.

Blocare conciliantă (*advisory*) este atunci când sistemul ştie care fişiere au fost blocate şi de către cine, iar procesele cooperează, prin funcţii de tip `my_lock` şi `my_unlock`, la accesarea fişierului. Nucleul sistemului de operare nu previne situaţia în care un proces indisciplinat scrie în fişierul blocat.

Blocarea obligatorie (*mandatory*) are loc atunci când sistemul verifică la fiecare scriere şi citire dacă fişierul este blocat sau nu. Pentru a permite blocarea obligatorie, trebuie invalidat dreptul de execuţie al grupului şi bitul set-group-ID să fie 1 pentru fişierul respectiv, iar la montarea sistemului de fişiere (prin `mount`) să se permită blocare obligatorie.

Blocarea unui fişier înseamnă blocarea accesului la orice octet din fişier.

Blocarea unui articol înseamnă blocarea accesului la un număr de octeţi consecutivi din fişier.

Mecanismele de blocare a fişierelor sub Unix specificate de standardul POSIX (şi implementate de biblioteca C de la GNU) permit, în esenţă, două tipuri de blocări, fiecare dintre ele putând fi conciliante sau obligatorii (implicit conciliante) în funcţie de îndeplinirea sau nu a condiţiilor de mai sus:

- blocare *exclusivă*, atunci când un singur proces are acces la fişier sau porţiunea din fişier. De obicei, este vorba aici de o operaţie de scriere în fişier, iar pe durata pregătirii operaţiei şi a scrierii propriu-zise, nici un alt proces nu poate nici scrie, nici citi porţiunea rezervată.
- blocare *partajată*, atunci când porţiunea rezervată poate fi citită, simultan, de către mai multe procese, dar nici un proces nu scrie.

Aceste tipuri de blocare corespund rezolvării unei probleme celebre de programare concurentă - problema cititorilor şi a scriitorilor.

#### 4.6.3 Blocarea conciliantă prin *fcntl*

Nucleele Unix oferă mai multe apeluri sistem destinate blocării. De exemplu, sub Linux și FreeBSD există o funcție sistem `flock` pentru blocare. De asemenea, biblioteca C de la GNU oferă funcțiile `flock` și `lockf`. Toate fac uz de serviciile apelului sistem `fcntl` pentru realizarea blocării. Reamintim cea de-a treia formă a prototipului `fcntl`, folosită pentru blocare:

```
int fcntl(int fd, int cmd, struct flock *lock);
```

Parametrul `fd` este descriptorul fișierului supus blocării.

Parametrul `lock` este un pointer la o structură de tip `flock`, prin care sunt indicați parametrii operației de blocare. Această structură ce conține cel puțin următoarele câmpuri:

```
struct flock {
    - - -
    short l_type;
    short l_whence; /* cum se interpretează l_start */
    off_t l_start; /* offset-ul de început */
    off_t l_len; /* numărul de octeți */
    pid_t l_pid; /* PID-ul procesului care blochează blocajul
                  curent */
    - - -
};
```

Câmpul `l_type` indică tipul operației de blocare:

- `F_WRLCK` indică faptul că este vorba de o blocare exclusivă - blocare la scriere. Un singur proces poate efectua scriere în porțiunea rezervată, iar toate celelalte procese solicitante sunt în așteptare.
- `F_RDLCK` indică faptul că este vorba de o blocare partajată, procesele ce solicită porțiunea rezervată pentru citire își pot efectua operațiile, în timp ce procesele scriitori stau în așteptare.
- `F_UNLCK` indică ridicarea stării de blocare, lăsând eventualele alte procese ce doresc blocarea să o facă.

Câmpul `l_whence` indică locul începând de unde se determină localizarea porțiunii rezervate. Valorile lui sunt aceleași cu cele de la funcția `lseek`:

- `SEEK_SET` indică reperarea porțiunii rezervate față de începutul fișierului.
- `SEEK_CUR` indică reperarea porțiunii rezervate față de poziția curentă a pointerului fișierului.
- `SEEK_END` indică reperarea porțiunii rezervate față de sfârșitul fișierului.

Câmpul `l_start` indică începutul porțiunii rezervate, numărat în octeți față de locul indicat prin `l_whence`. Acest câmp poate avea atât valori pozitive, cât și valori negative. În consecință, începutul zonei se specifică prin combinația parametrilor `l_whence` și `l_start`.

Câmpul `l_len` indică lungimea în octeți a zonei rezervate, un număr nenegativ. Dacă valoarea este strict pozitivă atunci sunt rezervați pentru blocare un număr de `l_len` octeți de la începutul zonei rezervate. Dacă `l_len` are valoarea zero, atunci sunt rezervați toți octeții de la începutul zonei până la sfârșitul fișierului.

Câmpul `l_pid` conține PID-ul procesului care ține blocată regiunea (are sens numai în cazul comenzii `F_GETLK`).

Parametrul `cmd`, în acțiunile de blocare, poate avea următoarele valori:

- `F_SETLKW` indică faptul că procesul cooperează la blocarea fișierului. Dacă fișierul este blocat exclusiv (la scriere) de către un alt proces, atunci procesul curent intră în așteptare până la deblocarea fișierului. Acțiunea de blocare pe care o execută, după eventuala ieșire din starea de așteptare, este cea dictată de parametrul `lock`.
- `F_SETLK` este similar cu `F_SETLKW`, numai că în cazul unei blocări exclusive procesul nu mai intră în așteptare, ci întoarce `-1` și setează `errno` la `EACCES` sau `EAGAIN`.
- `F_GETLK` întoarce la adresa `lock` starea de blocare în care se află fișierul.

În fig. 4.12 sunt prezentate sursele funcțiilor `my_lock` și `my_unlock` în care se folosește blocarea și deblocarea prin `fcntl`.

```
void my_lock(int fd) {
    struct flock lock;
    lock.l_type = F_WRLCK;      //blocat exclusiv (la scriere)
    lock.l_whence = SEEK_SET;   //baza blocarii (inceputul fisierului)
    lock.l_start = 0;           //offset blocare fata de baza
    lock.l_len = 0;             //lungimea blocarii, tot fisierul
    fcntl(fd, F_SETLKW, &lock); //comanda blocarea
} //my_lock

void my_unlock(int fd) {
    struct flock lock;
    lock.l_type = F_UNLCK;      //deblocheaza
    lock.l_whence = SEEK_SET;   //incepand de la inceputul fisierului
    lock.l_start = 0;           // pozitia 0 fata de whence
    lock.l_len = 0;             //tot fisierul
    fcntl(fd, F_SETLKW, &lock); //comanda deblocarea
} //my_unlock
```

**Figura 4.12 funcțiile `my_lock` și `my_unlock` folosind `fcntl`**

#### 4.6.4 Blocare prin `lockf` și `flock`

Prototipurile celor două funcții sunt:

```
int lockf(int fd, int actiune, long lungime);
int flock(int fd, int actiune);
```

Apelul `lockf` realizează, în funcție de condițiile specificate mai sus, fie blocare conciliantă, fie obligatorie (implicit, blocare conciliantă). De asemenea, poate bloca un fișier întreg sau numai o parte din el. Apelul `flock` blochează numai conciliant și numai întregul fișier.

`actiune` pentru `lockf` este una dintre următoarele:

- `F_ULOCK` deblocarea unei regiuni blocate
- `F_LOCK` blocarea unei regiuni
- `F_TLOCK` testarea și blocare dacă nu este deja blocată
- `F_TEST` testarea unei regiuni dacă este blocată sau nu

`actiune` pentru `flock` este una dintre următoarele:

- `LOCK_SH` face blocare partajată (acces din mai multe procese)

- `LOCK_EX` face blocare exclusivă
- `LOCK_UN` face deblocare
- `LOCK_NB` în general, dacă se cere blocarea și fișierul este deja blocat, atunci programul așteaptă la apelul `flock` până când poate realiza blocarea cerută. Legarea acestui atribut, cu *sau* (`()`), de una dintre acțiunile anterioare, face ca în această situație să nu se mai aștepte.

`fd` este descriptorul de fișier.

`lungime` spune câți octeți vor fi blocați în fișier începând cu poziția curentă. Dacă `lungime=0` atunci se blochează întregul fișier.

În caz de erori, funcția întoarce `-1` și poziționează variabila globală sistem `errno`.

Reamintim că dintre cele două funcții, numai `flock` reprezintă o funcție sistem implementată de nucleu (pe care o apelăm prin intermediul funcției POSIX cu același nume implementată de biblioteca C de la GNU – menționăm din nou că funcția `flock` a bibliotecii C de la GNU este o funcție *wrapper* pentru funcția sistem `flock` implementată de nucleu), apelul funcției `lockf` traducându-se de fapt într-un apel sistem `fcntl`.

Folosirea lui `F_TEST` este relativ nesigură, din cauză că secvența:

```
if (lockf (fd, F_TEST, size) == 0)
    rc = lockf (fd, F_LOCK, size)
```

nu este echivalentă cu secvența:

```
rc = lockf (fd, F_TLOCK, size)
```

deoarece între `F_TEST` și `F_LOCK` ar putea interveni un alt proces care să blocheze.

În fig. 4.13 sunt prezentate sursele `my_lock` și `my_unlock` realizate folosind `lockf`.

```
void my_lock (int fd) {
    lseek (fd, 0L, 0);
    lockf (fd, F_LOCK, 0L);
} //my_lock

void my_unlock (int fd) {
    lseek (fd, 0L, 0);
    lockf (fd, F_UNLOCK, 0L);
} //my_unlock
```

**Figura 4.13 funcțiile `my_lock` și `my_unlock` folosind `lockf`**

## 8 Sisteme de operare – prezentare generală

### 8.1 *Tipuri de sisteme de operare (SO); clasificări*

În prezent există un mare număr de SO în funcțiune și numărul acestora este în continuă creștere. Nu există un criteriu unitar de comparare a acestora. Noi vom da mai multe clasificări, în funcție de următoarele criterii:

- după gradul de partajabilitate a resurselor;
- după tipurile de interacțiuni permise;
- după organizarea internă a programelor ce compun SO.

#### 8.1.1 *Clasificare după gradul de partajabilitate a resurselor*

După acest criteriu, distingem trei categorii de sisteme de operare [10]:

- Sisteme monouser și monotasking;
- Sisteme monouser și multitasking;
- Sisteme multiuser.

Sisteme de operare monouser sunt acelea care permit, la un moment dat, unui singur utilizator să folosească sistemul. Un sistem este monotasking dacă admite la un moment dat execuția unui singur program.

În forma cea mai simplă, un sistem monouser și monotasking execută un singur program la un moment dat și acesta rămâne activ din momentul lansării lui și până la terminarea lui completă. Un astfel de sistem este de exemplu cel care deservește un telefon mobil sau un terminal PDA (Personal Digital Assistant).

Un sistem monouser și multitasking este acela în care un singur utilizator este conectat la sistem, dar el poate să își lanseze simultan mai multe procese în lucru. Tipice în acest sens sunt SO din familia Windows 9X, ca și cele din familia NT care nu sunt servere. Tot în această categorie putem îngloba sistemele Linux instalate pe calculatoare neconectate în rețea. Pentru sistemele de acest fel care au memoria internă relativ mică, se aplică așa-numita tehnică swapping, adică evacuarea temporară a unui program, din memoria internă pe discul magnetic. În timpul evacuării, în memoria internă este încărcat un alt program, care la rândul lui este executat parțial și este apoi evacuat și el. Cât timp un program este în memorie, el are acces la toate resursele sistemului. Programele sunt astfel executate pe porțiuni, momentele de execuție alternând cu cele de evacuare. După unii autori ([4]), SO de acest tip se află la granița dintre sistemele monoutilizator și cele multiutilizator.

Sisteme de operare multiutilizator permit accesul simultan al mai multor utilizatori la sistem. Ele trebuie să aibă în vedere partajarea procesorului, a memoriei, a timpului, a perifericelor și a altor tipuri de resurse, între utilizatorii activi la un moment dat. Vom reveni cu detalii asupra metodelor de partajare a procesorului, timpului și a memoriei. La astfel de sisteme sunt necesare tehnici mai sofisticate, de gestiune și protecție a utilizatorilor. La un moment dat aici există mai multe procese active care se execută concurent sub controlul SO.

### 8.1.2 Clasificare după tipurile de interacțiuni permise

În funcție de resursele pe care le are un sistem și în funcție de destinația acestuia, este necesară stabilirea unor strategii de interacțiune cu utilizatorul. Conform acestui criteriu de diferențiere, distingem [34]:

- Sisteme seriale (care prelucrează loturi de lucrări);
- Sisteme cu timp partajat (time – sharing);
- Calculatoare personale și stații de lucru (workstations);
- Sisteme autonome (embedded systems);
- Sisteme portabile, destinate comunicațiilor (telefoane mobile, PDA-uri etc.);
- Sisteme conectate în rețea.

**Sisteme seriale.** În cadrul acestor sisteme lucrările se execută pe loturi pregătite în prealabil. Practic, din momentul predării unei lucrări la ghișeul centrului de calcul și până la eliberarea ei, utilizatorul nu poate interveni spre a influența execuția programului său. SO afectate acestor sisteme pot lucra atât monoutilizator, cât și multiutilizator. Din punct de vedere istoric sunt printre primele sisteme. În prezent, în această categorie intră supercalculatoarele, precum și o serie de sisteme medii – mari cum ar fi AS-400. Sistemul de operare trebuie să gestioneze loturi de lucrări pe care să le pregătească spre a fi introduse în sistem și eventual să pregătească livrarea rezultatelor (listare etc.). În general gradul de concurență la acest tip de sisteme este relativ redus. În [10] sunt date mai multe detalii despre acest tip de sisteme.

**Sistemele cu timp partajat** trebuie să suporte, într-o manieră interactivă, mai mulți utilizatori. Fiecare utilizator este în contact nemijlocit cu programul său. În funcție de unele rezultate intermediare, el poate decide modul de continuare a activității programului său. SO trebuie să gestioneze printre altele și terminalele de teletransmisie la capătul cărora se află utilizatori care-i pot transmite diverse comenzi. Sistemul trebuie să dispună de mecanisme mai sofisticate decât cele seriale pentru partajarea procesorului, a memoriei și a timpului. De regulă, acest tip de sisteme practică partajarea timpului, astfel. Se fixează o cuantă de timp. Pe durata cuantei, un singur proces are acces la procesor (și la celelalte resurse aferente). După epuizarea cuantei, procesul este suspendat, pus la sfârșitul cozii de procese și un alt proces ocupă procesorul pentru următoarea cuantă ș.a.m.d. Dimensiunea cuantei și strategiile de comutare sunt astfel alese încât să se reducă timpii de așteptare și să se realizeze o servire rezonabilă a proceselor. Tipice acestui tip de sisteme sunt sistemele Unix, dar și Windows din familia NT / server.

**Calculatoare personale și stații de lucru.** La acest tip de sisteme, toate resursele mașinii sunt dedicate utilizatorului care este conectat. Așa că nu se pune problema partajării resurselor între utilizatori, ci doar a datelor între procesele pe care userul le lansează simultan. Practic, sarcina principală a SO este aceea de a oferi userului o interfață prietenoasă de exploatare a sistemului.

**Sistemele autonome** sunt sisteme dedicate unui anumit proces industrial, cum ar fi: funcționarea unui robot, coordonarea activităților dintr-un satelit geostaționar, supravegherea unui baraj de apă, a unei stații de radiolocație etc. Fiind conectate la diverse procese tehnologice, aceste sisteme trebuie să fie capabile să deservească în timp prestabilit fiecare serviciu care i se cere. Dacă sistemul nu este capabil să dea un răspuns, este posibilă oprirea procesului supravegheat. Astfel de sisteme sunt în prezent în plină dezvoltare, mai ales din cauza dezvoltării rapide a tehnologiilor multimedia.

Sisteme portabile, destinate comunicațiilor. Acestea formează cea mai nouă categorie de sisteme. Exemplele tipice sunt telefoane mobile și PDA-urile. Sunt mașini portabile, de dimensiuni mici și au facilități puternice de comunicare. Facilitatea lor principală este conexiunea wireless printr-una dintre tehnologiile existente: radio, infraroșu, Bluetooth etc. Evident, sistemul de operare aferent trebuie să fie capabil să gestioneze eficient aceste comunicații. Din cauza formatului mic, memoria internă, dar mai ales memoria pe suport extern este foarte limitată. Sistemul de operare trebuie să facă față acestei penurii de resurse. Securizarea accesului la un astfel de sistem este esențială. Fiind echipamente mici pot fi ușor pierdute / furate, iar această securizare trebuie să le facă de nefolosit de către persoanele neautorizate.

Sisteme conectate în rețea. Practic, astăzi marea majoritate a sistemelor nu mai sunt independente, ci sunt conectate la rețele de calculatoare, inclusiv la rețeaua publică Internet. Din această cauză, una dintre sarcinile fundamentale ale sistemelor de operare din această categorie trebuie să fie gestiunea accesului la rețea. De asemenea, accesul din exterior la resursele locale trebuie să fie bine protejat, spre a interzice accesele neautorizate.

În legătură cu acest criteriu de clasificare, trebuie arătat că el este relativ. Astfel, există sisteme seriale care permit interacțiunea cu programele, sistemele interactive au facilități puternice de lucru serial etc. Mai mult, același SO permite în același timp să lucreze, spre exemplu, în timp real, dar în același timp să servească, evident cu o prioritate mai mică, și alte programe de tip serial sau interactiv, în regim de multiprogramare clasică.

### **8.1.3 Clasificare după organizarea internă a programelor ce compun SO.**

Să vedem cum arată un SO văzut din "interior". Vom vedea patru tipuri de structuri, în evoluția lor istorică.

- Sisteme monolitice;
- Sisteme cu nucleu de interfață hardware;
- Sisteme cu structură stratificată;
- Sisteme organizate ca mașini virtuale

Sisteme monolitice. Un astfel de SO este o colecție de proceduri, fiecare dintre acestea putând fi apelată după necesități. Execuția unei astfel de proceduri nu poate fi întreruptă, ea trebuie să-și execute complet sarcina pentru care a fost apelată. Un program utilizator rulat sub un sistem monolitic se comportă ca o procedură apelată de SO. La rândul lui, programul poate apela în maniera apelurilor sistem (vezi 4.5 pentru Unix) diverse rutine din SO. Singura structură posibilă aici este legată de cele două moduri de lucru: nucleu și user (vezi pentru Unix 4.5 și pentru Windows 7.1.3). Esențial pentru aceste tipuri de sisteme este legătura de tip apel-revenire, fără posibilitatea de întrerupere a fluxului normal al execuției. În prezent, aceste tipuri de SO sunt pe cale de dispariție.

Sisteme cu nucleu de interfață hardware. Un astfel de sistem concentrează sarcinile vitale de nivel cel mai apropiat de hardware într-o colecție de rutine care formează nucleul SO. În acesta sunt înglobate inclusiv rutinele de întrerupere. Componentele nucleului se pot executa concurrent. Toate acțiunile utilizatorului asupra echipamentului hard trec prin acest nucleu. Unele SO mai plasează între nucleu și utilizator încă o interfață. Exemple de nuclee ale SO pot fi considerate componentele BIOS ale mașinilor PC deși fac parte din hardware.

Sisteme cu structură stratificată este o generalizare a organizării cu nucleu. Un astfel de sistem este construit nivel după nivel, componentele fiecărui nivel folosind toate serviciile oferite de nivelul inferior. În prezent aceste sisteme sunt cele mai răspândite. "Părinții" lor sunt SO THE (Dijkstra, 1968) și SO MULTICS. Ultimul dintre ele oferă o organizare mai interesantă, și unică în felul ei, aceea a inelelor de protecție. Pentru detalii se pot consulta [48], [49]. În secțiunea următoare vom ilustra structura stratificată a unui SO ipotetic.

Sisteme organizate ca mașini virtuale. Echipamentul hard al acestor tipuri de sisteme servește, în regim de multiprogramare, eventual în time-sharing, un număr de procese. Fiecare proces dispune, în mod exclusiv, de o serie de resurse, cea mai importantă fiind memoria. Fiecare dintre procesele deservite este un sistem de operare, care are la dispoziție toate resursele alocate procesului respectiv de către echipamentul hard. În acest mod, pe același echipament hard se poate lucra simultan sub mai multe SO. Primele sisteme de acest tip au fost VM/370 (Virtual Machines for IBM-370), care coordonează mai multe SO conversaționale monoutilizator de tip CMS (Conversational Monitor System). Fiecare utilizator lucrează sub propriul CMS, pentru el fiind transparent faptul că lucrează sub CMS singur, pe un calculator mic, sau că este legat, împreună cu alții la același echipament hard. Sistemele actuale Windows oferă mașini virtuale (ferestre) pentru lucru sub DOS. Implementările Linux, FreeBSD și Solaris actuale oferă mașina virtuală Wine care emulează platforme Windows. Pentru detalii, se pot consulta [4], [49].

## **8.2 Structura și funcțiile unui sistem de operare**

### **8.2.1 Stările unui proces și fazele unui program**

#### **8.2.1.1 Stările unui proces**

În secțiunile precedente am folosit destul de des termenul de proces, privit intuitiv ca un program în execuție. Am văzut în capitolele precedente noțiunea de proces sub Unix și sub Windows. În fig. 5.18 din 5.2.3 sunt prezentate stările unui proces Unix. Există [62] și stări în care se poate afla un proces Windows.

Sintetizând stările proceselor sub diverse sisteme de operare, putem defini stările unui proces într-un sistem de operare. Aceste stări sunt:

- HOLD – proces pregătit pentru intrarea în sistem;
- READY – procesul este în memorie și este gata de a fi servit de (un) procesor;
- RUN – un procesor execută efectiv instrucțiuni (mașină) ale procesului;
- WAIT – procesul așteaptă terminarea unei operații de intrare / ieșire;
- SWAP – imaginea procesului este evacuată temporar pe disc;
- FINISH – procesul s-a terminat, trebuie livrate doar rezultatele.

HOLD este starea unui proces la un sistem serial în care procesului îi sunt citite datele de lansare din cadrul lotului din care face parte, se face o analiză preliminară a corectitudinii lor și apoi procesul este depus într-o coadă pe disc (numită HOLD) spre a fi preluat spre prelucrare. De obicei, aceste sisteme au o componentă specializată pentru astfel de prelucrări – SPOOL-ing de intrare (Simultaneous Peripheral Operations OnLine – vezi [10]).

FINISH este, de asemenea, stare a unui proces la un sistem serial. Procesul se execută și rezultatele lui sunt plasate pe disc într-o coadă FINISH. După terminare, acestea vor fi listate,



fie pe o imprimantă locală, fie pe una aflată la distanță. Componenta din SO care face preia din coada FINISH și listează se numește SPOOL-ing de ieșire (vezi [10]).

READY este starea în care un proces se află în memoria internă dar nu este servit de procesor. Este posibil ca el să fi fost servit parțial, dar pe moment să fie suspendat în defavoarea altui proces, urmând să i se continue execuția mai târziu.

SWAP este starea în care un proces este evacuat temporar într-un spațiu rezervat pe disc (SWAP) pentru a face temporar loc altui proces în memoria internă. Ulterior procesul va fi reîncărcat în memoria internă și i se va continua execuția.

WAIT este starea în care un proces a cerut executarea unei operații de intrare / ieșire. Cât timp procesul așteaptă să se termine operația, cedă altor procese procesorul.

RUN este starea principală, aceea în care procesorul execută efectiv instrucțiuni mașină ale programului procesului. Într-un sistem monoprocesor doar un singur proces se află în starea RUN. Dacă sistemul are  $n$  procesoare, atunci maximum  $n$  procese se vor afla în starea RUN.

Aceste stări nu se vor întâlni la toate tipurile de sisteme de operare. De exemplu, la un sistem monoutilizator (vezi 8.1) nu va fi starea SWAP, numai sistemele seriale (vezi [10]) au stările HOLD și FINISH etc.

### 8.2.1.2 Fazele unui program

Fazele unui program reprezintă etapele prin care trece acesta din momentul în care s-a început proiectarea lui și până când devine cod mașină executabil spre a fi integrat într-un proces. Aceste faze, bine cunoscute de către programatori, sunt:

1. Editarea textului sursă într-un limbaj de programare, executată de către programator cu ajutorul unui editor de texte.
2. Compilarea, executată de un compilator specializat în limbajul respectiv. Ca rezultat se obține un fișier obiect.
3. Editarea de legături, fază în care se grupează mai multe module obiect rezultate din compilări și se obține un fișier executabil. Editarea se face cu un program specializat, numit editor de legături.
4. Execuția programului înseamnă mai întâi încărcarea fișierului executabil în memorie cu ajutorul unui program încărcător (loader) și lansarea lui în execuție.
5. Testarea și depanarea programului, fază care alternează cu cele de mai sus până când proiectantul are certitudinea că programul lui este corect. Uneori această depanare poate fi asistată de către un program specializat numit depanator.

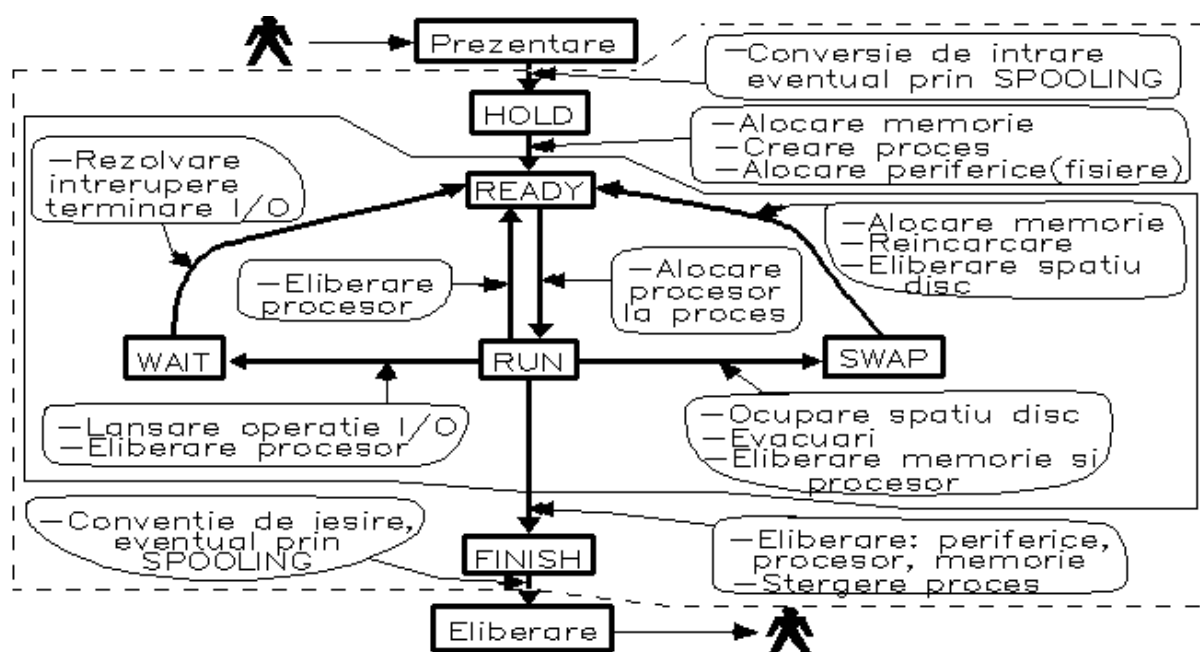
În secțiunile următoare vom reveni cu detalii, atât asupra stărilor unui proces cât și asupra fazelor unui program.

Menirea unui SO, indiferent de tipul lui, este pe de o parte de a facilita accesul la sistem a unuia sau mai multor utilizatori, iar pe de altă parte de a asigura o exploatare eficientă a echipamentului de calcul. Vom încerca să arătăm mai în detaliu care sunt funcțiile posibile ale unui sistem, independent de tipul acestuia. Desigur, dependentă de tip va fi ponderea în cadrul sistemului a uneia sau alteia dintre funcții. Din motive care se vor clarifica mai târziu, vom folosi în loc de termenul program termenul de proces sau task. Printr-un proces vom înțelege un calcul care poate fi efectuat concurent cu alte calcule (asupra acestei noțiuni vom reveni).

S-au proiectat și se proiectează încă multe tipuri de SO. Fiecare dintre acestea trebuie să aibă ca obiective fundamentale:

- optimizarea utilizării resurselor;
- minimizarea efectului uman de programare și exploatare;
- automatizarea operațiilor manuale în toate etapele de pregătire și exploatare a SC;
- creșterea eficienței utilizării SC prin scăderea prețului de cost al prelucrării datelor.

Principalele funcții ale SO sunt legate de acțiunile acestuia pentru tranziția între diverse stări ale procesului. În fig. 8.1 [10], [47] sunt schematizate stările unui proces. În dreptul fiecărei tranziții sunt trecute principalele acțiuni efectuate de către SO. Caracteristicile fiecăreia dintre cele șase stări posibile (HOLD, READY, RUN, SWAP, WAIT, FIHISH) le-am prezentat în secțiunea precedentă.



**Figura 8.1 Stările unui proces și acțiunile SO aferente**

Acțiunile SO ilustrate în fig. 8.1 sunt executate de către componenta nucleu a sistemului de operare. La nivelul componentei user, o mare parte dintre acțiunile SO sunt focalizate mai ales în sprijinul utilizatorului pentru dezvoltarea de programe, adică de asistare a userului la trecerea programului prin diverse faze (vezi secțiunea precedentă 8.2.1.2).

Am exclus de aici, de fapt am amânat numai până la o secțiune viitoare, o funcție primordială: punerea în lucru a SO. După cum vom vedea, în mod (aparent) paradoxal, această sarcină revine tot SO! La prima vedere pare a fi vorba de o antinomie, cum este posibil ca ceva (cineva) să-și dea viață singur? Vom vedea.

## 8.2.2 Structura generală a unui sistem de operare

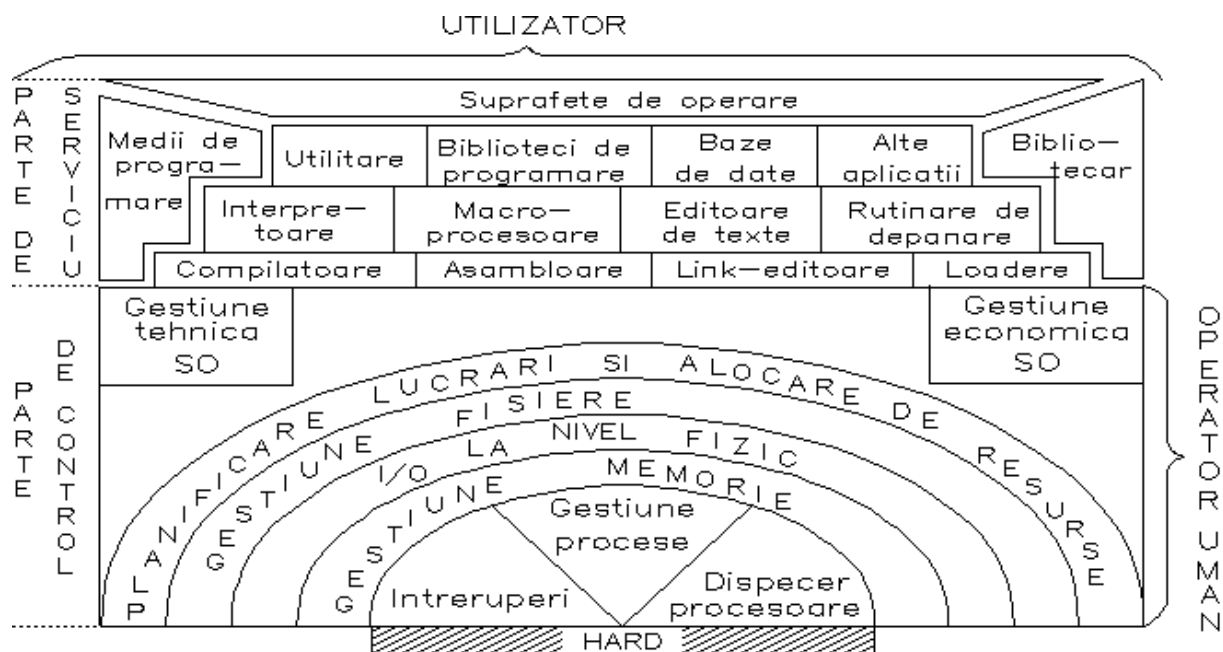
În această secțiune încercăm să oferim o imagine asupra structurii unui SO ipotetic. Un astfel de sistem nu va fi regăsit nicăieri în practică, însă structura fiecărui SO real va fi apropiată de a acestuia. Evident, în funcție de tipul de SO real, unele dintre prezentele componente vor fi atrofiate sau vor lipsi cu desăvârșire, iar altele vor putea fi eventual mult mai dezvoltate decât

aici. Fără pretenția de "gen-proxim" și "diferență-specifică", încercăm să definim principalele noțiuni, structuri și componente cu care operează orice SO.

Să trecem acum la structura propriu-zisă.

### 8.2.2.1 Structura unui SO

Schema generală a unui SO este ilustrată în fig. 8.2.



**Figura 8.2 Structura generală a unui SO**

Orice SO conține două părți mari: partea de control și partea de servicii.

**PARTEA DE CONTROL** este în legătură directă cu operatorul uman și indirectă cu utilizatorul. SO execută în mod master componentele acestei părți. Principalul ei rol este de a realiza legătura cu SC, deci cu echipamentul hard.

**PARTEA DE SERVICIU** lucrează în mod slave, folosind facilitățile părții de control. Este în legătură directă cu utilizatorul și indirectă cu operatorul uman. Sarcina de bază a ei este de a ține legătura cu utilizatorul, facilitându-i acestuia accesul la resursele SC și SO.

Ideea schemei din fig. 8.2 este preluată din [10]. În continuare vom descrie succint rolul fiecărei componente.

### 8.2.2.2 Partea de control

**INTRERUPERI.** Un ansamblu de rutine, fiecare dintre ele fiind activată la apariția unui anumit semnal fizic de întrerupere. Corespunzător, SO efectuează acțiunile aferente întreruperii respective. Semnificația conceptului de întrerupere este descrisă în fiecare manual de arhitectura calculatoarelor. Pentru întreruperile 8086, se poate consulta [52].

Funcționarea concretă a componentei (sistemului) de întreruperi depinde foarte mult de particularitățile SC, deci nu se pretează prea bine la abordări cu caracter general. Pe de altă parte, detalii despre aceste componente interesează în primul rând pe programatorii și inginerii de sistem, mai puțin pe utilizatorii obișnuiți., cărora ne adresăm noi. Din aceste motive vom schița numai mecanismul de funcționare a întreruperilor.

Fiecare tip de întrerupere are asociată o locație fixă de memorie. În această locație se află o adresă. Această adresă indică locul din memorie la care se află o secvență de instrucțiuni, numită handler, secvență care deservește întreruperea respectivă.

La apariția semnalului de întrerupere, după ce instrucțiunea mașină în curs s-a executat, se execută, în această ordine, următoarele activități:

- se salvează într-o zonă de memorie starea programului în curs de desfășurare;
- se activează handlerul asociat întreruperii respective;
- handlerul execută acțiunile necesare servirii întreruperii;
- se restaurează starea programului care a fost întrerupt.

**GESTIUNE PROCESE.** Creează procese și rezolvă problemele privind cooperarea și concurența acestora. La terminarea unui proces, acesta este șters din sistem. Asupra conceptului de proces vom mai reveni.

**DISPECER PROCESOARE.** La sistemele multiprocesor, repartizează sarcinile de calcul solicitate de procese, procesoarelor care sunt libere și care pot executa sarcinile cerute.

**GESTIUNEA MEMORIEI.** Alocă necesarul de memorie internă solicitat de procese și asigură protecția memoriei interprocese. Componenta este parțial realizată prin hard și parțial este realizată soft, de către o componentă a SO. Vom reveni și asupra problemelor legate de gestiunea memoriei.

**I/O LA NIVEL FIZIC.** Asigură efectuarea operațiilor elementare de I/O cu toate tipurile de periferice din sistem realizând, acolo unde este posibil, desfășurarea simultană a uneia sau mai multor operații I/O, cu activitatea procesorului central. Această componentă este, de asemenea, puternic dependentă hard, în special de multitudinea de echipamente periferice conectate la sistem. Așa după cum am văzut, există niște rutine de interfață numite drivere, care realizează compatibilizarea convențiilor de I/O între SC și echipamentele periferice. Driverele printre altele, fac parte din această componentă. Vom reveni într-o secțiune viitoare cu detalii.

**GESTIUNEA FISIERELOR.** O colecție de module prin intermediul cărora se asigură: deschiderea, închiderea și accesul utilizatorului la fișierele rezidente pe diferite suporturi de informații. Componentă de bază a SO, este cel mai des invocată de către utilizatori și de către operator. Vom reveni și asupra ei într-un capitol special.

**PLANIFICAREA LUCRARILOR SI ALOCAREA RESURSELOR.** Resursele unui SC se împart în: resurse fizice (procesoare, memorie internă, periferice etc.) și resurse logice (proceduri comune, tabele ale sistemului etc.). Fiecare proces solicită la un moment dat anumite resurse. Prin această componentă a SO se asigură planificarea proceselor astfel, încât ele să poată obține de fiecare dată resursele necesare, în mod individual sau partajate cu alte procese. Elemente ale planificării se răsfrâng, după cum vom vedea mai târziu, asupra gestiunii proceselor, gestiunii memoriei, intrărilor și ieșirilor etc. Atunci când vom discuta despre procese, memorie și I/O la nivel fizic vom aborda și problematica planificărilor și a alocărilor de resurse.

**GESTIUNEA TEHNICA** a SO. Tine evidența erorilor hard apărute la echipamentele SC. La cerere, furnizează informații statistice asupra gradului de utilizare a componentelor SC.

**GESTIUNEA ECONOMICA** a SO. Tine evidența utilizatorilor care folosesc sistemul, lucrările executate de către aceștia și resursele consumate de aceste lucrări (timp de rulare, memorie internă ocupată, echipamente și suporti folosiți etc). Periodic (zilnic, lunar, anual) editează liste cuprinzând lucrările rulate și facturi de plată către beneficiari.

### 8.2.2.3 *Partea de servicii*

**COMPILATOARELE** sunt programe care traduc texte sursă scrise în limbaje de nivel înalt (C, C++, Java, FORTRAN, COBOL, Pascal, Ada etc.) în limbaj mașină. La această oră există o mare varietate de tehnici de compilare, având ca nucleu teoria limbajelor formale. Există de asemenea sisteme automate de elaborare a diverselor componente ale compilatorului. Rezultatul unei compilări este un modul obiect memorat într-un fișier obiect. Detalii privind compilatoarele se pot obține consultând [33].

**ASAMBLORUL** este un program care traduce texte scrise în limbajul de asamblare propriu SC într-un modul obiect. Din punct de vedere al SO problemele legate de asamblor sunt incluse în problemele compilatoarelor. De exemplu, programarea în limbaj de asamblare 8086 este descrisă în [52].

**LINK-EDITORUL** sau **EDITORUL DE LEGATURI**. Grupează unul sau mai multe module obiect, împreună cu subprograme de serviciu din diverse biblioteci, într-un program executabil. De cele mai multe ori această componentă oferă posibilitatea segmentării programelor mari, adică posibilitatea ca două componente diferite a programului să poată ocupa în execuție, la momente diferite de timp, aceeași zonă de memorie. Rezultatul editării de legături este un fișier executabil. În 5.1 este descris formatul executabil de sub Unix, iar în [52] fișierele executabile COM și EXE de pe platformele Microsoft.

**LOADER (PROGRAM DE INCARCARE)** este un program al SO care are trei sarcini:

- citirea unui program executabil de pe un anumit suport;
- încărcarea acestuia în memorie;
- lansarea lui în execuție.

Dacă programul executabil este segmentat, atunci loaderul încarcă la început numai segmentul rădăcină al programului, după care, la cerere, în timpul execuției, încarcă de fiecare dată segmentul solicitat în memorie. În [52] este descris mecanismul de încărcare de sub MS-DOS.

**INTERPRETOR** este un program care execută pas cu pas instrucțiunile descrise într-un anumit limbaj. Tehnica interpretării este aplicată la diferite limbaje de programare, cum ar fi Java, BASIC, LISP, PROLOG etc. Pentru SO, cea mai importantă categorie de interpretoare sunt "interpretoarele de comenzi". În 2.1.1 am descris interpretoarele de comenzi Shell de sub Unix, iar în [10] [52] interpretorul de comenzi COMMAND.COM de pe platformele Microsoft.

**MACROPROCESOR (PREPROCESOR)** este un program cu ajutorul căruia se transformă o machetă de program, pe baza unor parametri, într-un program sursă compilabil. Spre exemplu, fie macheta:

```
MACRO      ADUNA      &1, &2, &3
LOAD       &1
ADD        &2
STORE     &3
ENDMACRO
```

Dacă într-un program (dintr-un limbaj de asamblare) se scrie:

```
ADUNA  A, B, C
```

atunci prin macroprocesoare se generează:

```
LOAD    A
ADD     B
STORE   C
```

Practic, macroprocesarea este tehnica prin care se realizează o scriere mai compactă a programelor și o posibilă parametrizare a acestora pentru realizarea comodă a unor implementări înrudite.

Astfel de macroprocesoare există la toate SO moderne. Remarcăm în mod deosebit interpretorul de comenzi "shell" de sub Unix, descris în 2.2. Acesta dispune, printre altele, de cele mai puternice facilități de macroprocesare dintre toate interpretoarele de comenzi cunoscute.

Unele limbaje au prevăzute în standardele lor specificații de macroprocesare. În acest sens amintim preprocesorul "C" care tratează în această manieră liniile care încep cu "#" și construcțiile "typedef", precum și facilitățile generice ale limbajului ADA .

EDITORUL DE TEXTE este un program interactiv destinat introducerii și corectării textelor sursă sau a unor texte destinate tipăririi. Orice SO interactiv dispune cel puțin de câte un astfel de editor de texte. Până vom ajunge la capitolul dedicat editoarelor de texte, să enumerăm câteva editoare "celebre" unele prin vechimea lor, iar altele prin performanțele lor:

- ed și vi sub Unix;
- Notepad și EDIT sub Windows;
- xedit și emacs de sub mediul X WINDOWS.

RUTINELE DE DEPANARE (DEPANATOARELE) asistă execuția unui program utilizator și-l ajută să depisteze în program cauzele erorilor apărute în timpul execuției lui. Intreaga execuție a programului utilizator se face sub controlul depanatorului. La prima apariție a unui eveniment deosebit în execuție (depășire de indici, adresă inexistentă, cod de operație necunoscut etc.) depanatorul primește controlul. De asemenea, el poate primi controlul în locuri marcate în prealabil de către utilizator. Un astfel de loc poartă numele de breakpoint. La primirea controlului, depanatorul poate afișa valorile unor variabile sau locații de memorie, poate reda, total sau parțial istoria de calcul, poate modifica valorile unor variabile sau locații de memorie etc. Unele depanatoare mai performante permit execuția reversibilă, concept tratat teoretic dar mai puțin aplicat practic. La această oră sunt cunoscute două tipuri de depanatoare:

- depanatoare mașină;
- depanatoare simbolice.

Depanatoarele mașină operează cu elemente cum ar fi: adrese fizice, configurații de biți, locații de memorie etc. Fiecare SO dispune de câte un astfel de instrument de depanare. Pentru

exemplificări se pot consulta documentațiile MS\_DOS pentru DEBUG, precum și documentația Unix despre depanatoarele adb și gdb [47].

Depanatoarele simbolice operează cu elemente ale textelor sursă ale programelor. Întâlnim aici termeni ca linie sursă, nume de variabilă, nume de procedură, stivă de apel a procedurilor etc. Evident, depanatoarele simbolice sunt mult mai tentante pentru utilizator, și este bine că este așa. De obicei, un astfel de depanator însoțește implementarea unui limbaj de programare de nivel înalt. Primul și cel mai răspândit limbaj la care au fost atașate depanatoare simbolice este limbajul Pascal. Amintim aici depanatorul Pascal de pe lângă implementarea TURBO Pascal, un depanator simbolic complet interactiv, cu o grafică destul de bună pentru momentul implementării lui și cu facilități puternice. Firma BORLAND, un adevărat vârf de lance în acest domeniu, a realizat printre altele produsul TURBO DEBUGGER, care recunoaște în vederea depanării toate limbajele pentru care s-au realizat medii TURBO: Pascal, "C", limbajul de asamblare TASM, Basic, PROLOG etc. Produsul lucrează atât simbolic, cât și la nivel mașină.

BIBLIOTECARUL este un program cu ajutorul căruia utilizatorul poate comanda păstrarea programelor proprii în fișiere speciale de tip bibliotecă, poate copia programe dintr-o bibliotecă în alta, șterge, inserează și modifică programe în aceste biblioteci. Exemple de programe biblioteci sunt LBR sub MS-DOS, ar (archive) sub Unix etc.

MEDIILE DE PROGRAMARE sunt componente complexe destinate în principal activității de dezvoltare de programe. O astfel de componentă înglobează, relativ la un limbaj, cel puțin următoarele subcomponente:

- editor de texte;
- compilator interactiv;
- compilator serial (clasic);
- editor de legături sau un mecanism echivalent acestuia destinat creerii de programe executabile;
- interpretor pentru execuția rezultatului compilării;
- depanator, de regulă simbolic;
- o componentă de legătură cu mediul exterior, de regulă cu SO sub care se lucrează.

De regulă, filozofia de manevrare a mediului are un caracter unitar, este simplă și puternică. Facilitățile grafice ale acestor medii sunt destul de mult utilizate, făcând deosebit de agreabil lucrul cu ele.

SUPRAFETELE DE OPERARE au apărut tot la calculatoarele IBM-PC și a celor compatibile cu ele. Ele au rolul principal de a "îmbrăca" sistemul de operare (Windows sau Unix), pentru a-i face mai ușoară operarea. Inițial au fost destinate utilizatorilor neinformaticieni. La această oră, nici profesioniștii nu se mai pot "dezlipi" de o suprafață de operare cu care s-a obișnuit să lucreze. De regulă, suprafețele de operare înlocuiesc limbajul de control (de comandă) al SO respectiv. Utilizatorul, în loc să tasteze o comandă a SO respectiv, poartă un dialog cu suprafața de operare, iar la finele dialogului suprafața generează sau execută comanda respectivă. Ponderea dialogului o deține suprafața însăși, utilizatorul răspunde de regulă prin apăsarea uneia sau maximum a două taste, sau unul sau două "clicuri" cu un mouse. Suprafața deține și un mecanism deosebit de puternic de HELP, conducând utilizatorul din aproape în aproape spre scopul dorit. Deși se adresează în primul rând neinformaticienilor, ele sunt din ce în ce mai mult folosite și de către specialiști.

O remarcă specială trebuie făcută pentru WINDOWS, care deține un mecanism propriu de dezvoltare a unor aplicații.

Din rațiuni impuse de standarde, sub Unix se folosește aproape peste tot suprafața X-WINDOWS [43].

Din scurta prezentare a componentelor principale ale unui SO rezultă caracterul modular al SO și organizarea lui ierarhică. Modularitatea este termenul prin care se descrie partiționarea unui program mare în module mai mici, cu sarcini bine precizate. Caracterul ierarhic este reflectat în faptul că orice componentă acționează sub controlul componentelor interioare acesteia, apelând primitivele oferite de acestea. Aceste facilități asigură o implementare mai ușoară a SO, o înțelegere mai facilă a SO de către diverse categorii de utilizatori, o întreținere și depanare mai ușoară a sistemului și creează premisele simulării lui pe alt sistem.

### **8.3 Incărcarea (lansarea în execuție) a unui sistem de operare**

Una dintre sarcinile de bază ale unui SO este aceea de a se putea autoîncărca de pe disc în memoria internă și de a se autolansa în execuție. Această acțiune se desfășoară la fiecare punere sub tensiune a SC, precum și ori de câte ori utilizatorul (operatorul uman) crede de cuviință să reîncarce SO. Acțiunea este cunoscută sub numele de încărcare a SO sau lansare în execuție a SO.

În cele ce urmează dorim să punem în evidență mecanismul de lansare. O vom face însă pentru un SO ipotetic, fără să ne legăm de un anumit sistem de operare concret, simplificând la maximum expunerea.

Pentru lansare se utilizează un mecanism combinat hard și soft, numit bootstrap. Pentru a putea înțelege mai bine modul lui de funcționare, să considerăm un exemplu.

Să notăm cu  $M[0]$ ,  $M[1]$ ,  $M[2]$ , . . . locațiile de memorie ale unui calculator ipotetic. Presupunem că în fiecare locație  $M[i]$  încapă un număr întreg. Presupunem, de asemenea, că sistemul nostru dispune de o memorie ROM și un disc (dischetă, CD, DVD, hard disc, partiție disc etc.) numit disc de boot, de la care se citesc, secvențial, începând din primul sector, numere întregi.

Prin  $\text{read}(x)$  vom înțelege că se citește următorul număr de la intrarea standard și valoarea lui se depune în locația  $x$ .

Prin  $\text{transferla}(y)$  înțelegem că următoarea instrucțiune mașină de executat este cea din (al cărei cod se află în) locația  $y$ .

Mecanismul bootstrap intră în lucru la apăsarea butonului de pornire <START>. Ca prim efect, sistemul extrage din ROM un șir de octeți reprezentând echivalentul în limbaj mașină al programului din fig. 8.3.

```
APASASTART:
  for (i = 0; i < 100 ; i++) {
    read(M[i]);
  }
  transferla(M[0]);
```

**Figura 8.3 Programul APASASTART**



Sirul de octeți citit este depus în memorie la o adresă pe care o vom nota APASASTART. După depunerea la adresa respectivă, sistemul execută instrucțiunea:

transferla(APASASTART);

Cu aceasta, sistemul transferă controlul programului din fig. 8.3. Execuția acestui program înseamnă că sistemul citește de pe discul de boot primele 100 de numere, le depune în primele 100 de locații de memorie, după care trece la executarea instrucțiunii al cărei cod se află în prima locație.

Este limpede că în 100 de locații nu se poate scrie decât un program simplu. Ori pentru lansarea în execuție a unui SO, probabil că este necesară punerea în lucru a unui program mai sofisticat. Un astfel de program poate fi scris pe mai multe sectoare consecutive, folosind, de exemplu reprezentarea din fig. 8.4. Să presupunem că într-un sector se pot înregistra  $S$  numere întregi.

AM1	n1	n1 întregi de program sau date	$S - n1 - 2$ întregi zero
AM2	n2	n2 întregi de program sau date	$S - n2 - 2$ întregi zero
---	---	---	---
AMk	Nk	nk întregi de program sau date	$S - nk - 2$ întregi zero
ADL	-1		$S - 2$ întregi zero

**Figura 8.4 Structura unui fișier ce conține un program executabil**

Deci fișierul executabil din fig. 8.4 poate avea un număr oarecare  $k$  de înregistrări a câte un sector, conținând numere întregi. Semnificațiile câmpurilor sunt:

- $AM_i$  adresa din memorie de unde începe introducerea porțiunii de cod de program sau date din înregistrarea  $i$ ;
- $n_i$  este numărul efectiv de locații ocupat de porțiunea de program sau date a înregistrării  $i$ . Valoarea -1 indică ultimul sector al a programului.
- ADL se află în ultimul sector al programului și este adresa de lansare în execuție a programului reprezentat: un număr între  $AM_1$  și  $AM_k + n_k - 1$ .

Să considerăm acum programul din fig. 8.5, în care presupunem că începând de la locația  $M[r]$  există  $S$  locații libere.

```

LOADERABSOLUT:
do {
  for ( i=0; i < S; i++) {
    read(M[r+i]);
  }
  A = M[r];
  n = M[r+1];
  if (n== -1)
    break;
  for ( i=0; i < n; i++) {
    M[A+i] = M[r+2+i]
  }
  while (FALSE);
LANSARE:
transferla(A);

```

### Figura 8.5 Incărcător (loader) de programe cu structura din fig. 8.4

Avem acum posibilitatea să citim de pe disc și să lansăm în execuție, cu programul din fig. 8.5, programe oricât de mari cu formatul din fig. 8.4. De fapt, în fig. 8.5 avem de-a face cu un exemplu simplu de program de încărcare (loader).

Cheia (saltul calitativ al) încărcării: În cele 100 de numere care vor fi citite în momentul startului cu programul APASASTART (din fig. 8.3), se trece codul mașină al programului LOADERABSOLUT (din fig. 8.5). Să vedem acum, pe etape, cum are loc încărcarea și lansarea în execuție a SO. În fig. 8.6 sunt ilustrate zonele de memorie folosite în faza de încărcare și lansare în execuție a SO.

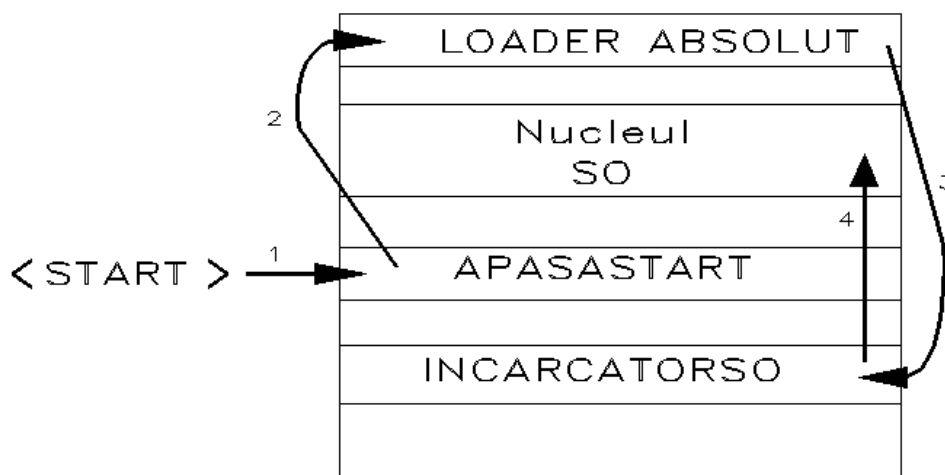


Figura 8.6 Incărcarea unui SO

Etapele încărcării sunt:

- 1) Se apasă <START>. Ca efect, programul APASASTART citește, în primele 100 de locații de memorie, codul mașină al programului LOADERABSOLUT.
- 2) După transferla(M[0]) intră în lucru programul LOADERABSOLUT.
- 3) Acesta citește la rândul lui un program, care de această dată poate fi oricât de sofisticat, cu formatul din fig. 8.4. Acest program îl vom numi INCARCATORSO.
- 4) După ce se execută și ultima instrucțiune transferla(A) a programului LOADERABSOLUT, intră în lucru programul INCARCATORSO. Dat fiind faptul că formatul de reprezentare a codului lui mașină este cel din fig. 8.4, acesta poate fi practic oricât de mare. Așa cum îi spune și numele, acest program va încărca de pe disc programele și rutinele (nucleului) sistemului de operare.

La încheierea încărcării SO, INCARCATORSO execută următoarele acțiuni:

- șterge din memorie codul programului APASASTART;
- șterge din memorie codul programului LOADERABSOLUT;
- eliberează spațiul de memorie ocupat de el, cu excepția unei scurte porțiuni, despre care vorbim mai jos;
- execută o instrucțiune transferla(ASSO), unde ASSO este adresa de start a sistemului de operare. Codul necesar acestei instrucțiuni este cel pe care nu-l poate, în mod firesc, elibera. Însă dacă acest cod este plasat într-o zonă folosită ca buffer, practic nu mai rămâne memorie ocupată.
- după aceasta, încărcătorul își încheie activitatea.

În sfârșit, toate SO moderne lansează la terminarea încărcării un fișier căruia o să-i spunem fișier de comenzi inițiale. Acesta este un fișier de comenzi, care la rândul lui poate lansa alte fișiere de comenzi ș.a.m.d. text în care sunt trecute toate comenzile pe care administratorul le dorește a fi executate la pornirea sistemului. Printre ultimele acțiuni desfășurate la lansarea SO amintim:

- montarea partițiilor locale;
- efectuarea unor operații de întreținere asupra sistemului de fișiere (filesystem cleanup);
- pornirea diferitelor servicii (de exemplu serviciile de printare);
- setare nume sistem și inițializare ceasuri în funcție de zona geografică;
- configurează rețeaua;
- montează sisteme de fișiere aflate la distanță (așa-numitele sisteme de fișiere remote)
- permite conectarea utilizatorilor;

Proiectantul SO va avea grijă ca programul încărcător să fie depus într-o zonă de memorie care să nu fie ocupată de către programele SO pe care el le încarcă.

Față de modelul de încărcare exemplificat de noi mai sus, realitatea, conceptual, este foarte aproape. Primele calculatoare executau APASASTART citind o cartelă, pe care era înregistrat programul LOADERABSOLUT. Acesta realiza citirea programului INCARCATORSO de pe bandă sau de pe disc. Multe calculatoare foloseau în acest scop banda perforată (care uneori trebuia "ajutată" cu mâna să defileze prin cititor ☺). Imaginea "arhaică" a unei astfel de porniri este similară -cei în vârstă probabil că își aduc aminte- cu pornirea în anii '50 a celebrului tractor KD, sau a unei drujbe, pornire care se face trăgând de o sfoară. ☺)

În prezent suportul magnetic de încărcare este primul sector (sectorul zero) al unui disc magnetic. Programele care realizează încărcarea se vor deosebi de cele de mai sus, numai prin înlocuirea lui 100 cu 128, 256, 512 sau 1024, deci cu lungimea unui sector.

Unele SO pretind o modalitate prin care i se indică sistemului de pe ce disc să încarce programele LOADERABSOLUT și INCARCATORSO: specificarea se face în memoria CMOS, la cheile pupitrului calculatorului etc. Alte sisteme stabilesc singure, în mod automat, de pe ce disc să facă încărcarea.

## 9 Teoria generală a sistemelor de operare

### 9.1 Procese

#### 9.1.1 Conceptul de proces

Facilitatea de a rula mai multe programe simultan în cadrul unui sistem de operare este considerată astăzi normală de către toți utilizatorii. Rularea unui navigator web (Mozilla Firefox, Internet Explorer, Safari, Konqueror etc.) simultan cu rularea unui program pentru citirea poștei electronice (Mozilla Thunderbird, Outlook Express, Eudora, Pine etc.) este o practică de zi cu zi a majorității utilizatorilor. Din punctul de vedere al unui proces de operare toate aceste programe rulând pe un sistem de operare sunt considerate procese sau task-uri.

Formal, un *proces* sau *task*, este un calcul care poate fi executat concurent (în paralel) cu alte calcule. Este o abstractizare a activității procesorului care pur și simplu execută instrucțiunile care i se transmit fără a face diferențiere între procesele cărora le aparțin. Revenind la exemplul inițial, un procesor nu “știe” dacă instrucțiunea executată aparține de Mozilla Firefox sau Outlook Express.

În terminologia actuală există o tendință de confuzie între noțiunea de proces și cea de program. Deși folosirea noțiunii de program pentru a ne referi la un proces nu creează în general confuzii, considerăm totuși necesară clarificarea situației. *Procesul are un caracter dinamic*, el precizează o secvență de activități în curs de execuție, iar *programul are un caracter static*, el numai descrie această secvență de activități. Cu alte cuvinte procesul este un program în execuție.

Existența unui proces este condiționată de existența a trei factori:

1. *procedură* (un set de instrucțiuni) care trebuie executată;
2. un *procesor* care să poată executa aceste instrucțiuni;
3. un *mediu* (memorie, periferice) asupra căruia să acționeze procesorul conform celor precizate în procedură.

Execuția unui program (și implicit transformarea acestuia în proces) implică execuția de către procesor a fiecărei instrucțiuni din programul (procedura) respectiv(ă). În mod inerent, aceste instrucțiuni vor avea efect asupra mediului în care procesul rulează, cum ar fi: alocarea de memorie, folosirea perifericelor, tipărirea unui mesaj pe terminal, consumul de timp etc.

Trebuie menționat că în arhitecturile actuale monoprocesor (vezi [52]), mai multe procese care rulează concomitent sunt de fapt deservite alternativ de către procesor. El (procesorul) execută alternativ grupuri de instrucțiuni din fiecare program de-a lungul unei cuante de timp, după care comută la alt grup de la alt program. Dimensiunea cuantei de timp este astfel aleasă (și așa de mică) încât momentele de staționare ale unui program nu sunt sesizabile de către utilizator. Astfel se creează impresia de simultaneitate în execuție.

Paralelismul este efectiv doar în cadrul sistemelor multiprocesor. Practic, două procese pot fi executate simultan în sens propriu doar dacă rulează pe o mașină bi-procesor.

### 9.1.2 Concurența între procese

Existența simultană a mai multor procese într-un sistem de operare ridică probleme în ceea ce privește accesul la resursele sistemului. Prin resursă înțelegem orice este necesar unui proces pentru a-și continua activitatea. Câteva resurse mai populare sunt:

- memoria (cereri de alocare de memorie)
- discul magnetic (acces la fișiere)
- terminalul (pentru interacțiunea cu utilizatorul)
- interfața de rețea (pentru interacțiunea cu alte mașini legate prin rețea)

#### 9.1.2.1 Secțiune critică; resursă critică; excludere mutuală.

O problemă legată de accesul la resurse este asigurarea corectitudinii operațiilor executate în regim de concurență. Pentru a vedea cum operații corecte pot da rezultate greșite în regim de concurență, vom analiza un exemplu clasic. Să considerăm un fișier care stochează numărul de locuri libere într-un sistem de vânzare de bilete. Vânzarea unui bilet va avea ca și consecință decrementarea numărului de locuri libere. De asemenea, să considerăm că există două procese care efectuează simultan operații asupra acestui fișier. Fragmentul de cod necesar pentru decrementarea numărului de locuri este prezentat în fig. 9.1.

```
int n;
int fd = open("locuri.db", "ORDWR");
read(fd, &n, sizeof(int));
lseek(fd, 0, SEEK_SET);
n--;
write(fd, &n, sizeof(int));
close(f);
```

**Figura 9.1 O secvență de decrementare**

În situația în care codul din fig. 9.1 este rulat simultan de două sau mai multe procese, și considerând faptul că procesorul execută alternativ grupuri de instrucțiuni din procese diferite, este posibil ca instrucțiunile de mai sus să fie executate în procese separate în ordinea din figura 9.2.

Procesul A	n = (în A)	Procesul B	n = (în B)
read(fd, &n, sizeof(int))	10		
	10	read(fd, &n, sizeof(int))	10
n--	9		10
lseek(fd, 0, SEEK SET)	9		10
write(fd, &n, sizeof(int))	9		10
		n--	9
		lseek(fd, 0, SEEK SET)	9
		write(fd, &n, sizeof(int))	9

**Figura 9.2 Programul din fig. 9.1 rulat "defavorabil" în două procese**

În mod evident, ambele procese citesc *aceeași valoare* din fișier, o decrementează și apoi o scriu la loc. Ca urmare, în final, în loc ca numărul de locuri să scadă cu două poziții, va scădea doar cu una și ca urmare operațiile executate asupra fișierului deși sunt greșite. În viața de zi cu zi o asemenea eroare ar conduce la vânzarea aceluiași loc în sală de două ori (o "bucurie" cu care mulți dintre cititori au avut probabil de-a face).

Trebuie remarcat că eroarea poate să apară doar la o execuție concurentă, altfel secvența de cod fiind perfect corectă. Eroarea apare din lipsa de sincronizare între accesele proceselor la fișier. Procesul B, în loc să opereze asupra numărului de locuri decrementat de procesul A, reușește să citească numărul de locuri prea devreme și va opera asupra *aceluiași număr de locuri* ca și procesul A. Pentru a obține rezultate corecte, secvența de execuție dorită este cea din fig. 9.3. Adică secvența de cod care decrementează numărul de locuri să fie executată la un moment dat doar de un singur proces.

Procesul A	n = (în A)	Procesul B	n = (în B)
read(fd, &n, sizeof(int))	10		
n--	9		
lseek(fd, 0, SEEK_SET)	9		
write(fd, &n, sizeof(int))	9		
		read(fd, &n, sizeof(int))	9
		n--	8
		lseek(fd, 0, SEEK_SET)	8
		write(fd, &n, sizeof(int))	8

**Figura 9.3** Programul din fig. 9.1 rulat ca secțiune critică

Același gen de situație se creează și în momentul în care procesele scriu date pe disc. Dacă sistemul de operare ar servi aceste cereri în stilul din figura 9.2 două procese ar putea primi același bloc pe disc pentru a scrie în el date aparținând la fișiere diferite. Este ușor de imaginat ce s-ar întâmpla dacă într-un script shell editat cu un editor de text ar apărea părți din conținutul unui fișier ".o" generat de gcc la compilare ...

Revenind la problema noastră, vom spune că porțiunea de program din fig. 9.1 începând cu read și terminând cu write este o *secțiune critică*, deoarece nu este permis ca ea să fie executată simultan de către două procese. Analog, vom spune că poziția pe care se stochează numărul de locuri în fișier este o *resursă critică*, deoarece nu poate fi accesată simultan de mai multe procese. În sfârșit, pentru a evita situația eronată de mai sus vom spune că procesele A și B trebuie să se *excludă reciproc*, deoarece trebuie să aibă acces exclusiv la secțiunea și la resursa critică [22].

Problema secțiunii critice a suscitat un interes deosebit în istoria SO. Vom expune și noi câteva aspecte ale ei, pentru a se vedea câteva dintre "subtilitățile" programării concurente. Să vedem mai exact cerințele problemei:

- la un moment dat, numai un singur proces este în secțiunea critică; orice alt proces care solicită accesul la ea, îl va primi după ce procesul care o ocupă a terminat de executat instrucțiunile secțiunii critice;
- vitezele relative ale proceselor sunt necunoscute;
- oprirea oricărui proces are loc numai în afara secțiunii critice;
- nici un proces nu va aștepta indefinit pentru a intra în secțiunea critică.

Pentru a fixa ideile, presupunem că fiecare proces execută în mod repetat, de un număr nedefinit de ori, întâi secțiunea critică și apoi restul programului. Descrierea problemei astfel puse este dată în limbaj pseudo-cod în fig. 9.4.

```
do {
    <secțiune critică>
    <rest program>
} while(false);
```

**Figura 9.4** Un program cu secțiune critică

Au existat multe încercări de rezolvare a problemei secțiunii critice. În [4] sunt date și comentate trei încercări de descriere, fiecare dintre ele încălcând unele dintre condițiile puse mai sus. Matematicianul Dekker a fost primul (1965) care a dat o soluție software (rezonabilă, care însă este destul de complicată și încălăcită). În 1981 Peterson a dat o soluție ceva mai simplă. Aceste soluții se referă la rezolvarea secțiunilor critice apărute la concurența în cadrul aceluiași proces (un subiect pe care nu l-am abordat încă). Totuși, pentru a evidenția complexitatea problemei, prezentăm soluția lui Peterson în fig 9.5 și lăsăm cititorului "plăcerea" de a demonstra că soluția este corectă.

Cod inițial	<pre>int c1 = 1; int c2 = 1; int schimb;</pre>	
Cod executat în paralel	<pre>do {     c1 = 0;     schimb = 1;     while((c2 == 0) &amp;&amp;           (schimb == 1)){         &lt;așteaptă&gt;;     }     &lt;secțiune critică&gt;;     c1 = 1;     &lt;rest program 1&gt; } while(false)</pre>	<pre>do {     c2 = 0;     schimb = 2;     while((c1 == 0) &amp;&amp;           (schimb == 2)){         &lt;așteaptă&gt;;     }     &lt;secțiune critică&gt;;     c2 = 1;     &lt;rest program 2&gt; } while(false)</pre>

**Figura 9.5 Soluția Peterson pentru secțiune critică**

Evident, rezolvarea nu este simplă. Dacă se cere ca secțiunea critică să fie accesibilă la mai multe procese, soluția din fig. 9.5 nu poate fi generalizată prea ușor. Mai mult, soluțiile software presupun o *așteptare activă*, adică toate procesoarele, în faza de așteptare execută în mod repetat o aceeași instrucțiune (ciclează încontinuu). Ar fi de dorit ca pe timpul cât un proces așteaptă să primească dreptul de execuție a secțiunii critice, procesorul lui să devină disponibil spre a servi alte procese. O asemenea soluție poate fi realizată folosind conceptul de *semafor* descris în secțiunea următoare.

### 9.1.2.2 Conceptul de semafor.

Pentru a evita situația de mai sus și pentru a permite o serie de operații cu procese (pe care le vom prezenta mai târziu), Dijkstra a introdus conceptul de semafor.

Un *semafor*  $s$  este o pereche

$$(v(s), c(s))$$

unde  $v(s)$  este valoarea semaforului, iar  $c(s)$  o coadă de așteptare. Valoarea  $v(s)$  este un număr întreg, care primește o valoare inițială  $v_0(s)$ . Coada de așteptare conține (pointeri la) procesele care așteaptă la semaforul  $s$ . Inițial coada este vidă, iar disciplina cozii depinde de sistemul de operare (FIFO, LIFO, priorități etc.).

Pentru gestiunea semafoarelor se definesc două operații indivizibile  $P(s)$  și  $V(s)$  ale căror roluri, exprimate laconic, sunt "a trece de resursă" și "a anunța eliberarea resursei". (Denumirile operațiilor au fost date de Dijkstra după primele litere din limba olandeză: P de la *prolaag* prescurtare de la *probeer te verlagen* – încercare de a descrește, respectiv V de la *verhoog* – a crește). În unele lucrări [4], operația  $P$  se mai numește **WAIT**, iar operația  $V$  se

mai numește **SIGNAL**. Indivizibilitatea operațiilor înseamnă că ele nu pot fi întrerupte în cursul execuției lor, deci nu pot exista două procese care să execute simultan  $P(s)$  sau  $V(s)$ , sau simultan două operații  $P$ , sau simultan două operații  $V$ .

Definițiile exacte ale operațiilor  $P$  și  $V$ , apelate de un proces A și aplicate asupra unui semafor sunt date în fig. 9.6.

<pre>// P(s) apelat de procesul A: v(s) = v(s) - 1; if(v(s) &lt; 0) {     STARE(A) = WAIT;     c(s) ← A;     //Procesul A intra în așteptare     &lt;Trece controlul la DISPECER&gt;; } else {     &lt;Trece controlul la procesul A&gt;; }</pre>	<pre>// V(s) apelat de procesul A v(s) = v(s) + 1; if(v(s) &lt;= 0) {     c(s) → B;     // Extrage din coadă alt proces B     STARE(B) = READY;     &lt;Trece controlul la DISPECER&gt;; } else {     &lt;Trece controlul la procesul A&gt;; }</pre>
---	--

**Figura 9.6 Operațiile  $P(s)$  și  $V(s)$  apelate de procesul A**

Pentru a sugera funcționarea (și denumirea) semaforului vom considera un exemplu concret. Fie G1 și G2 două gări legate prin  $n$  linii paralele, pe care se circulă numai de la G1 spre G2. Presupunem că în G1 intră mai mult de  $n$  linii. Procesul trecerii trenurilor va fi dirijat astfel:

Valoarea inițială:  $v_0(s) = n$

Procesul de trecere:

$P(s); \quad \text{<Trenul trece pe una din cele } n \text{ linii>;} \quad V(s);$

Să notăm cu  $np(s)$ ,  $nv(s)$ , numărul de primitive  $P(s)$ , respectiv  $V(s)$  efectuate până la un moment dat,  $v_0(s)$  valoarea inițială a semaforului  $s$ , iar  $nt(s)$  numărul proceselor care au trecut de semaforul  $s$ . Următoarele proprietăți au loc:

1.  $v(s) = v_0(s) - np(s) + nv(s)$  ;
2. dacă  $v(s) < 0$  atunci în  $c(s)$  există  $-v(s)$  procese;
3. dacă  $v(s) > 0$  atunci  $v(s)$  procese pot trece succesiv de semaforul  $s$  fără să fie blocate ;
4.  $nt(s) = \min(v_0(s) + nv(s), np(s))$

Demonstrația acestor proprietăți este simplă, prin inducție după numărul operațiilor  $P$  și  $V$  efectuate asupra semaforului. Aspectul critic al semaforului este că operațiile asupra lui ( $P$  și  $V$ ) sunt atomice. Astfel două procese nu le pot executa în același timp asupra aceluiași semafor  $s$ .

Folosirea semafoarelor rezolvă complet și elegant problema secțiunii critice prezentată în secțiunile anterioare. Rezolvarea presupune folosirea unui singur semafor  $s$  și este corectă indiferent de numărul de procese care folosesc în comun secțiunea critică. Descrierea soluției este dată în fig. 9.7.

Cheia este tocmai semaforul  $s$ . Dacă  $v_0(s) = 1$  și toate procesele care folosesc secțiunea critică sunt de forma:

$...; P(s); \text{ <secțiune critică>; } V(s); ...$



atunci se demonstrează ușor că  $v(s) < 1$ , indiferent de numărul proceselor care operează asupra lui  $s$  cu operații **P** sau **V**. Drept urmare, în baza proprietăților menționate mai sus, un singur proces poate trece de acest semafor, deci sînt îndeplinite condițiile cerute pentru a proteja o secțiune critică. Un astfel de semafor mai poartă numele de *semafor de excludere mutuală*. Sincronizarea este o operație fundamentală în programarea concurentă. În secțiunile care urmează vom întâlni mai multe aplicații ale ei.

Cod inițial	semaphore s; v0(s) = 1;	
Cod executat în paralel	do { P(s); <secțiune critică>; V(s); <rest program 1> } while(false)	do { P(s); <secțiune critică>; V(s); <rest program 2> } while(false)

**Figura 9.7 Soluția cu semafor pentru secțiune critică**

### 9.1.2.3 Problema producătorului și a consumatorului.

Să presupunem că există unul sau mai multe procese numite *producătoare*, și unul sau mai multe procese numite *consumatoare*. De exemplu, conceptele de pipe și FIFO sunt de această natură (vezi 5.3.2 și 5.3.4). Transmiterea informațiilor de la producători spre consumatori se realizează prin intermediul unui buffer cu  $n$  intrări pentru  $n$  articole. Fiecare producător depune câte un articol în buffer, iar fiecare consumator scoate câte un articol din buffer. Problema constă în a dirija cele două tipuri de procese astfel încât:

- să existe acces exclusiv la buffer;
- consumatorii să aștepte când bufferul este gol;
- producătorii să aștepte când bufferul este plin.

Rezolvarea prin semafoare este foarte simplă. Se vor folosi trei semafoare: *gol*, *plin* și *exclus*. Rezolvarea este descrisă în fig. 9.8. Semaforul *exclus* asigură accesul exclusiv la buffer în momentul modificării lui, pentru a nu permite depuneri și extrageri de date simultane. Semaforul *plin* indică numărul de poziții ocupate din buffer, iar semaforul *gol* numărul de poziții libere din buffer. Depunerea unui articol în buffer necesită decrementarea (**P**) semaforului *gol* și incrementarea (**V**) semaforului *plin* (practic o poziție este “mutată” din semaforul *gol* în semaforul *plin*). Extragerea unui articol din buffer este operația inversă depunerii.

Producător	Consumator
semaphore plin, gol, exclus; v0(plin) = 0; v0(gol) = n; v(exclus) = 1;	
do { <produce articol>; P(gol); P(exclus); <depone articol în buffer>; V(exclus); V(plin); } while(false);	do { P(plin); P(exclus); <extrage articol din buffer>; V(exclus); V(gol); <consuma articol>; } while(false);

**Figura 9.8 Problema producătorului și a consumatorului rezolvată cu semafoare**

#### 9.1.2.4 Regiuni critice condiționate.

Prin intermediul operațiilor de excludere și sincronizare pot fi rezolvate o serie de probleme clasice de concurență. Dintre acestea au fost prezentate problema secțiunii critice și problema producătorului și a consumatorului. Am văzut că descrierea lor cu ajutorul semafoarelor este destul de comodă. Din păcate, operațiile **P** și **V**, lăsate libere la îndemâna programatorului, pot provoca necazuri mari. Încercați să vedeți ce se întâmplă dacă la descrierea din fig. 9.8 se inversează între ele operațiile **P(gol)** și **P(exclus)** din procesul producător. Acesta este motivul pentru care au fost introduse construcții structurate de concurență. Prin intermediul lor se poate exercita un control riguros asupra încălcării regulilor de concurență.

Vom descrie o astfel de construcție structurată, numită *regiune critică condiționată*. Ea înglobează într-un mod elegant conceptele de secțiune critică și resursă critică despre care am vorbit mai sus, cu posibilitatea (eventuală) de a executa regiunea numai dacă este îndeplinită o anumită condiție.

Fiecărei regiuni critice  $i$  se asociază o *resursă* constând din toate variabilele care trebuie protejate în regiune. Declarația ei se face astfel:

```
resource  $r$  ::  $v_1, v_2, \dots, v_n$ 
```

unde  $r$  este numele resursei, iar  $v_1, \dots, v_n$  sunt numele variabilelor de protejat.

O *regiune critică condiționată* se definește astfel:

```
region  $r$  [ when  $B$  ] do  $S$ 
```

unde  $r$  este numele unei resurse declarate ca mai sus,  $B$  este o expresie booleană, iar  $S$  este secvența de instrucțiuni corespunzătoare regiunii critice. Dacă este prezentă opțiunea *when*, atunci  $S$  este executată numai dacă  $B$  este adevărată. Variabilele din resursa  $r$  pot fi folosite numai de către instrucțiunile din secvența de instrucțiuni  $S$ .

Descrierea prin semafoare a unei regiuni critice condiționate este dată în fig. 9.9. Sunt necesare două semafoare și un număr întreg. Semaforul *sir* reține în coada lui toate procesele care solicită acces la regiune. Semaforul *exclus* asigură accesul exclusiv la anumite secțiuni ale implementării (în special la modificarea numărului întreg *nr*). Întregul *nr* reține câte procese au cerut acces la regiune, la un moment dat. În secțiunea următoare vom da o aplicație interesantă a regiunii critice condiționate.

În fig 9.9 algoritmul propriu-zis începe de la linia 9. Dacă regiunea critică este folosită fără o condiție asociată, se va trece direct la execuția lui  $S$ . Altfel, incrementând variabila *nr*, se marchează faptul că încă un proces accesează regiunea critică cu o condiție asociată. Cum mai multe procese pot face acest lucru simultan, incrementarea trebuie protejată. Deși la prima vedere este doar o singură operație, în realitate incrementarea se traduce într-o secvență de instrucțiuni binare care pot fi executate într-o ordine similară cu cea din fig 9.2 și poate conduce la erori. Ca urmare, înainte ca incrementarea să aibă loc, se execută **P(exclus)** care dacă condiția  $B$  este satisfăcută va proteja și execuția lui  $S$ . Altfel, dacă condiția  $B$  nu este îndeplinită procesul trebuie să intre în așteptare. Aceasta se face prin **P(sir)**, dar nu înainte de a elibera semaforul *exclus* prin **V(exclus)**. Dacă **V(exclus)** este omis, practic orice alt proces care ar accesa regiunea critică se va opri la linia 9. În acest fel concurența pentru

evaluarea condiției  $B$  ar fi. Scopul nostru însă nu este serializarea evaluării condiției  $B$  și a secvenței de instrucțiuni  $S$ , ci serializarea accesului la  $S$  în momentul în care  $B$  este satisfăcută. Ca urmare trebuie să lăsăm fiecare proces să evalueze  $B$  fără restricții. Dacă condiția  $B$  este îndeplinită, se decrementează  $nr$  (marcând astfel intrarea în regiunea critică) și se execută instrucțiunile din  $S$ . Apoi (liniile 20-21), semaforul  $sir$  este incrementat până la zero în așa fel încât controlorul va alege unul din procesele din coadă pentru a fi executat. În final semaforul  $exclus$  este eliberat la linia 22. Procesul proaspăt reactivat va reevalua condiția  $B$  și dacă este îndeplinită va continua și va executa instrucțiunile din  $S$ .

```

1 // Valorile de pornire
2 semaphore sir, exclus;
3 int nr, i, cuConditie;
4 v0(sir) = 0;
5 v0(exclus) = 1;
6 nr = 0;
7
8 //Codul corespunzător regiunii
9 P(exclus);
10 if(cuConditie) {
11     nr++;
12     while(!B) {
13         V(exclus);
14         P(sir);
15         P(exclus);
16     }
17     nr--;
18 }
19 S;
20 for(i=0; i<nr; i++)
21     V(sir);
22 V(exclus);

```

**Figura 9.9 Implementarea unei regiuni critice condiționate**

Până în acest moment lucrurile ar părea clare, dar există o situație în care implementarea de mai jos nu funcționează! Este vorba de cazul în care toate procesele accesează regiunea critică condiționată în momentul în care  $B$  nu este satisfăcută și ca urmare vor intra în așteptare în coada semaforului  $sir$ . Întrebarea este: cum vor ieși aceste procese din coada semaforului? Este necesar ca cineva (un proces) să modifice starea condiției  $B$  și apoi să anunțe schimbarea proceselor care așteaptă. Aceasta poate avea loc în două moduri:

1. După câțva timp un alt proces schimbă starea condiției  $B$ , intră în regiunea critică și constată că  $B$  este satisfăcută. Ca urmare, nu va intra în coada semaforului  $sir$  și după execuție lui  $S$  îi va incrementa valoarea până la zero (liniile 20-21). În acest moment unul dintre procesele blocate în coada semaforului va fi reactivat și va trece prin secțiunea critică. La rândul lui, acesta va reincrementa semaforul până la zero astfel reactivând un nou proces.
2. Unul dintre procesele care folosesc regiunea critică nu o asociază cu condiția  $B$  și ca urmare va omite așteptarea la semaforul  $sir$ . La fel ca și în cazul anterior, după execuția lui  $S$ , acest proces va incrementa semaforul activând procesele din coada lui.

Un program care se bazează pe varianta 1 trebuie să se asigure ca întotdeauna va exista un proces care va găsi condiția  $B$  satisfăcută, altfel totul se blochează. Varianta a doua este mai puțin riscantă pentru că un proces apelând regiunea critică fără condiție nu va rămâne blocat. Pentru a cuprinde ambele într-un singur program, s-a adăugat condiția din linia 10 cu închiderea ei din linia 18.

### 9.1.2.5 Problema citirilor și a scrierilor.

Problema a fost formulată în 1971 de către Courtois, Heymans și Parnas [4] [50]. Se presupune că există două tipuri de procese: *cititor* și *scriitor*. Ele partajează împreună o aceeași resursă, să zicem un fișier. Un proces scriitor modifică conținutul fișierului, iar unul cititor consultă informațiile din el. Spre deosebire de problema producătorului și a consumatorului (unde toate procesele aveau acces exclusiv la resursă), la acest tip de problemă *orice proces scriitor are acces exclusiv la resursă, în timp ce mai multe procese cititor pot avea acces simultan la ea.*

Între cerințele problemei trebuie să se aibă în vedere ca nici un proces să nu fie nevoit să aștepte indefinit. De asemenea, trebuie să se precizeze, pentru cazurile de acces simultan, care tip de proces este mai prioritar. În prima versiune era stabilit ca procesele cititor să fie mai prioritare. Aceasta putea conduce la situația ca procesele scriitor să aștepte indefinit dacă o infinitate de procese cititoare cer accesul. Versiunile ulterioare au dat prioritate proceselor scriitor. Drept urmare, este posibil ca procesele cititor să aștepte indefinit.

În secțiunile precedente am văzut rezolvarea diferitelor probleme de programare concurentă, folosind cel mai des operațiile **P** și **V**. Ele sunt suficiente pentru descrierea acestei probleme și invităm cititorul să o facă. În secțiunea precedentă am definit o construcție mai elegantă și mai ușor de urmărit în programe, regiunea critică condiționată.

O posibilă rezolvare, folosind regiuni critice condiționate, este cea din fig. 9.10. Soluția aparține lui P. B. Hansen, 1973, și ea dă prioritate proceselor scriitor față de cele cititor. Un proces scriitor intră în lucru imediat după ce procesele cititor active și-au încheiat citirile curente. Resursa *f* este zona la care doresc accesul cele două tipuri de procese. Presupunem că ea este un fișier, iar operațiile de bază asupra lui sunt *read* și *write*. Resursa *c* conține două contoare care indică câte procese de fiecare tip au cerut acces la resursa *f*. Este de remarcat folosirea regiunii *c* cu și fără o condiție atașată.

Cititor	Scriitor
<pre> int nw, nr; resource f; resource c :: nr, nw; nr = 0; nw = 0; </pre>	
<pre> region c   when nw == 0     do nr++; &lt;citește date&gt;; region c   do nr--; </pre>	<pre> region c   do nw++; region c   when nr == 0     do &lt;așteaptă să se termine toți cititorii activi&gt;; region f   do &lt;scrie date&gt;; region c   do nw--; </pre>

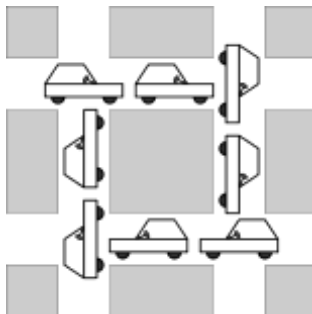
**Figura 9.10 Problema citirilor și a scrierilor**

### 9.1.3 Problema impasului

#### 9.1.3.1 Conceptul de impas.

Noțiunea de impas nu este specifică sistemelor de operare sau informaticii, ci este un fenomen întâlnit în multe situații din viața de zi cu zi. În general, impasul poate fi definit ca o blocare a activității normale a două sau mai multe entități ca urmare a efectului lor coroborat asupra mediului.

Un exemplu de impas ușor de vizualizat este cel al traficului în jurul unui cvartal. După cum se vede în fig 9.11, cele 8 mașini sunt blocate în mod circular, și traficul nu poate continua fără ca una dintre ele să dea înapoi.



**Figura 9.11 Un impas într-o intersecție**

Trecând la domeniul informatic, o modalitate ușoară de a crea o situație de impas este invesarea operațiilor **P** și **V** în procese concurente. Fig 9.12 prezintă codul celor două procese.

Proces A	Proces B
semaphore x, y; v0(x) = 1; v0(y) = 1;	
P(x); <instrucțiuni A>; P(y); V(y); V(x);	P(y); <instrucțiuni B>; P(x); V(x); V(y);

**Figura 9.12 Impas provocat de două semafoare**

Dacă instrucțiunile celor două procese sunt executate întâmplător în secvența din figura 9.13, cele două procese **A** și **B** se vor afla într-o stare de impas. Cauza este oarecum similară cu situația traficului prezentată anterior. Procesul **A** va deține semaforul binar **x**, procesul **B** va deține semaforul binar **y**, apoi ambele vor încerca să obțină semaforul deținut de celălalt, fără a elibera semaforul pe care îl deține înainte de a-l obține pe cel deținut de celălalt proces. Situația este ilustrată în fig. 9.13.

Proces A	Proces B
P(x);	
	P(y);
<instrucțiuni A>;	
	<instrucțiuni B>;
P(y);	
wait	P(x)
...	wait...

**Figura 9.13 O situație de impas**

Cauza principală care a generat această situație este faptul că procesele își ocupă resursele numai atunci când au nevoie efectiv de ele și ca urmare apare situații de așteptare circulară.

Apropiindu-ne mai mult de funcționarea unui sistem de operare, să presupunem că în sistem sunt numai două benzi magnetice **MM0** și **MM1**. Dacă un proces cere mai întâi **MM0** și apoi

**MM1** fără a elibera **MM0**, iar al doilea cere **MM1** și apoi **MM0** fără a elibera **MM0**, poate apare o astfel de așteptare circulară. Acest fenomen este cunoscut în literatură sub mai multe denumiri, și anume: *impas*, *interblocare*, *deadlock*, *deadly embrace* etc. [4], [10], [51].

Ultimele două exemple de *impas* au avut ca protagoniști câte două procese și două resurse. Impasul poate apărea însă în situații mai complexe care implică nu doar două procese, ci mai multe și nu neapărat resurse specifice, ci clase de resurse. Să presupunem că într-un anumit moment din funcționare, un sistem de operare dispune de 50MB de memorie liberă și de 10 poziții libere pentru crearea de conexiuni în rețea. În acest moment un proces **A** cere alocarea a 30MB de memorie și apoi a 7 conexiuni în rețea, iar un proces **B** cere întâi alocarea a 5 conexiuni de rețea și apoi alocarea a 40MB de memorie. Similar cu exemplul anterior, este posibil ca ambelor procese să li se satisfacă primele cereri, însă după aceea din lipsă de resurse va apare *impas*. Acest tip de *impas* este mai greu de detectat pentru că nu implică resurse punctuale ci clase de resurse, în cazul nostru memoria RAM și conexiunile la rețea.

Impasul rezultat din conflictul asupra câtorva resurse punctuale poate fi modelat teoretic pe baza *impasului* pe clase, considerând că fiecare clasă de elemente conține un singur element. Altfel spus, *impasul* pe clase este o generalizare a *impasului* pe resurse punctuale.

Impasul este o stare gravă care poate duce la un blocaj al sistemului de operare sau la distrugerea unor procese. În cazul în care două procese se blochează unul pe altul, doar ele vor fi în *impas* și situația se poate salva terminând forțat unul dintre ele. Dacă însă *impasul* se manifestă între un proces și sistemul de operare sau direct în interiorul sistemului de operare, parțial sistemul "îngheață" și ieșirea cea mai frecventă din asemenea situații este restartarea calculatorului!

Astfel de situații trebuie avute în vedere și prevenite de către proiectanții de sisteme de operare. Dar, după cum vom vedea, tratarea *impasului* este costisitoare și ca urmare proiectanții o tratează doar în situații izolate pentru a nu afecta performanța sistemului.

În legătură cu *impasul*, fiecare sistem de operare trebuie să dea rezolvări măcar la una dintre următoarele probleme:

- ieșirea din *impas*,
- detectarea unui *impas*,
- evitarea (prevenirea apariției) *impasului*.

Cele două rezultate de mai sus dau condiții prin care este indicată eventuala apariție a unui *impas*, însă nu furnizează metode practice de rezolvare a problemelor de mai sus. În continuare vom da metode și tehnici de rezolvare a problemelor generate de *impas*.

### 9.1.3.2 Ieșirea din *impas*.

Soluțiile ieșirii din *impas* pot fi clasificate în manuale și automate. Soluțiile manuale sunt cele care implică acțiuni ale operatorului uman. Cele automate sunt executate de către sistemul de operare.

Ieșirea din starea de *impas* automată (de către sistemul de operare) presupune în primul rând detectarea acestei stări. La rândul ei, detectarea stării de *impas* implică existența unei

componente a sistemului de operare care să întrețină și să analizeze starea alocărilor de resurse din sistem (vom discuta acest aspect mai jos). Pe baza acestei analize, sistemul de operare recunoaște starea de impas și trebuie să aleagă o soluție pentru a-l rezolva.

Indiferent dacă este automată sau manuală, soluția ieșirii din impas necesită acțiuni drastice. În cele ce urmează prezentăm trei astfel de soluții.

O primă variantă este *reîncărcarea sistemului de operare*. Soluția nu prea este de dorit, dar este cea mai la îndemâna operatorului. Dacă pierderile nu sunt prea mari se poate adopta și o astfel de strategie.

O a doua variantă este *alegerea un proces "victimă"*, care este fie cel care a provocat impasul, fie un altul de importanță mai mică, astfel încât să fie înlăturat impasul. Se distruge acest proces și odată cu el toți descendenții lui. Această alegere este o parte importantă a algoritmilor automați de rezolvare a impasului. În cazul abordării manuale, din păcate, criteriul de alegere a "victimei" este de multe ori departe de a fi obiectiv.

O a treia soluție este *creerea unui "punct de reluare"*, care este o fotografie a memoriei pentru procesul victimă și pentru procesele cu care el colaborează. Apoi se distruge procesul victimă și subordonații lui. Procesul victimă se reia ulterior. Crearea punctelor de reluare reclamă consum de timp și complicații, motiv pentru care de multe ori nu se execută.

### 9.1.3.3 Detectarea impasului

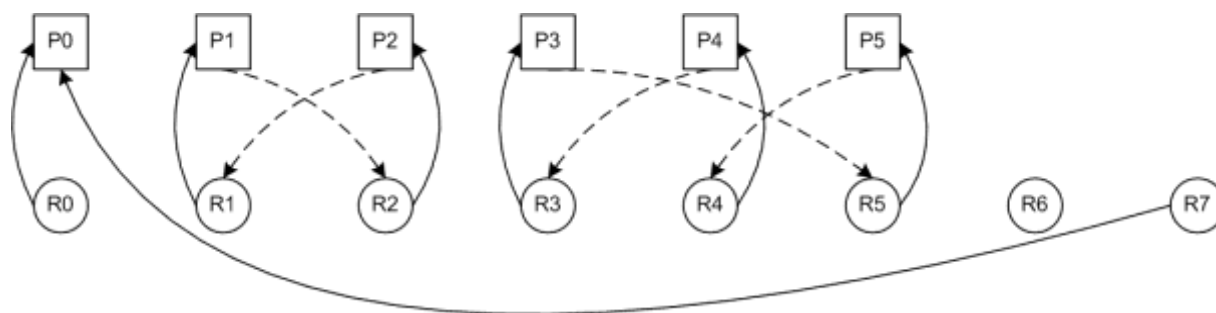
*Detectarea unui impas* se face atunci când **SO** nu are un mecanism de prevenire a impasului. Pentru a se reuși detecția este necesar ca **SO** să aibă o evidență clară, pentru fiecare proces, privind resursele ocupate, precum și cele solicitate dar neprimite. Pentru aceasta, cel mai potrivit model este *graful alocării resurselor*.

Facem ipoteza, că fiecare resursă este într-un singur exemplar și vom nota  $R1, R2, \dots, R_m$  aceste resurse. Fie acum  $P1, P2, \dots, P_n$  procesele din sistem. Se construiește un graf bipartit  $(X, U)$  astfel:

- $X = \{P1, P2, \dots, P_n, R1, R2, \dots, R_m\}$
- $(R_j, P_i)$  este arc în  $U$  dacă procesul  $P_i$  a ocupat resursa  $R_j$ .
- $(P_i, R_j)$  este arc în  $U$  dacă procesul  $P_i$  așteaptă să ocupe resursa  $R_j$ .

Dacă graful  $(X, U)$  definit mai sus este ciclic, atunci sistemul se află în stare de impas. Demonstrația este un exercițiu simplu de inducție. Detectarea ciclicității unui graf se face folosind algoritmi clasici din teoria grafelor.

În fig 9.14 prezentăm un posibil graf de alocare. Nodurile pătrate reprezintă procesele iar nodurile rotunde reprezintă resursele. Pentru o mai ușoară diferențiere arcele de așteptare pentru o resursă au fost desenate cu linie întreruptă.



**Figura 9.14** Un graf de alocare a resurselor

Dacă selectăm numai subgraful cu nodurile  $\{P1, P2, R1, R2\}$ , iar pentru resursele R1 și R2 considerăm unitățile de bandă MM0 și MM1, atunci se obține graful de alocare pentru unul dintre exemplele prezentate în 9.1.3.1.

Deși nu sunt evidente la prima vedere, graful din figura de mai sus conține două cicluri. Primul ciclu implică procesele P1 și P2, și resursele R1 și R2. Al doilea ciclu implică procesele P3, P4 și P5, și resursele R3, R4, și R5. În ambele cazuri aveam de-a face cu stări de impas.

#### 9.1.3.4 Evitarea (prevenirea apariției) impasului.

În 1971, Coffman, Elphic și Shoshani [4], [19] au indicat patru condiții *necesare* pentru apariția impasului:

1. procesele solicită controlul exclusiv asupra resurselor pe care le cer (condiția de *excludere mutuală*);
2. procesele păstrează resursele deja ocupate atunci când așteaptă alocarea altor resurse (condiția de *ocupă și așteaptă*);
3. resursele nu pot fi șterse din procesele care le țin ocupate, până când ele nu sunt utilizate complet (condiția de *nepreempție*);
4. există un lanț de procese în care fiecare dintre ele așteaptă după o resursă ocupată de altul din lanț (condiția de *așteptare circulară*).

Evitarea impasului presupune împiedicarea apariției uneia dintre aceste patru condiții. Vom lua fiecare condiție pe rând și vom analiza modalitățile în care ea poate fi împiedicată și efectele asupra funcționării sistemului

**Condiția 1 – excluderea mutuală.** În absența excluderii mutuale între procese, dispare noțiunea de așteptare după o resursă critică blocată de alt proces. Dispărând noțiunea de așteptare, evident dispare și noțiunea de impas din moment ce nici un proces nu se va opri din execuție. Eliminarea acestei condiții însă nu este o alegere fericită. După cum am demonstrat anterior, excluderea mutuală este esențială pentru pastrarea corectitudinii datelor. Practic impasul a apărut ca și un efect secundar al aplicării acestei condiții indiscutabil necesare.

**Condiția 2 – ocupă și așteaptă.** Această condiție spune că un proces poate aștepta după o resursă în timp ce deține (blochează) alte resurse. Eliminarea acestei condiții poate fi făcută în două moduri:

- Procesul trebuie să elibereze toate resursele blocate înainte de a solicita o nouă resursă.
- Procesul trebuie să blocheze toate resursele de care are nevoie la pornire.



Prima variantă este destul de forțată pentru că cere procesului să concureze pentru toate resursele ori de câte ori ar avea nevoie de una nouă.

A doua variantă nu este aplicabilă sistemelor de operare interactive (Windows, Unix etc.) pentru că ar restrânge posibilitățile de utilizare a sistemului. Să luăm spre exemplu cazul comun al unui editor de text care suportă deschiderea simultană a mai multor fișiere. Dacă sistemul de operare ar impune blocarea tuturor resurselor la încărcarea procesului, editorul nostru ar fi obligat să solicite o cantitate finită de memorie pe care să o folosească la încărcarea fișierelor editate. Efectul acestui comportament are două tăișuri:

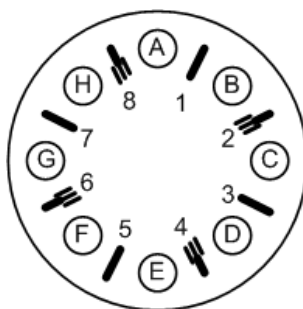
- Dacă editorul deschide un singur fișier mic, practic majoritatea memoriei solicitate va fi ținută blocată fără motiv. Evident ar fi mult mai util ca sistemul de operare să dispună de această memorie pentru a servi alte procese.
- Dacă editorul trebuie să deschidă multe fișiere de dimensiuni mari, va ajunge să ocupe toată memoria alocată și să fie nevoit să refuze deschiderea altor fișiere.

**Condiția 3 – nepreempție.** Eliminarea condiției de nepreempție permite sistemului de operare să aleagă procesele care ar trebui oprite pentru a ieși din impas și chiar să le oprească. Evident simpla terminare a unui proces nu este acceptabilă și ca urmare sistemul trebuie să și repornească acel proces. Deși această abordare rezolvă situația de impas, nu garantează că procesele vor funcționa conform așteptărilor. Teoretic este posibil ca după restartarea unui proces dintr-un grup implicat într-un impas, procesele rămase să nu progreseze suficient pentru a ieși din regiunea critică înainte ca programul repornit să reajungă acolo. În acel moment impasul va reapare și din nou unul dintre procese va fi repornit. Această secvență se poate repeta la nesfârșit și deși nici unul dintre procese nu este blocat, totuși niciunul nu progresează cu execuția. Acest fenomen este cunoscut sub numele de *impas activ* (livelock).

**Condiția 4 – așteptarea circulară.** O soluție simplă pentru împiedicarea așteptării circulare este impunerea unei ordini de blocare a resurselor. Sistemul de operare decide o ordine absolută a tuturor resurselor și nici un proces nu are voie să blocheze o resursă cu număr de ordine mai mare înaintea uneia cu un număr de ordine mai mic. Pentru a vedea efectiv efectul blocării ordonate, vom prezenta o *problema filozofilor*.

Să considerăm o masă circulară la care sunt așezați opt filozofi. Pe masă, în loc să fie șaisprezece tacâmuri (furculițe și cuțite) sunt doar opt, așezate în așa fel încât lângă fiecare farfurie este un cuțit și o furculiță (vezi fig 9.15). Fiecare filozof are nevoie de două tacâmuri pentru a mânca, și odată ce a luat un tacâm de pe masă nu-l va elibera până când nu termină de mâncat. Dacă toți filozofii iau de pe masă tacâmul din stânga, ne aflăm se creează o situație de impas, pentru că nici un filozof nu va avea două tacâmuri pentru a putea mânca.

Impunerea blocării resurselor (luării tacâmurilor) în ordine, va preveni impasul pentru că filozoful H nu va putea ridica tacâmul 8 înaintea tacâmului 7. În consecință filozoful G va avea la dispoziție două tacâmuri pentru a mânca. În momentul în care G va termina, va elibera tacâmul 6 și filozoful F va putea mânca și el. Procesul se va continua astfel până când toți filozofii vor fi mâncat.



**Figura 9.15 Problema filozofilor**

În realitate însă, această soluție este dificil de implementat pentru programatorii unui sistem de operare. Având în vedere cantitatea mare de resurse dintr-un sistem (de ordinul miilor) respectarea unei ordini prestabilite devine practic imposibilă. În plus, ordinea resurselor va fi de cele mai multe ori diferită de ordinea logică în care un program are nevoie de resurse, cauzând astfel scrierea de algoritmi mult mai complecși sau mai puțin performanți.

#### **9.1.3.5 Alocarea controlată (conservativă) de resurse; algoritmul bancherului.**

Evitarea impasului se poate face și impunând o modalitate restrictivă de alocare de resurse. În esență, înainte de a servi o cerere de resurse, sistemul de operare analizează alocările de resurse. Dacă servirea cererii conduce la o stare cu potențial de a genera un impas, cererea de resurse este amânată și reanalizată în momentul în care are loc o modificare în starea alocării de resurse în sistem. Soluțiile pe care le vom discuta vor adresa cazul general al impasului în care lucrăm cu clase de resurse.

Pentru a evalua dacă o stare a alocărilor poate conduce la impas este necesară cunoașterea pentru fiecare proces a necesităților maxime de resurse din fiecare clasă. Astfel, la încărcare, fiecare proces va transmite sistemului de operare cantitățile de resurse din fiecare clasă de care va avea nevoie în cel mai rău caz.

**Soluție ineficientă.** Folosind aceste maxime declarate de fiecare proces la încărcare, este foarte simplu să evităm impasul. Dacă adunăm maximum declarat de un nou proces pentru o resursă la cantitatea din acea resursă deja alocată și suma depășește maximumul disponibil în sistem, pur și simplu punem procesul în așteptare până când se vor elibera suficiente resurse.

Aspectul pozitiv al acestei soluții este că garantează servirea imediată a tuturor cererilor de resurse venite din partea proceselor încărcate. Acest lucru este posibil pentru că algoritmul menține întotdeauna un număr redus de procese încărcate pentru a avea resurse suficiente pentru ele.

Soluția este ineficientă pentru că este foarte puțin probabil ca un proces să folosească imediat după încărcare maximumul declarat de resurse. În general resursele sunt folosite treptat și rareori se va atinge maximumul. Ca urmare, în majoritatea covârșitoare a cazurilor un proces poate să-și înceapă lucrul fără ca sistemul de operare să poată acoperi întreaga plajă de resurse pe care o poate solicita. Algoritmul bancherului evită impasul fără a impune condiții atât de restrictive.

**Algoritmul Bancherului** [4], [19] datorat lui Dijkstra este probabil cel mai renumit pentru alocarea controlată. Spre diferență de soluția precedentă algoritmul bancherului permite tuturor proceselor să se încarce și să înceapă execuția. În schimb, algoritmul nu garantează servirea imediată a tuturor cererilor proceselor încărcate (un lucru de altfel rezonabil). Tot ce

garantează algoritmul bancherului este că orice proces încărcat poate fi servit (nu va rămâne blocat la nesfârșit) cu condiția să aștepte până când se mai eliberează resurse.

Pentru a clarifica numele algoritmului, funcționarea algoritmului este similară cu procedura prin care un bancher (sistemul de operare) decide când și cât (resurse) să împrumute unui client (procesul) dintr-o sumă maximă cu care îl poate credita. Diferențele față de lumea reală sunt:

1. Bancherul dispune de o sumă limitată de împrumut, mai mică decât suma creditelor maxime ale clienților
2. Bancherul nu percepe nici dobândă și nici comision
3. Toți clienții returnează banii în cele din urmă

Înainte de a prezenta algoritmul propriu-zis trebuie să prezentăm structurile de date pe care le vom folosi.

- `resurseExistente`
  - Cantitățile de resurse existente în sistem pentru fiecare clasă de resurse.
  - Tabel unidimensional. Fiecărei clase de resurse îi corespunde o poziție
- `maxResursePerProces`
  - Cantitățile maxime de resurse pe care fiecare proces le poate solicita
  - Tabel bi-dimensional. Fiecare linie corespunde unui proces. Fiecare celulă dintr-o linie conține cantitatea corespunzătoare unei clase de resurse.
- `resurseAlocatePerProces`
  - Cantitățile de resurse alocate fiecărui proces
  - Tabel bi-dimensional. Fiecare linie corespunde unui proces. Fiecare celulă dintr-o linie conține cantitatea corespunzătoare unei clase de resurse.
- `solicitant`
  - Procesul care solicită resurse
- `cerere`
  - Cererea de analizat emisă de procesul solicitant
  - Tabel unidimensional. Fiecărei clase de resurse îi corespunde o poziție
- `resurseDisponibile`
  - Resursele disponibile în sistem
  - Tabel unidimensional. Fiecărei clase de resurse îi corespunde o poziție
- `proceseImplicate`
  - Tablou pentru marcarea proceselor implicate în algoritm
  - Fiecărui proces îi corespunde o poziție

Folosind notațiile de mai sus, algoritmul bancherului este prezentat în limbaj pseudocod în fig 9.16.

```
<Fă o copie tabelului "resurseAlocatePerProces" în tabelul "alocari">;

<Dacă nici un proces nu mai este implicat în algoritm, cererea poate
fi servită. Sfârșit algoritm>;

<Populează tabelul "resurseDisponibile" cu diferențele
dintre "resurseExistente" și sumele resurselor alocate fiecărui
proces în tabelul "alocari">;

<Găsește un proces Pi pentru care sunt suficiente resurse pentru ai
servi cererea maximă>;

<Dacă nu există nici un astfel de proces, înseamnă că cererea inițială
```

```
nu poate fi servită și procesul solicitant este pus în
așteptare. Sfârșit algoritm>;
```

```
<Simulează terminarea procesului Pi scoțându-l din lista proceselor
implicate și eliberându-i resursele (populează cu 0 intrarea lui Pi
din tabelul "alocari")>;
```

```
<Mergi la pasul 2>;
```

### Figura 9.16 Algoritmul bancherului

În esență, algoritmul bancherului consideră cererea servită și apoi verifică dacă există o ieșire din cel mai rău caz posibil, adică atunci când toate procesele își cer maximum de resurse. Aceasta se face găsimd rând pe rând câte un proces căruia i se poate servi maximum de resurse solicitat, și apoi simulându-i terminarea. Dacă folosind acest procedeu, se ajunge la terminarea tuturor proceselor, înseamnă că servirea cererii nu implică riscul unui impas.

În cele ce urmează vom oferi o implementare a algoritmului bancherului. Pentru simplitate vom considera că sistemul de operare dispune de patru tipuri de resurse și că asupra acestor resurse vor acționa cinci procese. Implementarea conține și date de test pentru care algoritmul poate fi verificat.

```
#include<stdio.h>

#define PROCESE 5
#define RESURSE 4

// Date initiale care descriu alocarea de resurse din sistem
int resurseExistente[PROCESE] = {7, 7, 6, 5};

int maxResursePerProces[PROCESE][RESURSE] = {{3, 2, 4, 2}, // P0
                                              {3, 2, 3, 1}, // P1
                                              {4, 1, 2, 3}, // P2
                                              {3, 3, 4, 4}, // P3
                                              {6, 5, 5, 1}}; // P4

int resurseAlocatePerProces[PROCESE][RESURSE] = {{1, 1, 0, 1}, // P0
                                                  {1, 0, 1, 0}, // P1
                                                  {1, 0, 1, 1}, // P2
                                                  {0, 1, 0, 2}, // P3
                                                  {0, 1, 1, 0}}; // P4

void tiparesteStare(int procesAles,
                   int proceseImplicate[PROCESE],
                   int resurseDisponibile[RESURSE],
                   int alocairi[PROCESE][RESURSE]) {
    int i, j;

    printf("\nProcese implicate: ");
    for(i=0; i<PROCESE; i++) {
        if(proceseImplicate[i] == 1) {
            printf("P%d ", i);
        }
    }

    printf("\nProces ales: P%d\n", procesAles);

    printf("Cereri maxime de resurse ale lui P%d:", procesAles);
    for(j=0; j<RESURSE; j++) {
        printf(" %d", maxResursePerProces[procesAles][j]);
    }
}
```

```

    }

    printf("\nResurse disponibile:");
    for(j=0; j<RESURSE; j++) {
        printf("  %d", resurseDisponibile[j]);
    }

    printf("\nResurse alocate lui P%d:", procesAles);
    for(j=0; j<RESURSE; j++) {
        printf("  %d", alocari[procesAles][j]);
    }
    printf("\n");
}

int verificaCerere(int solicitant, int *cerere) {
    // Variabile de lucru
    int resurseDisponibile[RESURSE];
    int proceseImplicate[PROCESE] = {1, 1, 1, 1, 1};
    int alocari[PROCESE][RESURSE];
    int i, j, procesAles, gasit, n;

    // Copiem tabela de alocariare de resurse la procese si adaugam la copie
    // noua cerere
    for(i=0; i<PROCESE; i++) {
        for(j=0; j<RESURSE; j++) {
            alocari[i][j] = resurseAlocatePerProces[i][j];
            if(i == solicitant) {
                alocari[i][j] += cerere[j];

                // Procesul cere mai multe resurse decat maximul permis
                if(alocari[i][j] > maxResursePerProces[i][j])
                    return -1; //Sfarsit algoritm
            }
        }
    }

    while(1) {
        // Calculam cantitatile de resurse disponibile
        for(j=0; j<RESURSE; j++) {
            resurseDisponibile[j] = resurseExistente[j];
            for(i=0; i<PROCESE; i++) {
                resurseDisponibile[j] -= alocari[i][j];

                // Insuficiente resurse disponibile. Starea curenta conduce la impas
                if(resurseDisponibile[j] < 0)
                    return 0; // Sfarsit algoritm
            }
        }

        // Calculam numarul de procese inca implicate in algoritm
        n = 0;
        for(i=0; i<PROCESE; i++) {
            n += proceseImplicate[i];
        }

        // Toate procesele au trecut cu bine prin analiza. Cererea va fi servita
        if(n == 0)
            return 1; // sfarsit algoritm

        // Cautam un proces pentru care avem suficiente resurse pentru a-i
        // servi cererea maxima
        procesAles = -1;
        for(i=0; i<PROCESE; i++) {

```

```

    if(proceseImplicate[i] != 1)
        continue;

    gasit = 1;
    for(j=0; j<RESURSE; j++) {
        if(maxResursePerProces[i][j]-alocari[i][j] > resurseDisponibile[j]) {
            gasit = 0;
            break;
        }
    }
    if(gasit == 1) {
        procesAles = i;
        break;
    }
}

// Nici unul dintre procesele implicate in algoritm nu poate fi
// servit cu maximul posibil, deci exista posibilitatea unui impas
if(procesAles == -1)
    return 0; // sfarsit algoritm

tiparesteStare(procesAles,proceseImplicate, resurseDisponibile, alocairi);

// Simulam faptul ca procesul procesAles ajunge sa se termine cu bine,
// adica resursele ii sunt dealocate si este marcat ca neimplicat in
// algoritm
proceseImplicate[procesAles] = 0;
for(j=0; j<RESURSE; j++)
    alocairi[procesAles][j] = 0;
}

int main() {
    int cerere[RESURSE] = {1, 2, 2, 1};
    int solicitant = 1;

    switch(verificaCerere(solicitant, cerere)) {
        case -1:
            printf("\n\nEROARE. Procesul cere mai mult decat maximul declarat.\n");
            break;
        case 0:
            printf("\n\nCererea nu poate fi servita fara risc de impas. "
                "Procesul va fi pus in asteptare.\n");
            break;
        case 1:
            printf("\n\nCererea poate fi servita fara risc de impas.\n");
            break;
    }

    return 0;
}

```

**Figura 9.17 Implementarea algoritmului bancherului**

Pentru datele de test din figura 9.17 algoritmul decide că cererea poate fi servită, eliminând toate procesele în următoarea secvență: P1, P0, P2, P3, P4. Dacă însă în tabelul `resurseAlocatePerProces` se modifică alocarea resursei R2 pentru procesul P2 din 1 în 2, algoritmul nu va servi cererea din cauza apariției riscului de impas. Lăsăm pe seama cititorului să ruleze programul din fig. 9.17 și să vadă rezultatele obținute.

## 9.2 Conceptul de multiprogramare.

Cel mai important concept, introdus mai întâi la sistemele seriale și preluat apoi la toate sistemele de calcul moderne, este conceptul de *multiprogramare*. El reprezintă modul de exploatare a unui sistem de calcul cu un singur procesor central, care presupune existența simultană în memoria internă a mai multor procese care se execută concurrent. Multiprogramarea duce la o mai bună utilizare a procesorului și a memoriei.

Primele sisteme de calcul ce lucrează în multiprogramare au apărut cu ceva mai înainte de anul 1965. Tehnologic, în aceeași perioadă au apărut plăcile cu circuite integrate. De acum se poate vorbi de generația a III-a de calculatoare.

Pe scurt, lucrul în multiprogramare se desfășoară astfel: în fiecare moment procesorul execută o instrucțiune a unui proces. Vom spune că acest proces este în starea **RUN**. Restul proceselor, fie că așteaptă apariția unui eveniment extern (terminarea unei operații I/O, scurgerea unui interval de timp etc.), fie că sunt pregătite pentru a fi servite, în orice moment, de către procesor. Despre cele care așteaptă un eveniment extern spunem că sunt în starea **WAIT**, iar cele care sunt pregătite de execuție spunem că sunt în starea **READY**.

Pentru ca un sistem de operare să poată lucra în multiprogramare, este necesar ca:

Sistemul de calcul să dispună de un sistem de întreruperi pentru a semnaliza apariția evenimentelor externe;

Sistemul de operare să gestioneze, să aloce și să protejeze resursele între utilizatori. Prin resurse înțelegem: memoria, dispozitivele periferice, fișierele, timpul fizic etc.

### 9.2.1 Trecerea unui proces dintr-o stare într-alta

Trecerea unui proces din starea **RUN** în starea **WAIT** este (evident) comandată de către procesul însuși; el (procesul) știe, în funcție de rezultatele obținute până în prezent, când trebuie să urmeze o instrucțiune de I/O.

Trecerea proceselor din starea **RUN** în starea **READY** și invers poate fi comandată sau de proces (cedare voluntară de procesor) sau poate fi comandată de către sistemul de operare (cedare involuntară de procesor).

Trecere unui proces dintr-o stare într-alta este făcută pe baza unui *algoritm de planificare* rulat de o componentă a sistemului de operare numită *planificator*. Planificatorul este o componentă centrală a oricărui sistem de operare modern.

#### 9.2.1.1 Cedarea voluntară a procesorului

Acceastă modalitate de cedare a procesorului se bazează în general pe programatorul care dezvoltă programul să introducă instrucțiuni de cedare în codul sursă. În alte situații sistemul de calcul în sine introduce instrucțiuni de cedare printre instrucțiunile binare ale programului. Un planificator care folosește această metodă se numește *nepreemptiv*.

Cedarea voluntară a procesorului este o practică din vremurile de început ale sistemelor de operare. Deși simplifică mult munca planificatorului, acest gen de cedare dă unui proces posibilitatea de a bloca întregul sistem de operare (o problemă foarte serioasă care se poate

remedia doar prin re-startarea sistemului de calcul). Pentru a bloca sistemul, un proces nu trebuie decât să ruleze un ciclu infinit în care să nu facă nici o cerere de resurse care să-l întrerupă și nici să nu apeleze instrucțiuni de cedare a procesorului.

### 9.2.1.2 Cedarea involuntară a procesorului

Cedarea involuntară a procesorului presupune intervenția planificatorului pentru a întrerupe procesul curent. Un astfel de planificator se numește *preemptiv*. Întreruperea procesului curent se face în general folosind o întrerupere activată de sistemul de calcul la intervale constante (în general) de timp. Procedura de servire a întreruperii este aceeași pentru toate procesele și cheamă instrucțiunea de cedare a procesorului. În urma acestei cedări, planificatorul sistemului având o prioritate absolută, preia procesorul (devine activ) și decide care dintre procesele existente va primi procesorul.

### 9.2.2 Funcționarea unui planificator

În fig. 9.18 este prezentat algoritmul după care funcționează un planificator. În el se consideră că sunt  $n$  procese, iar prioritățile lor sunt numerele 1, 2, ...,  $n$ , cu 1 prioritatea ca mai mică și  $n$  prioritatea ca mai mare.

Fiecare proces este încărcat într-o *partiție* de memorie. Fiecare partiție are un *număr de prioritate*. În fig. 9.19 se prezintă o posibilă evoluție a trei procese în regim de multiprogramare. Fiecare dintre cele trei procese este încărcat într-o partiție proprie.

Procesul **P3**, are nevoie să efectueze până la terminare două operații I/O și are cea mai mare prioritate dintre cele trei procese. Procesul **P2** are nevoie de trei operații I/O. Procesul **P1**, cel mai puțin prioritar, are nevoie să efectueze până la terminare o singură operație I/O. Presupunem că cele trei procese intră în sistem în același timp. Desenul ilustrează pe orizontală axa timpului, iar porțiunile în care este scris "**RUN**" indică faptul că procesorul execută succesiv instrucțiuni mașină ale procesului respectiv. În funcție de momentele în timp în care cele trei procese solicită operații de I/O, fig. 9.19 ilustrează 12 momente eveniment, momente în care se schimbă starea unora dintre cele trei procese.

```
do {
  for(i=n; i>=1; i--) {
    if(<Pi este în starea RUN>) {
      EXECUTA(<urmatoarea instructiune masina din Pi>);
      continue;
    } else if(<Pi este în starea READY>) {
      // Exista, în mod sigur, cel puțin un proces Pj cu i>j în starea RUN
      STARE(Pj) = READY;
      STARE(Pi) = RUN;
      EXECUTA(<urmatoarea instructiune masina din Pi>);
      continue;
    }
  }
} while(true);
```

**Figura 9.18 Implementarea unui planificator**

Observăm că **P1** din cauza priorității mici, reușește să intre în starea **RUN** numai când **P2** și **P3** sunt în starea **WAIT** sau **FINISH** (procesul s-a terminat). Acestea sunt momentele 2, 6, 8 și 11. Să mai observăm, de asemenea, că sistemul de calcul este bine exploatat, singura perioadă de timp în care procesorul este neîncărcat este între momentele 7 și 8. Acesta este de fapt scopul principal al multi programării: valorificarea procesorului de către alt proces câtă vreme cel curent așteaptă un eveniment extern.



P3	RUN	WAIT			RUN	WAIT			RUN				
P2	READY	RUN	WAIT	RUN	WAIT	RUN	WAIT			RUN			
P1	READY		RUN	READY		RUN	WAIT	RUN	READY		RUN		
	0	1	2	3	4	5	6	7	8	9	10	11	12

**Figura 9.19 Exemplu de evoluție în multiprogramare**

Evoluțiile din fig. 9.19 descriu o situație particulară. Astfel, toate procesele sunt lansate în execuție simultan, iar după terminarea lui **P1** și **P2** nu se mai lansează nimic. În realitate, fiecare partiție primește spre execuție un nou proces (dacă există) după terminarea celui curent. Deci, spre exemplu, la momentul 10 în partiția 1 se va lansa un nou proces.

### 9.3 Planificarea proceselor

#### 9.3.1 Sarcinile planificatorului de procese.

Planificatorul de procese este responsabil pentru trecerea proceselor din starea **READY** în starea **RUN** și invers. Pentru aceasta planificatorul trebuie să îndeplinească următoarele sub-sarcini:

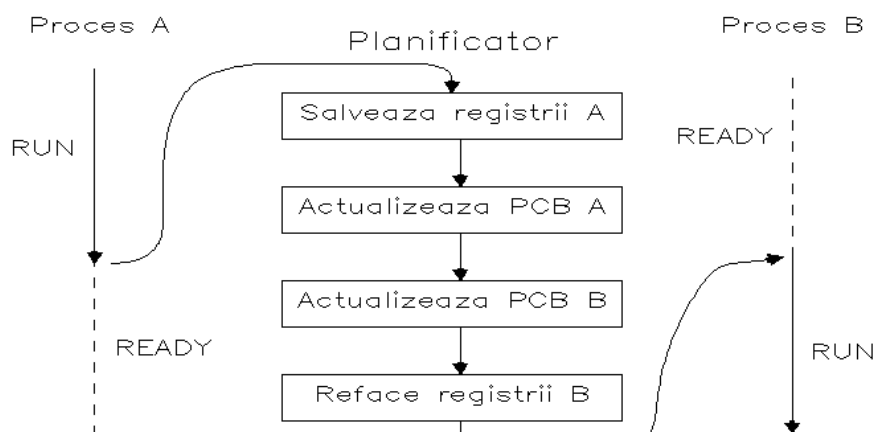
- ținerea evidenței tuturor proceselor din sistem;
- alegerea procesului căruia i se va atribui procesorul și pentru cât timp;
- alocarea procesului un procesor;
- eliberarea procesorului la ieșirea procesului din starea **RUN**.

Pentru fiecare proces din sistem planificatorul întreține o structură de date numită Process Control Block (**PCB**) în care trebuie să apară, printre altele, următoarele informații:

- starea procesului;
- pointer la următorul **PCB** în aceeași stare;
- numărul procesului (acordat de planificator);
- contorul de proces (adresa instrucțiunii mașină care urmează a fi executată);
- zonă pentru copia regiștrilor generali ai mașinii;
- limitele zonei (zonelor) de memorie alocate procesului;
- lista fișierelor deschise etc.

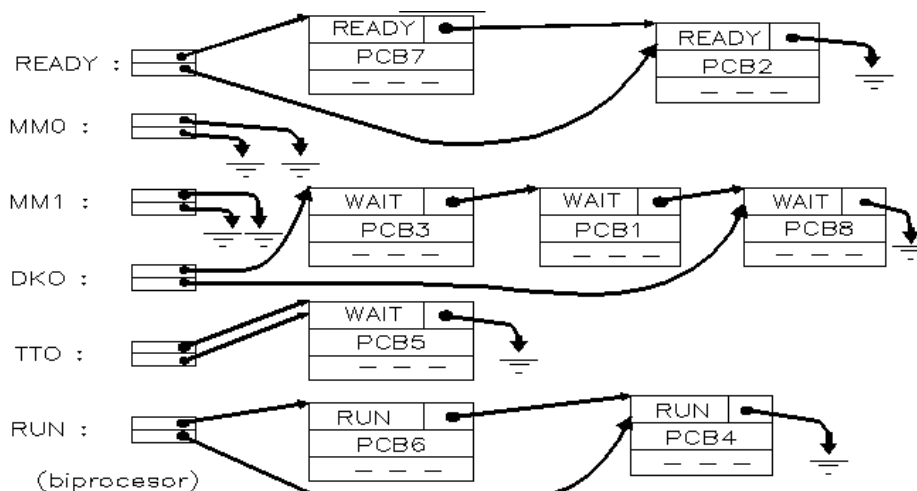
Toate structurile **PCB** sunt memorate într-o zonă de memorie proprie planificatorului de procese. Starea unui proces poate fi **RUN**, **READY**, **WAIT** sau **SWAP** (vezi 8.2.1.1). Procesele din aceeași stare sunt înlănțuite între ele.

Deoarece toate procesele folosesc regiștrii generali ai procesorului, este necesar ca la ieșirea din starea **RUN** conținutul acestora să fie salvat. La o nouă intrare în **RUN** se reface conținutul regiștrilor. Spre exemplu, în fig. 9.20 este prezentat un schimb **READY** ↔ **RUN** între două procese **A** și **B**. Evident, la trecerea din **B** în **A** operațiile se inversează.



**Figura 9.20 Schimbul RUN ↔ READY între două procese**

În sarcina planificatorului de procese stă și organizarea cozilor la perifericele de I/O. Deoarece mai multe procese pot solicita simultan accesul la un periferic, planificatorul întreține cozi de așteptare pentru toate perifericele. În fig. 9.21 este exemplificat un ansamblu de cozi proprii planificatorului de procese. Cozile **MM0**, **MM1**, **DKO** și **TTO** sunt cozi de așteptare asociate cu perifericele cu același nume. Coada **READY** conține procese aflate în starea **READY** (pregătite pentru a intra în execuție în orice moment). Coada **RUN** conține procese aflate în execuție deja. Cum un singur proces poate fi în execuție pe un procesor, se deduce că fig. 9.21 se referă la un sistem de calcul bi-procesor.



**Figura 9.21 Cozi gestionate de planificatorul proceselor**

### 9.3.2 Algoritmi de planificare.

O mare parte a literaturii a pus un accent pe algoritmi de planificare. Iată numai câteva lucrări de referință în domeniul sistemelor de operare, care tratează pe larg această problemă: [4], [10], [45], [49]. Motivele principale pentru care s-a dat importanță mare planificării sunt că performanța unui sistem de operare este influențată decisiv de funcționarea planificatorului și nu în ultimul rând pentru că planificare se pretează ușor la modelarea matematică.

Alte lucrări, cum ar fi [22], [29], enumeră foarte pe scurt câțiva algoritmi, trecând repede peste acest capitol. Punctul de vedere al celor din urmă este justificat prin faptul că *ponderea planificării în sistemele de operare este destul de mică*. În cele ce urmează, enumerăm cei mai utilizați algoritmi de planificare a proceselor.

### 9.3.2.1 FCFS (*First Come First Served*)

Numele algoritmului înseamnă în traducere directă înseamnă *primul venit - primul servit*. Este cunoscut și sub numele de *FIFO (First In First Out)*. Procesele sunt servite în ordinea lor cronologică. Este un algoritm simplu, dar nu prea eficient.

Spre exemplu, presupunem că există 3 procese cu următorii timpi de execuție:

- procesul 1 durează 24 minute;
- procesul 2 durează 3 minute;
- procesul 3 durează 3 minute.

Dacă acestea sosesc la sistem în ordinea 1, 2, 3, atunci timpul mediu de servire este:

$$\frac{24 + 27 + 30}{3} = 27 \text{ minute}$$

Dacă sosesc în ordinea 2, 3, 1 timpul mediu de servire este:

$$\frac{3 + 6 + 30}{3} = 13 \text{ minute}$$

### 9.3.2.2 SJF (*Shortest Job First*).

Se execută primul, jobul (procesul) care consumă cel mai puțin timp procesor. Probabil că acest algoritm este cel mai bun, însă are un mare dezavantaj: presupune cunoașterea exactă a timpului procesor necesar execuției procesului.

### 9.3.2.3 Algoritm bazat pe priorități.

Este algoritmul cel mai des folosit. Toți ceilalți algoritmi descriși în literatură sunt cazuri particulare ale acestuia. Fiecare proces primește un număr de prioritate. Procesele se ordonează după aceste priorități, apoi se execută în această ordine.

Se disting două metode de stabilire a priorităților:

1. procesul primește prioritatea la intrarea în sistem și și-o păstrează până la sfârșit;
2. Sistemul de operare calculează prioritățile după reguli proprii și le atașează, dinamic, proceselor în execuție.

Varianta 1) este folosită destul de des dar are dezavantajul că dacă apar multe procese cu prioritate mare, atunci cele cu prioritate mică așteaptă practic indefinit. Acest fenomen este cunoscut sub numele de *starvation*. Pentru a evita acest fenomen, de obicei cele două metode sunt combinate. Astfel, unui proces cu prioritate mică i se va mări prioritatea cu cât așteaptă mai mult și se readuce la valoare inițială după ce intră în execuție.

#### 9.3.2.4 Algoritm bazat pe termene de terminare (deadline scheduling)

Acest algoritm de planificare este folosit în principal la sistemele cu cerințe de timp real foarte stricte. În asemenea sisteme este critic pentru un task sau proces să se termine la timp. Astfel fiecărui task  $i$  se atașează un termen de terminare. Un caz special sunt task-urile repetitive cărora li se asociază termene de terminare asociate fiecărei iterații. Planificatorul folosește aceste termene pentru a decide când și cărui task și pentru cât timp să-i acorde procesorul în așa fel încât să se termine la timp. Evident pentru acest gen de planificare este esențială cunoașterea exactă a duratei fiecărui task. La apariția unui nou task în sistem, planificatorul îl va accepta pentru execuție doar dacă poate găsi o alocare prin care termenele noului task să fie respectate fără a le afecta pe cele ale task-urilor deja existente.

O strategie posibilă pentru calcularea secvenței de execuție a task-urilor este aducerea în stare **RUN** întotdeauna a task-ului cu termenul de terminare cel mai apropiat.

#### 9.3.2.5 Round-Robin (planificare circulară).

Acești algoritmi sunt destinați în special sistemelor de operare care lucrează în timesharing. Pentru realizare, se definește o *cuantă de timp*, de obicei între 10-100 milisecunde. Coada **READY** a proceselor este tratată circular. Pe durata unei cuante se alocă procesorul unui proces. După epuizarea acestei cuante, procesul este trecut la sfârșitul cozii, al doilea proces este preluat de către procesor ș.a.m.d.

Există și *variante de Round-Robin*. Spre exemplu, dacă un proces nu și-a consumat în întregime cuanta (spre exemplu din cauza unei operații de I/O), atunci locul lui în coada **READY** este invers proporțional cu partea consumată din cuantă. De exemplu, dacă și-a consumat toată cuanta, atunci procesul trece la sfârșitul cozii; dacă și-a consumat numai jumătate din cuantă, atunci el trece la mijlocul cozii etc.

O altă variantă urmărește echilibrarea sistemului folosind așa-zisul *Round-Robin cu reacție*. Când un proces nou este acceptat, el se rulează mai întâi atâtea cuante câte au consumat celelalte procese, după care trece la planificarea Round-Robin simplă.

#### 9.3.2.6 Algoritm de cozi pe mai multe nivele

Acest algoritm se aplică atunci când lucrările pot fi clasificate ușor în grupe distincte. Spre exemplu, la un sistem de operare de tip mixt, interactiv și serial, pot fi create 5 cozi distincte, înșirate mai jos, de la cea mai prioritară, la cea mai puțin prioritară:

Taskuri sistem.  
Lucrările interactive.  
Lucrări în care se editează texte.  
Lucrări seriale obișnuite.  
Lucrări seriale ale studenților.  
...



## 10 Gestiunea memoriei

### 10.1 Structură; calculul de adresă

#### 10.1.1 Problematika gestiunii memoriei

Pentru a fi executat, un program are nevoie de o anumită cantitate de memorie. Dacă se lucrează în multiprogramare este necesar ca în memorie să fie prezente simultan mai multe programe. Fiecare program folosește zona (zonele) de memorie alocată (alocate) lui, independent de eventuale alte programe active.

Se știe că, în general, pe durata execuției unui program, necesarul de memorie variază. Acest necesar variabil de memorie apare din două surse: folosirea *variabilelor dinamice* și *segmentarea*. Variabilele dinamice sunt alocate de către proces, de regulă dintr-o zonă de memorie numită *heap* și spațiul este eliberat tot de către proces. Alocarea și eliberarea se face prin perechi de funcții `new - dispose`, `malloc - mfree`, `constructor - destructor` etc.

*Segmentarea* este o tehnică prin care un program executabil este decupat în entități distincte numite *segmente*. Segmentele pot fi: de cod, date sau de stivă. Pe durata vieții programului unele dintre segmente pot fi prezente în memorie, altele nu. Programul însuși poate cere încărcarea sau reîncărcarea unui segment, fie într-o zonă de memorie liberă, fie în locul altui (altor) segment(e). La construirea programului executabil se poate defini *structura de acoperire a segmentelor*. Spre exemplu, să presupunem că un program este segmentat și are forma din fig. 10.1a. Atunci segmentele active la un moment dat pot fi: A, AC, AB, AF, ACE sau ACD, așa cum reiese din fig.10.1b.

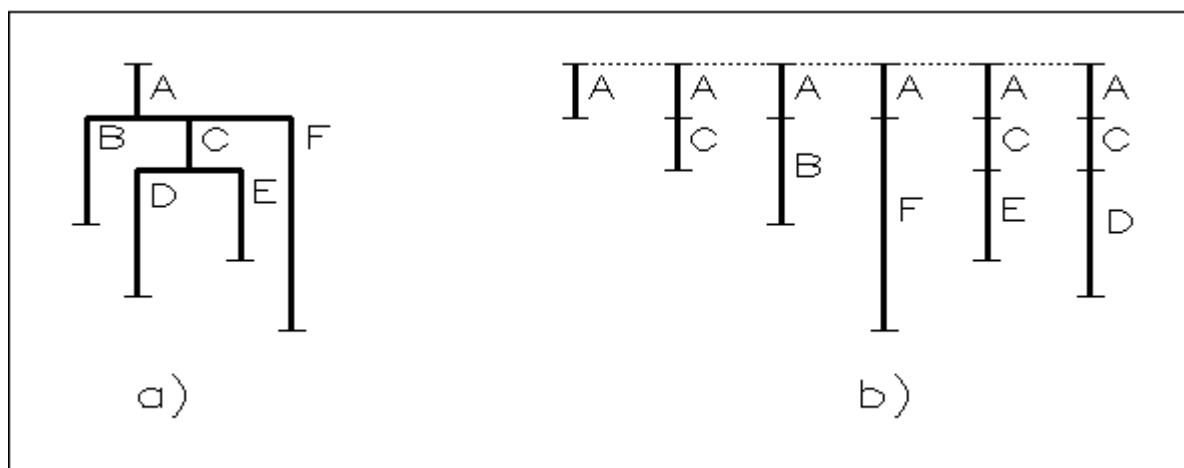


Figura 10.1 Instance posibile ale unui program segmentat

Evident, la fiecare din situațiile precedente se solicită o altă cantitate de memorie.

*Spațiul* de memorie principală al unui **SC**, adică ceea ce poate fi accesat în mod direct de către **CPU**, este în prezent încă *limitat*. Chiar dacă spațiul de memorie oferit de actualele **SC** este mult mai mare decât cel oferit acum 20 de ani sau chiar 10 ani (iar prețul de cost al memoriei a scăzut drastic), limitarea totuși rămâne. Din această cauză, **SO** cu sprijinul **SC** trebuie să gestioneze eficient folosirea acestui spațiu.

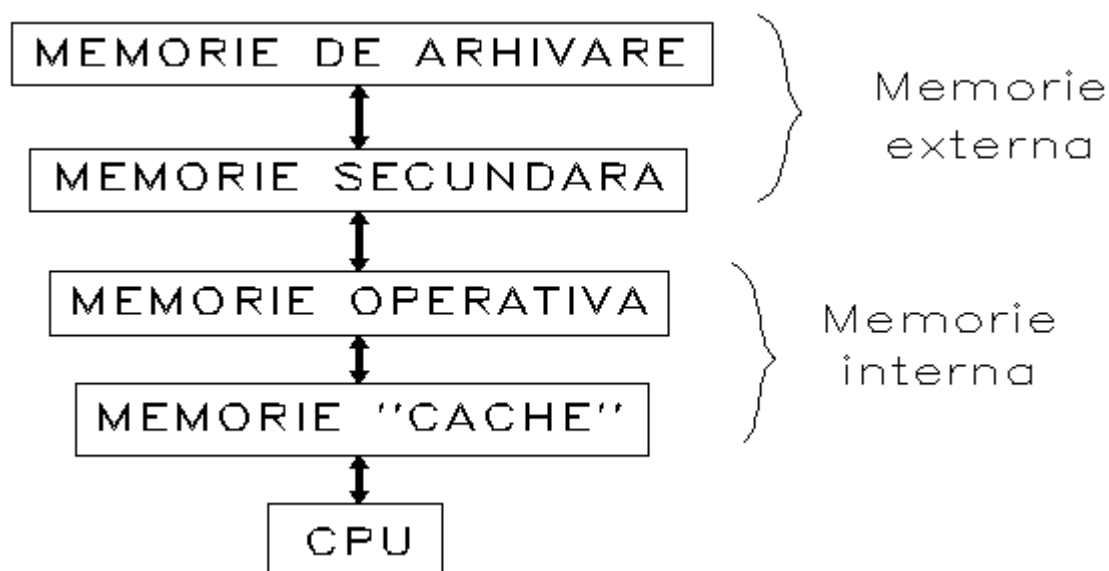
Problemele privind gestiunea memoriei sunt rezolvate la nivelul inferior de către **SC**, extins eventual cu o componentă de *management al memoriei*. La nivelul superior, rezolvarea se face de către **SO** în colaborare cu **SC**. Principalele obiective ale gestiunii memoriei sunt:

- Calculul de translație a adresei (relocare) ;
- Protecția memoriei;
- Organizarea și alocarea memoriei operative;
- Gestiunea memoriei secundare;
- Politici de schimb între proces, memoria operativă și memoria secundară.

Vom aborda fiecare dintre aceste obiective, cu excepția protecției memoriei. În general protecția memoriei este puternic dependentă de hardware, motiv pentru care nu o tratăm aici, într-un cadru general. Fiecare pereche sistem de calcul – sistem de operare are propria politică de protecție a memoriei, deci aceasta trebuie studiată în context concret.

### 10.1.2 Structura ierarhică de organizare a memoriei

În structura sa actuală, memoria unui sistem de calcul apare ca în fig. 10.2.



**Figura 10.2 Structura memoriei unui sistem de calcul**

Să prezentăm componentele de jos în sus.

*Memoria "cache"* conține informațiile cele mai recent utilizate de către **CPU**. Ea are o capacitate relativ mică, dar cu timp de acces foarte rapid. La fiecare acces, **CPU** verifică dacă data invocată se află în memoria "cache" și abia apoi solicită memoria operativă. Dacă este, atunci are loc schimbul între **CPU** și ea. Dacă nu, atunci data este căutată în nivelele superioare. Aplicând *principiul vecinătății* [40], data respectivă este adusă din nivelul la care se află, dar odată cu ea se aduce un număr de locații vecine ei, astfel încât, împreună să umple memoria "cache". Principiul de lucru este cel al memoriilor tampon temporare, prezentat în 11.1. Într-o secțiune viitoare vom prezenta în detaliu cum funcționează memoria cache.

În ce constă principiul vecinătății? P. Denning afirmă, pe baza unor studii de simulare temeinice, că dacă la un moment dat se solicită o dată dintr-un anumit loc atunci solicitarea din momentul următor se va face, cu mare probabilitate, la o dată din apropierea precedentei.

Memoria "cache" este practic prezentă la toate sistemele de calcul moderne. De exemplu, un notebook actual, Intel Pentium 2GHz poate avea o memorie cache de până la 1Mo.

*Memoria operativă* conține programele și datele pentru toate procesele existente în sistem. În momentul în care un proces este terminat și distrus, spațiul de memorie operativă pe care l-a ocupat este eliberat și va fi ocupat de alte procese. Capacitatea memoriei operative variază azi de la 128Mo până la 8 Go sau chiar mai mult. Viteza de acces este foarte mare, dar mai mică decât a memoriei "cache".

*Memoria secundară* apare la **SO** care dețin mecanisme de memorie virtuală. De asemenea, tot în cadrul memoriei secundare poate fi inclus spațiul disc de swap (vezi 8.2.1.1). Această memorie este privită ca o extensie a memoriei operative. Suportul ei principal este discul magnetic. Accesul la această memorie este mult mai lent decât la cea operativă.

*Memoria de arhivare* este gestionată de utilizator și constă din fișiere, baze de date ș.a. rezidente pe diferite suporturi magnetice (discuri, benzi, etc.).

Memoria "cache" și memoria operativă formează ceea ce cunoaștem sub numele de *memoria internă*. Accesul **CPU** la acestea se face în mod direct. Pentru ca **CPU** să aibă acces la datele din memoria secundară și de arhivare, acestea trebuie mai întâi mutate în memoria internă.

Din punct de vedere a performanțelor, privind fig. 10.2 de jos în sus se disting următoarele caracteristici ale componentelor memoriei:

- viteza de acces la memorie scade;
- prețul de cost pe unitatea de alocare scade;
- capacitatea de memorare crește.

### 10.1.3 Mecanisme de traducere a adresei

Adresarea memoriei constă în realizarea legăturii între un obiect al programului și adresa corespunzătoare din memoria operativă a **SC**. Pentru a explica mecanismele de traducere, să adoptăm câteva *notații*:

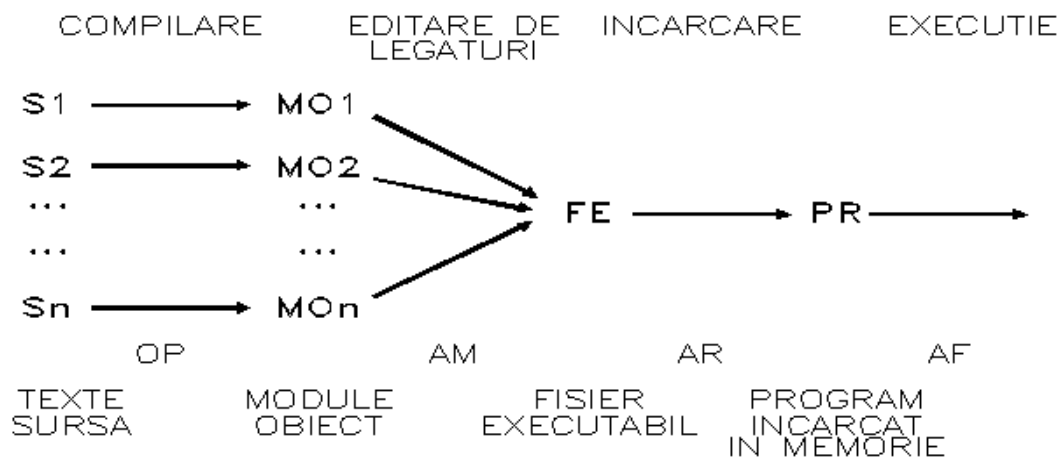
- **OP** notăm spațiul de nume al obiectelor din programul sursă: nume de constante, de variabile, de etichete, de proceduri etc.
- **AM** este adresa relativă din cadrul unui modul compilat.
- **AR** este adresa relocabilă – adresă relativă din cadrul unui segment dintr-un fișier executabil.
- **AF** notăm mulțimea adreselor fizice din memoria operativă la care face referire programul în timpul execuției.

*Calculul de adresă* este modalitatea prin care se ajunge de la un obiect sursă din **OP** la adresa lui fizică din **AF**. După cum se vede, acest calcul necesită trei faze, corespunzătoare fazelor în care se poate afla un program (vezi 8.2.1.2). Matematic vorbind, calculul de adresă se realizează prin compunerea a trei funcții: **c**, **l**, **t**, astfel:

$$OP \xrightarrow{c} AM \xrightarrow{l} AR \xrightarrow{t} AF$$

În fig. 10.3 sunt ilustrate fazele prin care trece un program componentele invocate și etapele calculului de adresă, de la textul sursă până la execuție.





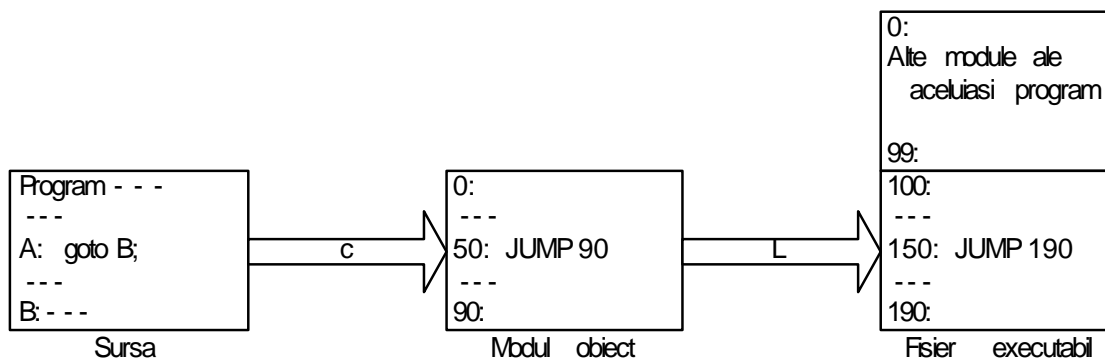
**Figura 10.3 Fazele traducerii unui program**

### 10.1.3.1 Faza de compilare

Faza de *compilare* transformă un text sursă  $S_i$  într-un *modul obiect*  $MO_i$ . Corespunzător, numele obiectelor program sunt transformate în numere reprezentând  $AM$ , adică adrese în cadrul modulului obiect. În cadrul fiecărui modul obiect aceste adrese încep de la 0. Deci prima funcție din calculul de adresă:

$$c: OP \rightarrow AM$$

este executată de către compilator sau asamblor. Modul ei de evaluare depinde de limbajul, de compilatorul și de  $SO$  concret. Teoria compilării [33] are în vedere această evaluare, ea nu intră în scopurile noastre actuale. În fig. 10.4, prima parte, este prezentat un exemplu de aplicare a acestei funcții.



**Figura 10.4 Traducere de la sursă la fișier executabil**

### 10.1.3.2 Faza editării de legături

Faza de *editare de legături* grupează mai multe module formând un *fișier executabil*. Editorului de legături îi revine sarcina evaluării celei de-a doua funcții de calcul a adresei. Această funcție transformă adresele din cadrul modulelor în așa zisele *adrese relocabile* (*relocatable*). Funcția de legare este:

$$l: AM \rightarrow AR$$

Particularitățile de evaluare a acestei funcții sunt proprii editorului de legături. În fig. 10.4 sugerăm modul în care acționează compunerea celor două funcții.

### 10.1.3.3 Faza de încărcare și execuție

Funcția

$t: AR \rightarrow AF$

este *funcția de translatare (relocare) a adresei*. În mod obișnuit ea este executată de către **CPU**. Translatarea depinde de tipul sistemului de calcul, în particular de existența dispozitivului de management al memoriei. Pentru a vedea principiul ei de lucru, presupunem că fișierul executabil conține înregistrări succesive cu instrucțiuni de forma celei din fig. 10.5.

ARI CO A1 A2 - - - An

**Figura 10.5** Formatul unei instrucțiuni mașină

Semnificațiile câmpurilor instrucțiunii sunt:

- *ARI* este adresa relocabilă a instrucțiunii;
- *CO* este codul operației,
- *A1, ..., An* sunt argumentele instrucțiunii mașină: nume de regiștri, constante (argumente *immediate*), adrese relocabile din memorie

Presupunem că fiecare instrucțiune încapă într-o locație de memorie și că instrucțiunile sunt plasate una după cealaltă în locații succesive. Atunci încărcarea unui astfel de fișier se poate face folosind un încărcător analog lui **LOADERABSOLUT**, descris în 8.3 fig. 8.5.

Să presupunem că mașina dispune de un singur registru general pe care-l vom numi *A* (de la *registru acumulator*). De asemenea, presupunem că între *A1, ..., An* există o singură adresă relocabilă. Evident că în realitate structura unui fișier executabil este mai complexă, dar pentru moment aceasta ne este suficientă.

Să adoptăm următoarele notații:

- $M[0..m]$  conține locațiile memoriei operative;
- *pc* (Program Counter) indică adresa fizică a instrucțiunii care urmează a fi executată;
- *w* este conținutul instrucțiunii curente;
- *Opcode(w)* este o funcție care furnizează codul operației din instrucțiunea curentă. Pentru fixarea ideilor, să presupunem că:
  - 1 este codul adunării;
  - 2 este codul operației de memorare (depunere din *A*) într-o anumită locație;
  - 3 este codul instrucțiunii de salt necondiționat;
  - ș.a.m.d.
- *Adress(w)* este o funcție care furnizează valoarea adresei relocabile aflată între argumentele instrucțiunii curente.

Cu aceste notații, în fig. 10.6 este descris modul de funcționare a **CPU** și de acțiune a funcției de translatare *t*. De fapt, în fig. 10.6 este schițat un *interpretor* al limbajului care are instrucțiuni de forma celei din fig. 10.5.

```

pc = t(adresa-de-start-a-programului);
do {
    w = M[pc];                // operația fetch
    co = Opcode(w);
    adr = Adress(w);
    pc = pc+1;
    switch (co) {
        1: A:=A+M[t(adr)];    // adunare
        2: M[t(adr)]:=A;      // memorare
        3: pc:=t(adr);        // salt necondiționat
        - - -
    }
} while (false);

```

**Figura 10.6 Funcționarea CPU și calculul funcției de traducere**

## 10.2 Scheme simple de alocare a memoriei

### 10.2.1 Clasificarea tehnicilor de alocare

Problema alocării memoriei se pune în special la sistemele multiutilizator, motiv pentru care în continuare ne vom ocupa aproape exclusiv numai de aceste tipuri de sisteme. Tehnicile de alocare utilizate la diferite **SO** se împart în două mari categorii, fiecare categorie împărțindu-se la rândul ei în alte subcategorii, ca mai jos [19] [10].

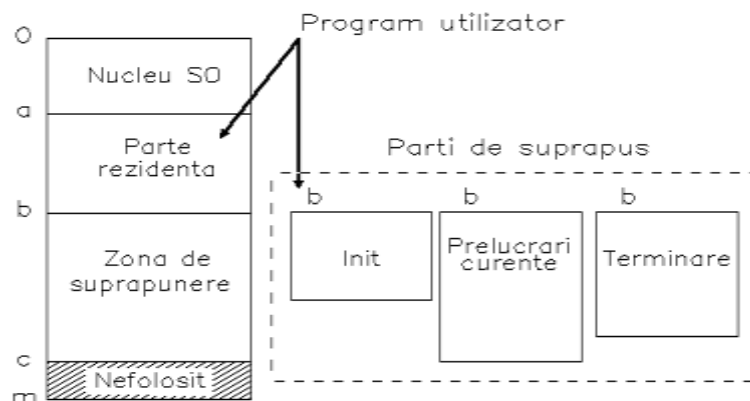
- alocare reală:
  - la **SO** monoutilizator;
  - la **SO** multiutilizator:
    - cu partiții fixe (statică):
      - absolută;
      - relocabilă;
    - cu partiții variabile (dinamică);
- alocare virtuală:
  - paginată;
  - segmentată;
  - segmentată și paginată.

### 10.2.2 Alocarea la sistemele monoutilizator

La sistemele monoutilizator este disponibil aproape întreg spațiul de memorie. Gestiunea acestui spațiu cade exclusiv în sarcina utilizatorului. El are la dispoziție tehnici de *suprapunere (overlay)* pentru a-și putea rula programele mari. În fig. 10.6 este ilustrat acest mod de lucru.

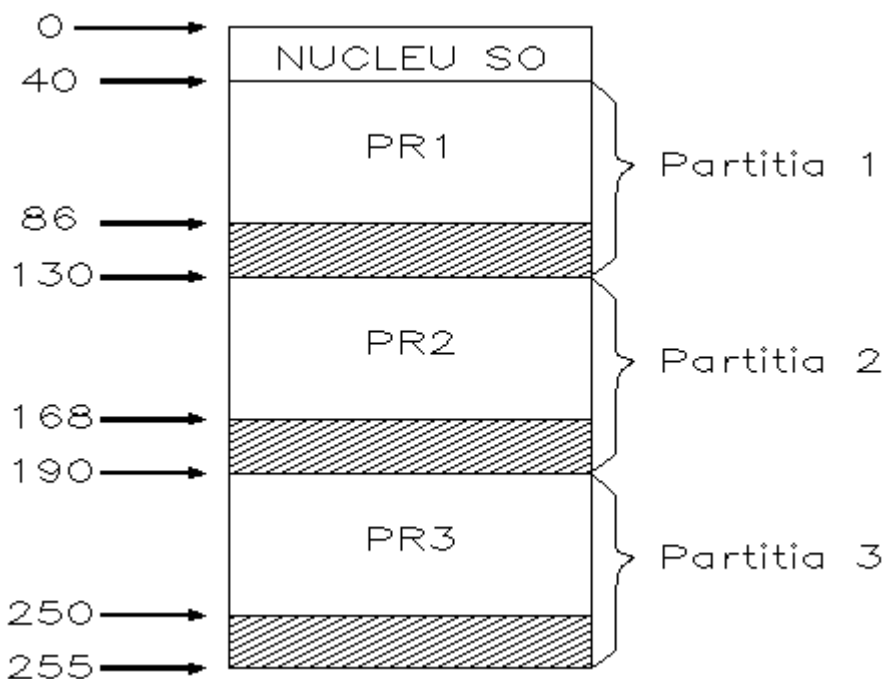
Porțiunea dintre adresele 0 și  $a-1$  este rezervată nucleului **SO**, care rămâne acolo de la încărcare și până la oprirea sistemului. Între adresele  $c$  și  $m-1$  (dacă memoria are capacitatea de  $m$  locații) este spațiu nefolosit de către programul utilizator activ. Evident, adresa  $c$  variază de la un program utilizator la altul.

### 10.2.3 Alocarea cu partiții fixe



**Figura 10.7 Alocarea memoriei la sistemele monoutilizator**

Acest mod de alocare mai poartă numele de alocare statică sau alocare MFT - Memory Fix Tasks. El presupune decuparea memoriei în zone de lungime fixă numite *partiții*. O partiție este alocată unui proces pe toată durata execuției lui, indiferent dacă o ocupă complet sau nu. Un exemplu al acestui mod de alocare este descris în fig. 10.8. Zonele hașurate fac parte din partiții, dar procesele active din ele nu le utilizează.



**Figura 10.8 Exemplu de alocare cu partiții fixe**

**Alocarea absolută** se face pentru programe pregătite de editorul de legături pentru a fi rulate într-o zonă de memorie prestabilă și numai acolo.

Mult mai folosită este **alocarea relocabilă**, la care adresarea în partiție se face cu bază și deplasament. La încărcarea unui program în memorie, în registrul lui de bază se pune adresa de început a partiției. În cazul sistemelor seriale cu multiprograme, dacă un proces este plasat spre execuție într-o partiție insuficientă, el este eliminat din sistem fără a fi executat.

De obicei, partițiile au lungimi diferite. Una dintre problemele cele mai dificile este fixarea acestor dimensiuni. Dificultatea constă în faptul că nu se pot prevedea în viitor cantitățile de memorie pe care le vor solicita procesele încărcate în aceste partiții. Alegerea unor dimensiuni

mai mari scade probabilitatea ca unele procese să nu poată fi executate, dar scade și numărul proceselor active din sistem.

Acest mod de alocare este utilizat preponderent de către sistemele seriale. La fiecare partiție există un șir de procese care așteaptă să fie executate. Modul în care se organizează acest sistem de așteptare poate influența performanțele de ansamblu ale sistemului și poate eventual atenua efectul unei dimensionări defectuoase a partițiilor. În general există două moduri de legare a proceselor la partiții:

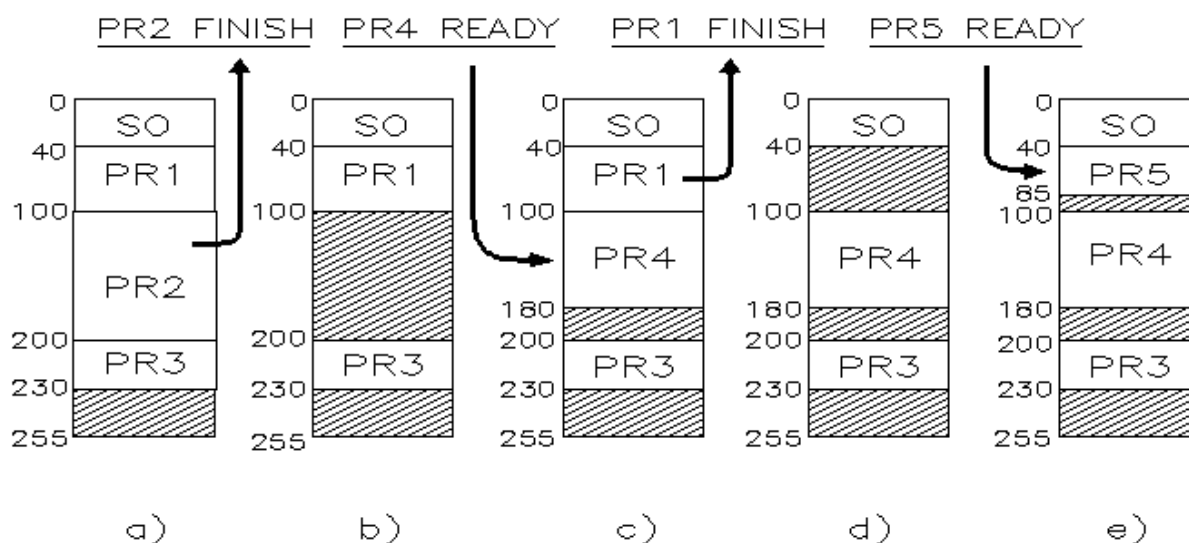
- *Fiecare partiție are coada proprie*; operatorul stabilește de la început care sunt procesele care vor fi executate în fiecare partiție.
- *O singură coadă pentru toate partițiile*; **SO** alege, pentru procesul care urmează să intre în lucru, în ce partiție se va executa.

Legarea prin cozi proprii partițiilor este mai simplă din punctul de vedere al **SO**. Primele sisteme multiutilizator au adoptat acest mod de legare. În schimb, legarea cu o singură coadă este mai avantajoasă, pentru faptul că se poate alege partiția cea mai potrivită pentru plasarea unui proces.

#### 10.2.4 Alocarea cu partiții variabile

Acest mod de legare mai este cunoscut și sub numele de alocare dinamică sau alocare MVT - Memory Variable Task. El reprezintă o extensie a alocării cu partiții fixe, care permite o exploatare mai suplă și mai economică a memoriei **SC**. *În funcție de solicitările la sistem și de capacitatea de memorie încă disponibilă la un moment dat, numărul și dimensiunea partițiilor se modifică automat.*

În fig. 10.9 sunt prezentate mai multe stări succesive ale memoriei, în acest mod de alocare.

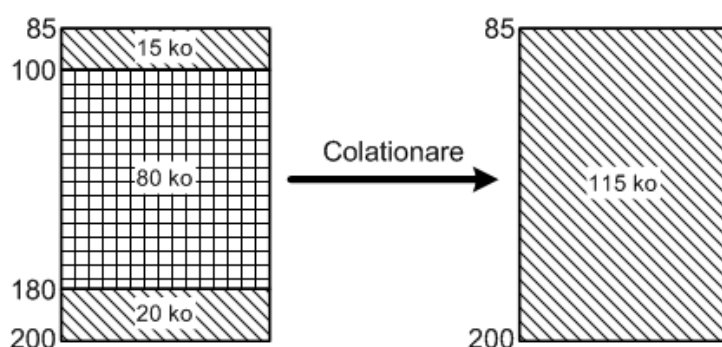


**Figura 10.9 Evoluția proceselor la alocarea cu partiții variabile**

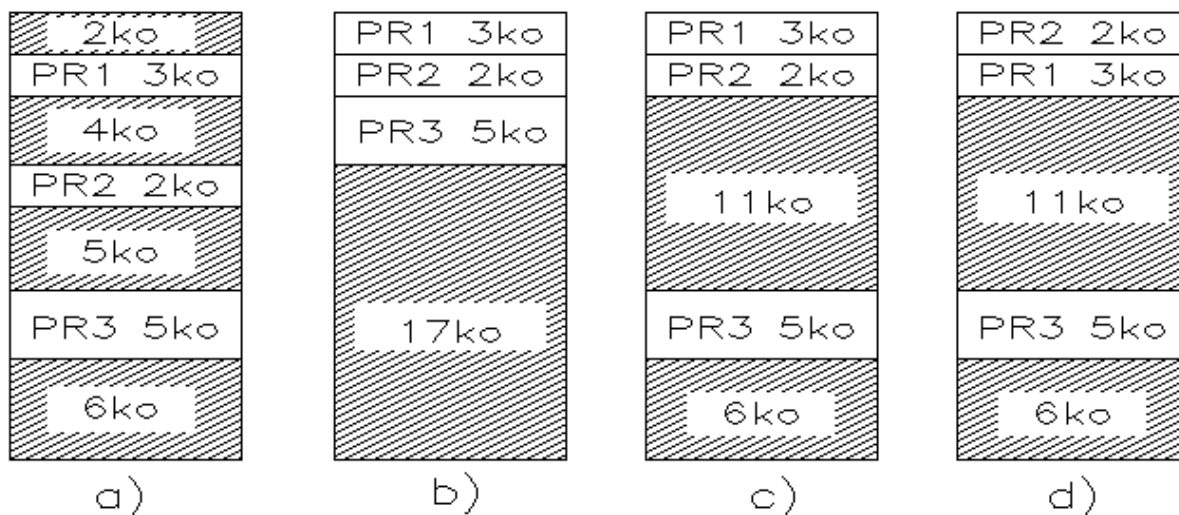
În momentul în care procesul intră în sistem, el este plasat în memorie într-un spațiu în care încapă cea mai lungă ramură a sa. Spațiul liber în care a intrat procesul este acum descompus în două partiții: una în care se află procesul, iar cealaltă într-un spațiu liber mai mic. Este ușor de observat că dacă sistemul funcționează timp îndelungat, atunci numărul spațiilor libere va crește, iar dimensiunile lor vor scădea. Fenomenul este cunoscut sub numele de *fragmentarea*

*internă a memoriei.* După cum se va vedea, acest fenomen poate avea efecte neplăcute. În momentul în care un proces nu are spațiu în care să se încarce, **SO** poate lua una din următoarele trei decizii:

- Procesul așteaptă* până când i se eliberează o cantitate suficientă de memorie.
- SO** încearcă *alipirea unor spații libere vecine - colaționare*, în speranța că se va obține un spațiu de memorie suficient de mare. Spre exemplu, dacă momentul care urmează după cel din fig. 10.9d este terminarea procesului *PR4*, atunci în gestiunea sistemului apar trei zone libere adiacente: prima de 15Ko, a doua de 80Ko, iar a treia de 20Ko (vezi fig. 10.10). **SO** poate (nu întotdeauna o și face în mod automat) să formeze din aceste trei spații unul singur de 115Ko, așa cum se vede în fig. 10.10.
- SO** decide efectuarea unei operații de *compactare a memoriei (relocare)* adică de deplasare a partițiilor active către partiția monitor pentru a se absorbi toate "fragmentele" de memorie neutilizate. Este posibil ca spațiul astfel obținut să fie suficient pentru încărcarea procesului. În fig. 10.11b este dat un exemplu de compactare.



**Figura 10.10 Colaționarea de spații libere vecine**



**Figura 10.11 Posibilități de compactare prin relocare totală sau parțială**

De regulă, compactarea este o operație costisitoare și în practică se aleg soluții de compromis, cum ar fi:

- Se lansează periodic compactarea (de exemplu la 10 secunde), indiferent de starea sistemului. În intervalul dintre compactări memoria apare ca un mozaic de spații ocupate care alternează cu spații libere. Procesele care nu au loc în memorie așteaptă compactarea sau terminarea altui proces.

- Se realizează o compactare parțială pentru a asigura loc numai procesului care așteaptă. Spre exemplu, dacă harta memoriei este cea din fig. 10.11a și un proces cere 10Ko, se poate realiza numai compactarea parțială din fig. 10.11c.
- Se încearcă numai mutarea unora dintre procese cu colaționarea spațiilor rămase libere. Dacă reluăm exemplul de mai sus, este posibilă mutarea procesului *PR2* în primul spațiu liber și problema este rezolvată, harta fiind cea din fig. 10.11d.

Între alocările de tip MFT și MVT nu există practic diferențe hard. Alocarea MVT este realizată de cele mai multe ori prin intermediul unor rutine specializate, eventual microprogramate, deci de către **SO**.

Alocarea MVT a fost utilizată mai întâi la **SC IBM-360** sub **SO OS-360 MVT**, apoi la **PDP 11/45**.

### 10.3 Mecanisme de memorie virtuală

Termenul de *memorie virtuală* este de regulă asociat cu capacitatea de a adresa un spațiu de memorie mai mare decât este cel disponibil la memoria operativă a **SC** concret. Conceptul este destul de vechi, el apărând odată cu **SO ATLAS** al Universității Manchester, Anglia, 1960 [19]. Se cunosc două metode de virtualizare, mult înrudite după cum vom vedea în continuare. Este vorba de alocarea paginată și alocarea segmentată. Practic, toate sistemele de calcul actuale folosesc, într-o formă sau alta, mecanisme de memorie virtuală.

#### 10.3.1 Alocarea paginată

Alocarea paginată a apărut la diverse **SC** pentru a evita fragmentarea excesivă, care apare la alocarea MVT, și drept consecință, la evitarea aplicării relocării. Această alocare presupune cinci lucruri și anume:

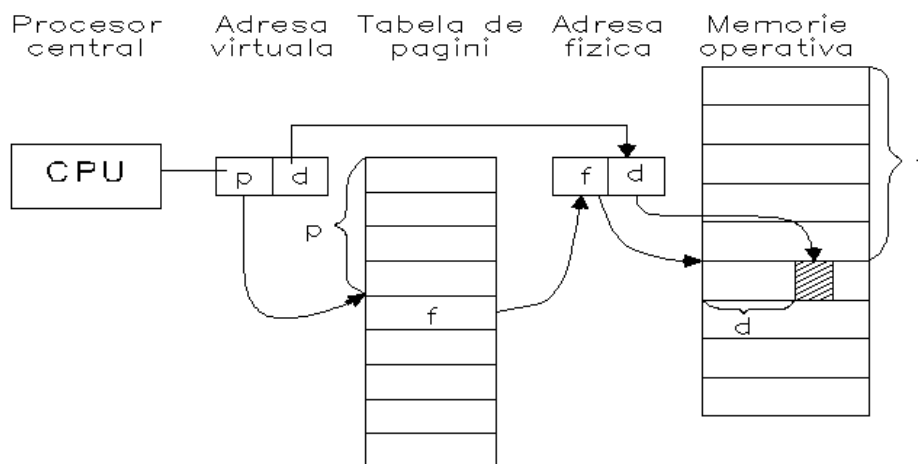
- a) Instrucțiunile și datele fiecărui program sunt împărțite în zone de lungime fixă, numite *pagini virtuale*. Fiecare **AR** (notațiile din 10.1.3) aparține unei pagini virtuale. Paginile virtuale se păstrează în memoria secundară.
- b) Memoria operativă este împărțită în zone de lungime fixă, numite *pagini fizice*. Lungimea unei pagini fizice este fixată prin hard. Paginile virtuale și cele reale au aceeași lungime, lungime care este o putere a lui 2, și care este o constantă a sistemului (de exemplu 1Ko, 2Ko etc).
- c) Fiecare **AR** este o pereche de forma:  $(p, d)$ , unde  $p$  este numărul paginii virtuale, iar  $d$  adresa în cadrul paginii.
- d) Fiecare **AF** este de forma  $(f, d)$ , unde  $f$  este numărul paginii fizice, iar  $d$  adresa în cadrul paginii.
- e) Calculul funcției de translație  $t : AR \rightarrow AF$  se face prin hard, conform schemei din fig. 10.12.

Dacă prin  $M[0..m]$  notăm memoria operativă, prin  $k$  puterea lui 2 (numărul de biți) care dă lungimea unei pagini, prin **TP** adresa de start a tabelului de pagini, atunci algoritmul calculului funcției  $t$  este:

$$t(p, d) = M[TP + p] * 2^k + d$$

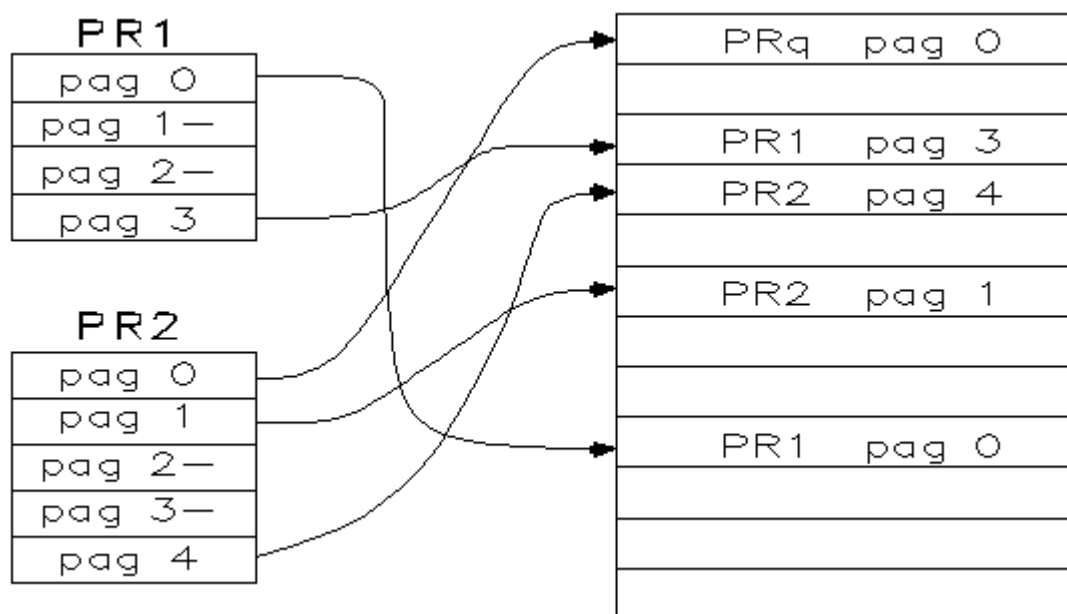
Acest calcul este valabil atunci când tabela de pagini ocupă un spațiu în memoria operativă. Există însă **SC** ce dispun, hard, de o memorie specială de capacitate mică, numită *memorie asociativă*. Calitatea ei fundamentală este *adresarea prin conținut*, ceea ce înseamnă că

găsește locația care are un conținut specificat, căutând simultan în toate locațiile ei. Memoria asociativă conține atâtea locații câte pagini fizice are. În fiecare locație a memoriei asociative este trecut numărul paginii virtuale care se află în pagina fizică având numărul de ordine identic cu numărul de ordine al locației de memorie asociativă. Atunci când se dă un număr de pagină virtuală, se obține automat numărul paginii fizice care o găzduiește.



**Figura 10.12** Translatarea unei pagini virtuale într-una fizică

Fiecare proces are propria lui tabelă de pagini, în care este trecută adresa fizică a paginii virtuale, dacă ea este prezentă în memoria operativă. La încărcarea unei noi pagini virtuale, aceasta se depune într-o pagină fizică liberă. Deci, în memoria operativă, paginile fizice sunt distribuite în general necontiguu, între mai multe procese. Spunem că are loc o *proiectare a spațiului virtual peste cel real*. Este posibil, la un moment dat, să existe situația din fig. 10.13.



**Figura 10.13** Două procese, într-o alocare paginată

Acest mecanism are avantajul că folosește mai eficient memoria operativă, fiecare program ocupând numai memoria strict necesară la un moment dat. Un avantaj colateral, dar nu de neglijat, este folosirea în comun a unei porțiuni de cod.

Să presupunem că avem un editor de texte al cărui cod (instrucțiuni pure, fără date) ocupă două pagini. Mai presupunem că fiecare utilizator consumă câte o pagină pentru datele proprii



de editat. Dacă sunt trei utilizatori, ei trebuie în mod normal să "consume" nouă pagini din memoria operativă. În fig. 10.14 se arată cum pot fi consumate numai cinci pagini.

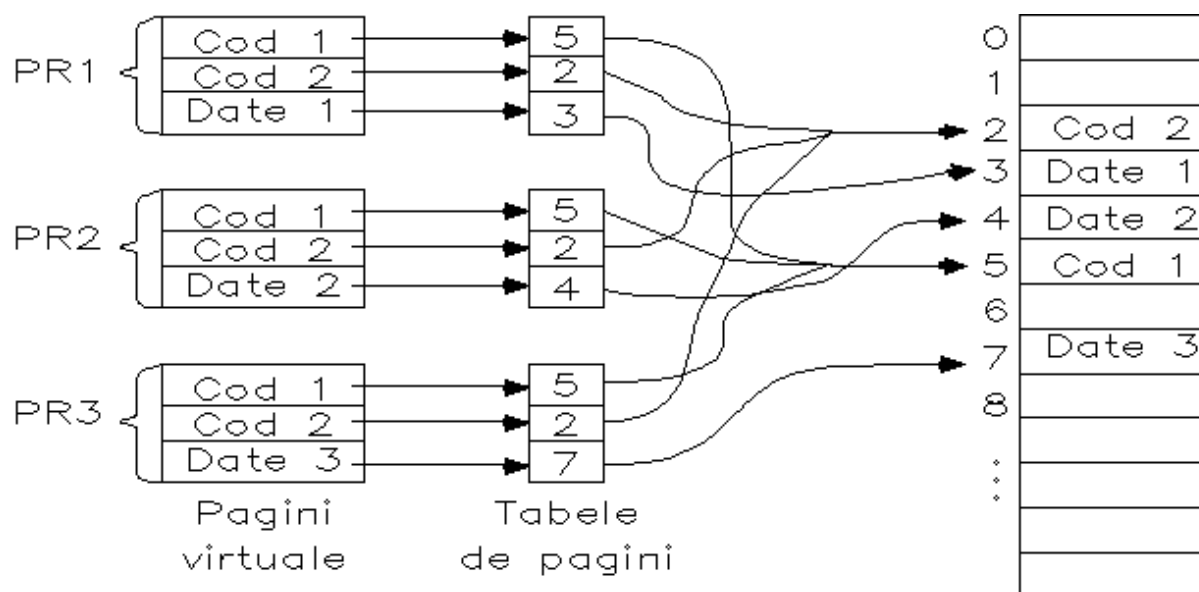


Figura 10.14 Folosirea în comun a unui cod

### 10.3.2 Alocare segmentată

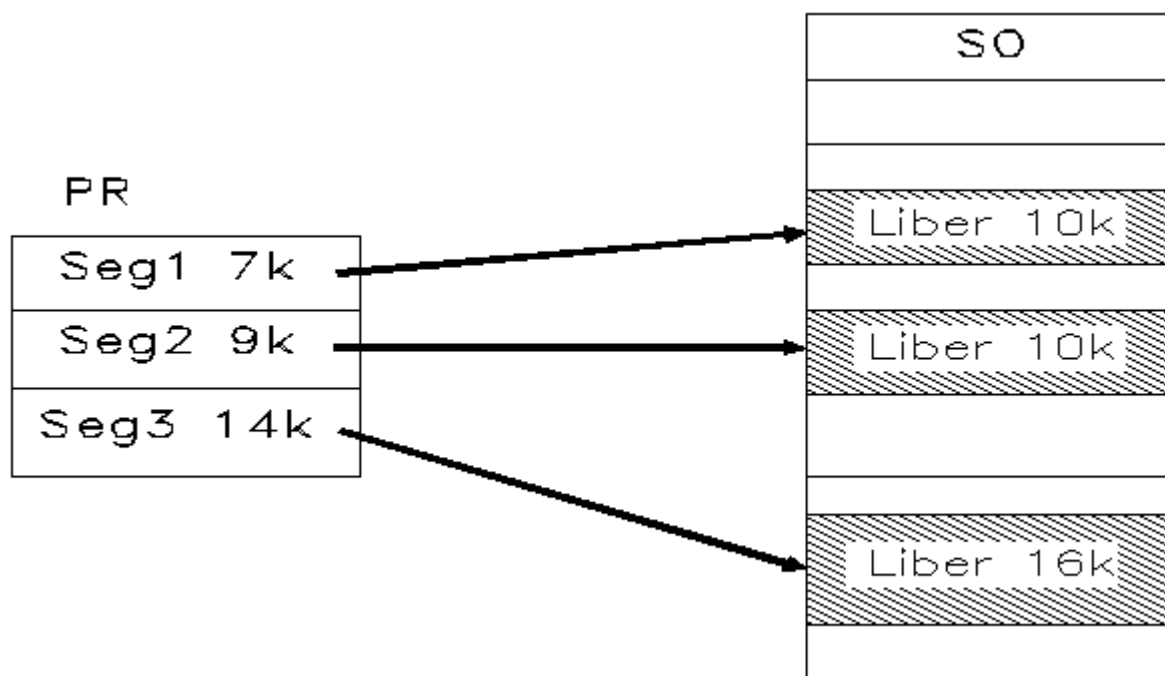
Atunci când am vorbit despre alocarea reală, am văzut că fiecare proces trebuia să ocupe un spațiu contiguu de memorie, numit partiție. Ceea ce introduce nou mecanismul de alocare segmentată este faptul că textul unui program poate fi plasat în zone de memorie distincte, fiecare dintre ele conținând o bucată de program numită *segment*. Singura deosebire principală dintre alocarea paginată și cea segmentată este aceea că segmentele *sunt de lungimi diferite*. În fig. 10.15 am ilustrat această situație cu locurile ce vor fi ocupate de un program format din trei segmente.

În mod analog cu alocarea paginată, o adresă virtuală este o pereche  $(s, d)$ , unde  $s$  este numărul segmentului, iar  $d$  este adresa în cadrul segmentului. Adresa reală (fizică) este o adresă obișnuită. Fiecare proces activ are o *tabelă de segmente*. Fiecare intrare în această tabelă conține *adresa de început a segmentului*. Calculul de adresă se face analog celui de la alocarea paginată. Presupunem că la adresa **TS** se află începutul tabelii de segmente. Cu notațiile obișnuite, funcția  $t$  de traducere a adresei se calculează astfel:

$$t(s, d) = M[TS + s] + d$$

Pe lângă avantajul net față de alocările pe partiții, alocarea segmentată mai prezintă încă două avantaje:

- Se pot crea *segmente reentrante*, - cod pur - care pot fi folosite în comun de către mai multe procese. Pentru aceasta este suficient ca toate procesele să aibă în tabelele lor aceeași adresă pentru segmentul pur. A se vedea aceeași problemă discutată la alocarea paginată.
- Se poate realiza o *foarte bună protecție a memoriei*. Fiecare segment în parte poate primi alte drepturi de acces, drepturi trecute în tabela de segmente. La orice calcul de adresă se pot face și astfel de verificări. Pentru detalii se pot consulta lucrările [10], [19].



**Figura 10.15 Alocare necontiguă prin segmentare**

### 10.3.3 Alocare segmentată și paginată

La alocarea segmentată am arătat că adresa fizică este una oarecare. Este deci posibil să apară fenomenul de fragmentare, despre care am vorbit la alocarea cu partiții variabile (10.2.4). Ideea alocării segmentate și paginate este aceea că alocarea spațiului pentru fiecare segment să se facă paginat.

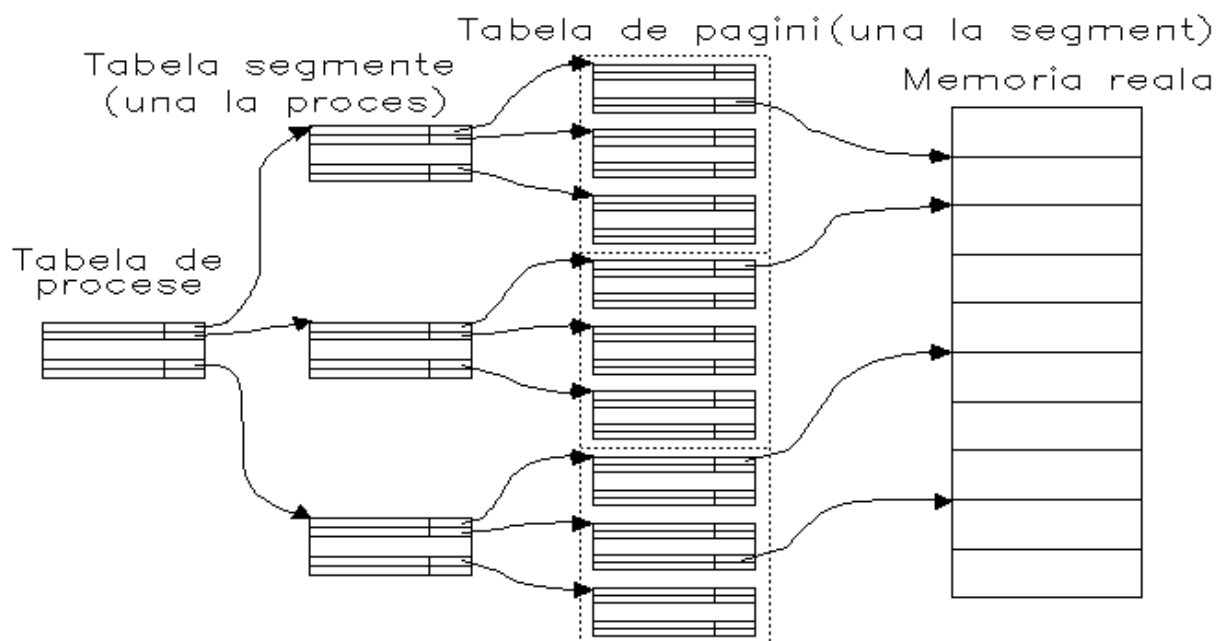
Pentru aceasta, mai întâi *fiecare proces* activ are *propria lui tabelă de segmente*. Apoi, *fiecare segment* dintre cele încărcate în memorie are *propria lui tabelă de pagini*. Fiecare intrare în tabela de segmente are un câmp rezervat adresei de început a tabelii de pagini proprii segmentului respectiv, așa cum se vede fig. 10.16.

O adresă virtuală este de forma:  $(s, p, d)$ , în care  $s$  este numărul segmentului,  $p$  este numărul paginii virtuale în cadrul segmentului, iar  $d$  este deplasamentul în cadrul paginii. O adresă fizică este de forma:  $(f, d)$ , unde  $f$  este numărul paginii fizice, iar  $d$  este deplasamentul în cadrul paginii.

Fie  $k$  constanta ce dă dimensiunea unei pagini ( $2^k$ ),  $TS$  adresa de început a tabelii de segmente a unui proces și presupunem că primul câmp al fiecărei intrări din tabela de segmente este pointerul spre tabela lui de pagini, atunci funcția  $t$  de traducere se calculează astfel:

$$t(s, p, d) = M[M[TS + s] + p] * 2^k + d$$

Printre **SO** remarcabile care utilizează acest mod de alocare trebuie amintit în primul rând MULTICS, cel care a introdus de fapt acest mod de alocare. De asemenea, **SO VAX/VMS** adoptă acest mod de alocare.



**Figura 10.16 Alocare segmentată și paginată**

Calculatoarele IBM-PC cu microprocesor cel puțin 80386, dispun de un mecanism hard de gestiune paginată și segmentată a memoriei extinse. Conceptual el funcționează așa cum am arătat mai sus, completat evident cu o serie de particularități legate de adresarea și protecția memoriei la acest tip de mașină.

## **10.4 Planificarea schimburilor cu memoria**

### **10.4.1 Întrebările gestiunii memoriei și politici de schimb.**

Pe lângă mecanismele de alocare descrise în 10.2 și 10.3, sistemul de operare trebuie să rezolve o serie de probleme care se pot ivi la modul de alocare respectiv. Rezolvarea acestor probleme înseamnă răspunsul la întrebările: CAT? UNDE? CAND? CARE?

*Întrebarea "CAT?"*, apare atunci când se pune problema cantității de memorie alocată. La alocarea pe partiții, se alocă la început toată cantitatea cerută; avem deci o *alocare statică*. La alocarea paginată avem o *alocare dinamică*, fiecare program consumă numai memoria necesară la un moment dat. În 10.3 am prezentat aceste tehnici.

*Întrebarea "UNDE?"*, apare la alocarea cu partiții variabile. Atunci când un program cere intrarea în sistem, trebuie luată decizia: dintre locurile goale pe care le poate ocupa, unde va fi plasat programul? Rezolvările posibile sunt cunoscute în literatură sub numele de *politici de plasare*. Acestea au o utilizare răspândită, depășind cadrul sistemului de operare. Din acest motiv (cât și datorită farmecului lor) le vom dedica o secțiune specială.

*Întrebarea "CAND?"*, apare cel puțin în două situații: la sistemele cu paginare și la alocarea cu partiții variabile. În sistemele cu paginare, trebuie stabilit momentul în care o pagină virtuală este depusă într-o pagină fizică. În literatură, tehnicile de răspuns sunt cunoscute sub numele de *politici de încărcare (fetch)*. Le vom dedica și acestora o secțiune specială.

Aceeași întrebare "*CAND?*", apare pentru a decide momentele de *compactare (relocare)* a memoriei la alocarea cu partiții variabile. În 10.2.4 am arătat în ce constă compactarea și am dat trei posibilități de acțiune în acest sens.

*Intrebarea "CARE?"*, apare la sistemele cu paginare. Să presupunem că la un moment dat toate paginile fizice sunt ocupate. Dacă un program mai cere încărcarea unei pagini, atunci este necesar ca una din paginile fizice să fie evacuată, în memoria secundară, pentru a i se face loc noii pagini. Această manevră poartă numele (după cum am mai spus) de *swapping*. Alegerea paginii care va fi înlocuită face obiectul unor *politici de înlocuire (replacement)*, cărora le vom dedica o secțiune separată.

Deși, face parte dintre metodele implementate hard, deci nu implică deloc sistemul de operare, credem că este cazul să vedem puțin *cum funcționează o memorie cache?* Vom dedica și acesteia o secțiune specială.

## 10.4.2 Politici de plasare.

### 10.4.2.1 Metode de plasare și structuri de date folosite

Nu numai sistemul de operare, ci foarte multe programe de aplicații solicită în timpul execuției lor diverse cantități de memorie. Vom presupune că întreaga cantitate de memorie solicitată la un moment dat *este formată dintr-un șir de octeți consecutivi*. De asemenea, presupunem că există un *depozit de memorie* (numit *heap* în limbajele de programare) de unde se poate obține memorie liberă.

Rezolvarea cererilor presupune existența a două rutine. O primă rutină are sarcina de a *ocupa (aloca)* o zonă de memorie și de a întoarce adresa ei de început. De exemplu, funcția **malloc** din limbajul C și operatorul **new** din limbajul C++ au acest rol. O a doua rutină are rolul de a *elibera* spațiul alocat anterior, în vederea refolosirii lui.

Analogia dintre cele spuse mai sus și alocarea cu partiții variabile este mai mult decât evidentă. Problema politicilor de plasare nu lipsește practic din nici un curs de sisteme de operare. Dintre lucrările mai consistente în acest domeniu amintim [4], [10], [22], [49].

Dată fiind importanța acestor politici, vom da câteva metode de organizare a spațiului de memorie și algoritmii de ocupare și alocare corespunzători. Pe lângă cerințele de funcționalitate, este bine ca plasarea succesivă a programelor (zonelor alocate) să împiedice, pe cât posibil, fragmentarea excesivă a memoriei operative. Altfel, este posibil ca cererile de memorie mai mari să nu poată fi servite deși sistemul dispune per total de memoria necesară.

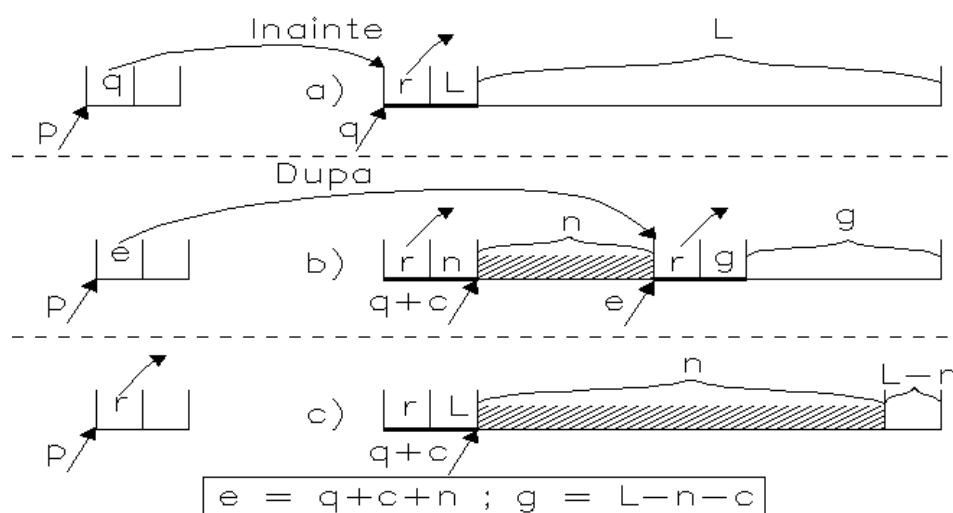
Dintre metodele de plasare, cele mai răspândite sunt următoarele patru:

- Metoda primei potriviri (First-fit).
- Metoda celei mai bune potriviri (Best-fit).
- Metoda celei mai rele potriviri (Worst-fit).
- Metoda alocării prin camarazi (Buddy-system).

Vom analiza pe rând fiecare dintre ele. Primele trei sunt foarte asemănătoare, iar ultima este oarecum deosebită. Pentru început ne vom ocupa de primele trei.

Data fiind fragmentarea inherentă a memoriei cea mai convenabilă structură de date pentru regăsirea zonelor libere este *lista înlanțuită*. Fiecare nod al listei va descrie o zonă de memorie liberă specificându-i adresa de început, lungimea și adresa nodului următor. Cum însă această listă înlanțuită trebuie și ea să fie stocată în undeva în memorie, cea mai convenabilă soluție este ca fiecare nod să fie stocat la începutul zonei de memorie pe care o descrie. În această abordare, un nod va conține doar lungimea zonei de memorie și adresa următoarei zone libere (adresa următorului nod). Un nod nu va mai conține adresa de început a zonei la care se referă pentru că este implicită prin adresa lui de memorie. Un nod al acestei liste înlanțuite se numește în literatura de specialitate *cuvânt de control*.

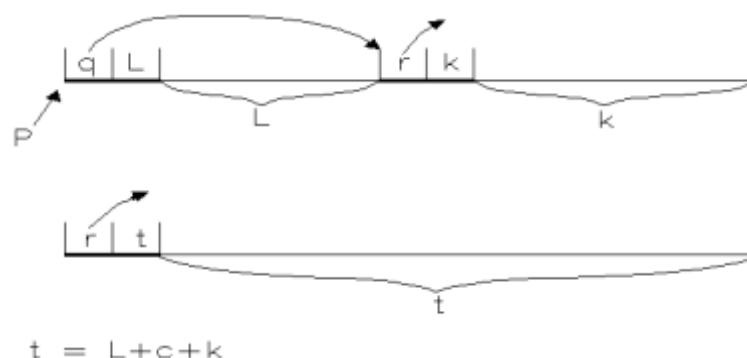
Vom prezenta în continuare câțiva algoritmi de plasare a cererilor de memorie. Pentru a trata unitar alocarea și eliberarea octeților, vom adopta următoarea convenție. Fiecare zonă liberă sau ocupată de memorie începe cu un cuvânt de control. Câmpul **lung** (aflat în a doua jumătate a cuvântului de control) al lui indică numărul de octeți liberi de *după* cuvântul de control. Pentru zonele libere de memorie pointerul **next** (aflat în prima jumătate a cuvântului de control) indică următoarea zonă liberă. Pentru zonele de memorie ocupate acest pointer nu este folosit. În fig. 10.17, 10.18 și 10.19, cuvintele de control sunt subliniate cu linii îngroșate, iar pointerii la zone sunt indicați prin săgeți. De asemenea, lungimea cuvântului de control este notată cu  $c$  și este o constantă specifică sistemului de operare. După cum vom vedea apar situații în care numărul de octeți alocați depășește (cu maximum  $c$ ) numărul de octeți solicitați. O zonă ocupată este reperată *după* cuvântul ei de control.



**Figura 10.17 Alocarea de octeți într-o zonă liberă**

În fig. 10.17a este prezentată o zonă liberă al cărei cuvânt de control începe la adresa  $q$  și are  $L$  octeți liberi. Presupunem că se cere alocarea de  $n$  octeți și că  $L > n$ . Ca rezultat al alocării, va apare una dintre situațiile din fig. 10.17b sau 10.17c. Zona hașurată reprezintă octeții ceruți pentru alocare. Situația din fig. 10.17b apare dacă  $L - n > c + 1$ , adică spațiul rămas în zonă permite crearea unei zone libere de cel puțin un octet. În cazul în care această condiție nu este satisfăcută, se alocă întreaga zonă și ultimii  $L - n$  octeți rămân nefolosiți fig. 10.17c).

Înainte de a descrie procedura inversă (de eliberare), să ne ocupăm de problema *comasării a două zone libere adiacente*. Pentru simplificare, vom presupune că lista zonelor libere este păstrată în ordinea crescătoare a adreselor. Situația în care este posibilă concatenarea a două zone libere este ilustrată în fig. 10.18.



**Figura 10.18 Comasarea a două zone libere adiacente**

În fig. 10.19 este ilustrat un exemplu de eliberare a unei zone. La eliberare, partiția devenită liberă se repune în lista de spații libere. De asemenea, se verifică dacă nu cumva partiția proaspăt eliberată *poate fi comasată cu una vecină din dreapta sau din stânga ei*. Verificând sistematic, la fiecare eliberare, dacă este posibilă concatenarea, *nu vor exista două zone libere adiacente*.

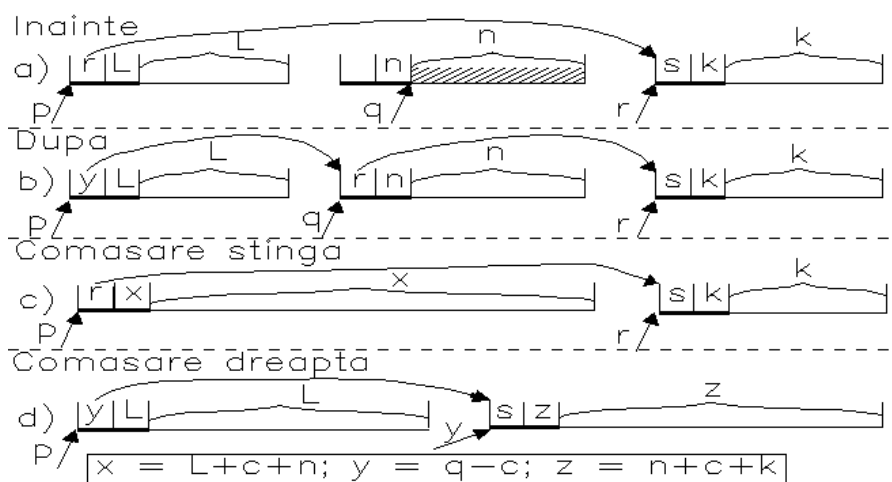
În fig. 10.19a este prezentată situația înainte de eliberarea zonei. În fig. 10.19b se prezintă situația în care nu este posibilă nici un fel de comasare. Posibilitatea de *comasare la stânga* are loc dacă:

$$p + 2 * c = q$$

Ca rezultat se obține zona din fig. 10.19c. O *comasare la dreapta* are loc dacă:

$$q + n = r$$

Ca efect se obține zona din fig. 10.19d.



**Figura 10.19 Eliberarea unei zone ocupate**

Acum este momentul să prezentăm cele patru politici de plasare. Primele trei le vom detalia utilizând descrierile procedurilor de mai sus. Punctul de pornire al fiecărui astfel de algoritm este faptul că în momentul pornirii sistemului, se alocă întreaga zonă de memorie sub forma unei singure zone libere.

#### 10.4.2.2 Metoda primei potriviri (First-fit).

Esența metodei constă în aceea că partiția solicitată este alocată în prima zonă liberă, în care încap. Principalul avantaj al metodei este simplitatea căutării de spațiu liber. Pentru această metodă, structura listei de spații libere prezentate mai sus este cea mai adecvată.

#### 10.4.2.3 Metoda celei mai bune potriviri (Best-fit).

Esența metodei constă în căutarea acelei zone libere, care lasă după alocare cel mai puțin spațiu liber. Metoda Best-Fit a fost larg utilizată mulți ani. Ea pare a fi destul de bună, deoarece economisește zonele de memorie mai mari astfel încât dacă ulterior va fi nevoie de ele vor fi disponibile. Există însă și obiecții, dintre care amintim două: timpul suplimentar de căutare și proliferarea blocurilor libere de lungime mică, adică *fragmentarea internă excesivă* (vezi 10.2.4).

Primul neajuns este eliminat parțial dacă lista de spații libere se păstrează, nu în ordinea crescătoare a adreselor, ci *în ordinea crescătoare a lungimilor spațiilor libere*. Din păcate, în acest caz problema comasării zonelor libere adiacente se complică foarte mult.

#### 10.4.2.4 Metoda celei mai rele potriviri (Worst-fit).

Metoda Worst-Fit este oarecum duală metodei Best-Fit. Esența ei constă în căutarea acelei zone libere care lasă după alocare cel mai mult spațiu liber.

Deși, numele ei sugerează că este vorba de o metodă slabă, în realitate nu este chiar așa. Faptul, că după alocare rămâne un spațiu liber mare, este benefic, deoarece în spațiul rămas poate fi plasată, în viitor, o altă partiție. Fragmentarea internă probabil că nu evoluează prea rapid, însă timpul de căutare este mai mare decât cel de la metoda primei potriviri.

Si în acest caz este posibil ca lista de spații libere să se păstreze nu în ordinea crescătoare a adreselor, ci *în ordinea descrescătoare a lungimilor spațiilor libere*. Din păcate, în acest caz problema comasării zonelor libere adiacente se complică, de asemenea, foarte mult.

#### 10.4.2.5 Metoda alocării prin camarazi (Buddy-system).

Metoda alocării prin camarazi (Buddy) este deosebit de interesantă. Ea exploatează reprezentarea binară a adreselor și faptul că din rațiuni tehnologice, dimensiunea memoriei interne este un multiplu al unei puteri a lui doi. Fie  $c \cdot 2^n$  dimensiunea memoriei interne. De exemplu, memoria unui IBM PC XT (unul dintre primele calculatoare personale) era 640 Ko, adică  $10 \cdot 2^{16}$  octeți.

Să notăm cu  $n$  cea mai mare putere a lui 2 prin care se poate exprima dimensiunea memoriei interne. În exemplul de mai sus  $n = 16$ . Din rațiuni practice, se stabilește ca unitate de alocare a memoriei tot o putere a lui 2. Fie  $m$  această putere a lui 2. Pentru exemplul de mai sus să considerăm că unitatea de alocare este 256 octeți, adică  $2^8$ .

La sistemele Buddy, dimensiunile spațiilor ocupate și a celor libere sunt de forma  $2^k$ , unde  $m \leq k \leq n$ . Ideea fundamentală este de a păstra liste separate de spații disponibile pentru

fiecare dimensiune  $2^k$  dintre cele de mai sus. Vor exista astfel  $n-m+1$  liste de spații disponibile. In exemplul de mai sus, vom avea **9** liste: lista de ordin **8** având dimensiunea unui spațiu de **256** octeți, lista de ordin **9** cu spații de dimensiune **512** etc. Ultima listă va fi de ordinul **16** și poate avea maximum **10** spații a câte **65536** ( $2^{16}$ ) octeți fiecare.

Prin definiție, fiecare spațiu liber sau ocupat de dimensiune  $2^k$  are adresa de început un multiplu de  $2^k$ .

Tot prin definiție, două spații libere de ordinul  $k$  se numesc *camarazi* (*Buddy*) de ordin  $k$ , dacă adresele lor  $A_1$  și  $A_2$  verifică:

$$A_1 < A_2, A_2 = A_1 + 2^k \text{ și } A_1 \bmod 2^{(k+1)} = 0$$

sau

$$A_2 < A_1, A_1 = A_2 + 2^k \text{ și } A_2 \bmod 2^{(k+1)} = 0$$

Atunci când într-o listă de ordin  $k$  apar doi camarazi, sistemul îi concatenează într-un spațiu de dimensiune  $2^{(k+1)}$ .

*Alocarea într-un sistem Buddy* se desfășoară astfel:

- Se determină cel mai mic număr  $p$ , cu  $m \leq p \leq n$ , pentru care numărul  $o$  de octeți solicitați verifică:  $o \leq 2^p$
- Se caută, în această ordine, în listele de ordin  $p, p+1, p+2, \dots, n$  o zonă liberă de dimensiune cel puțin  $o$ .
- Dacă se găsește o zonă de ordin  $p$ , atunci aceasta este alocată și se șterge din lista de ordinul  $p$ .
- Dacă se găsește o zonă de ordin  $k > p$ , atunci se alocă primii  $2^p$  octeți, se șterge zona din lista de ordin  $k$  și se crează în schimb alte  $k-p$  zone libere, având dimensiunile:

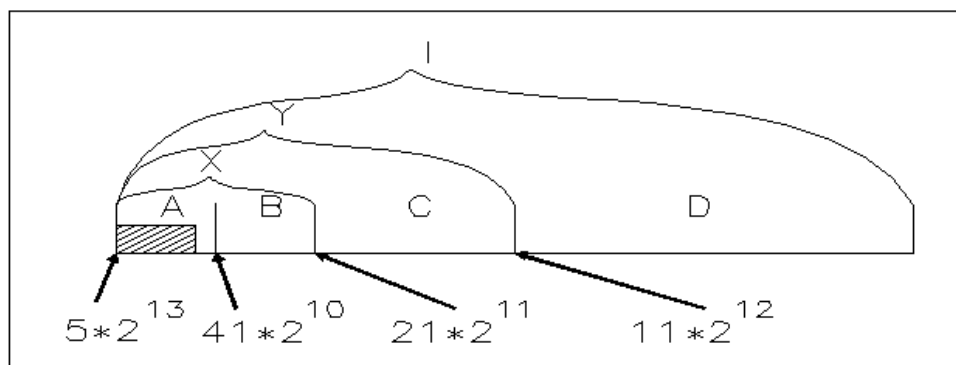
$$2^p, 2^{(p+1)}, \dots, 2^{(k-1)}$$

In fig. 10.20 este dat un astfel de exemplu. Se dorește alocarea a **1000** octeți, deci  $p = 10$ . Nu s-au găsit zone libere nici de dimensiune  $2^{10}$ , nici  $2^{11}$  și nici  $2^{12}$ . Prima zonă liberă de dimensiune  $2^{13}$  are adresa de început  $5 \cdot 2^{13}$  și am notat-o cu **I** în fig. 10.20. Ca rezultat al alocării a fost ocupată zona **A** de dimensiune  $2^{10}$  și au fost create încă trei zone libere: **B** de dimensiune  $2^{10}$ , **C** de dimensiune  $2^{11}$  și **D** de dimensiune  $2^{12}$ . Zonele **B**, **C**, **D** se trec respectiv în listele de ordine **10**, **11** și **12**, iar zona **I** se șterge din lista de ordin **13**.

*Eliberarea într-un sistem Buddy* a unei zone de dimensiune  $2^p$  este un proces invers alocării. Astfel:

1. Se introduce zona respectivă în lista de ordin  $p$ .
2. Se verifică dacă zona eliberată are un camarad de ordin  $p$ . Dacă da, atunci zona este comasată cu acest camarad și formează împreună o zonă liberă de dimensiune  $2^{(p+1)}$ . Atât zona eliberată cât și camaradul ei se șterg din lista de ordinul  $p$ , iar zona nou apărută se trece în lista de ordin  $p+1$ .
3. Se execută pasul 2 în mod repetat, mărinde de fiecare dată  $p$  cu o unitate, până când nu se mai pot face comasări.





**Figura 10.20 Alocare în sistem Buddy**

De exemplu, să presupunem că în fig. 10.20 sunt libere la un moment dat zonele **A**, **C**, **D**, iar zona **B** este ocupată. La momentul următor, se eliberează și zona **B**. În conformitate cu pașii descriși mai sus, se execută următoarele acțiuni:

- Se trece zona **B** în lista de ordin **10**.
- Se depistează că zonele **A** și **B** sunt camarazi. Drept urmare, cele două zone sunt comasate și formează o nouă zonă **X**. Zona **X** se trece în lista de ordin **11**, iar zonele **A** și **B** se șterg din lista de ordin **10**.
- Se depistează că zonele **X** și **C** sunt camarazi. Drept urmare, ele sunt comasate și formează o zonă **Y** care se trece în lista de ordin **12**, înlocuind zonele **X** și **C** din lista de ordin **11**.
- În sfârșit, se depistează că **Y** și **D** sunt camarazi. Ele sunt șterse din lista de ordin **12**, iar în lista de ordin **13** se introduce rezultatul comasării lor.

Alocarea Buddy are multe avantaje. Unul dintre ele, deloc de neglijat, se referă la manipularea comodă a adreselor de zone. Se știe că un număr binar este multiplu de **2** dacă el se termină cu **k** zerouri binare. Rezultă deci că adresele a doi camarazi de ordin **k** diferă doar prin bitul de pe poziția **k** (numerotarea începând cu 0) și ambele se termină prin **k** zerouri. Spre exemplu, adresele zonelor **A** și **B** din fig. 10.20 sunt:

```
A:  10100000000000000000
B:  10100100000000000000
      ^
```

Evident, astfel de teste se fac deosebit de ușor folosind instrucțiunile mașină care operează asupra biților.

Până în prezent nu s-a găsit un criteriu solid de comparare a acestor patru metode de plasare. Compararea lor se face empiric, eventual prin simulare. De multe ori, se adoptă o metodă mai simplă, poate și numai din rațiuni tehnologice. În orice caz, trebuie necondiționat să se aibă în vedere și concluziile, mai mult sau mai puțin empirice, rezultate din experiență. Una dintre acestea spune că nu întotdeauna o metodă mai sofisticată este și mai bună!

#### **10.4.3 Politici de încărcare.**

La sistemele cu alocare paginată, în momentul lansării în execuție a programului acesta nu are nici o pagină în memorie. La prima solicitare a programului, sistemul de operare îi va aduce în memorie numai pagina solicitată. Dacă este vorba de un program mare, acesta va funcționa normal un timp, după care va cere din nou o pagină care nu este în memorie etc. Întrebarea

care se pune este: *când* să se aducă o anumită pagină în memorie, pentru ca cererile de pagini să se reducă?

O soluție simplă, dar evident ineficientă, este *încărcarea la început a tuturor paginilor*. Prin aceasta va dispărea însuși efectul mecanismului de paginare! O altă modalitate constă în aducerea unei pagini *la cerere*, adică atunci când este ea solicitată. Această modalitate pare a fi cea mai naturală, și ea este într-adevăr și cea utilizată în sistemele de operare moderne.

Există însă metode de încărcare prin care se aduc pagini *în avans*. Astfel, odată cu o pagină se aduc și câteva pagini vecine, în ipoteza că ele vor fi invocate în viitorul apropiat. O evidență statistică a utilizării paginilor, poate furniza, cu o oarecare probabilitate, care ar fi paginile invocabile în viitor. Dacă se poate, acestea sunt aduse în avans în memoria operativă.

Legat de încărcarea în avans, încă din 1968, P.J. Denning [10] a emis *principiul vecinătății*: adresele de memorie solicitate de un program nu se distribuie uniform pe întreaga memorie folosită, ci se grupează în jurul unor centre. Apelurile în apropierea acestor centre sunt mult mai frecvente decât apelurile de la un centru la altul.

Acest principiu sugerează o politică simplă de încărcare în avans a unor pagini. Se stabilește o așa zisă *memorie de lucru* [48] compusă din câteva pagini. Atunci când se cere aducerea unei pagini de pe disc, în memoria de lucru sunt încărcate câteva pagini vecine acesteia. În conformitate cu principiul vecinătății, este foarte probabil ca următoarele referiri să fie făcute în cadrul memoriei de lucru.

#### 10.4.4 Politici de înlocuire.

În mod natural, numărul total al paginilor programelor active poate deveni mai mare decât numărul paginilor fizice din memoria operativă. Din această cauză, uneori apare situația că un program cere încărcarea unei pagini virtuale, dar nu există o pagină fizică disponibilă pentru a o găzdui. Întrebarea este *care* dintre paginile fizice va fi evacuată pentru a crea spațiul necesar?

Răspunsul optim este simplu, dar imposibil de realizat: *se evacuează acea pagină care va fi solicitată în viitor cel mai târziu* !?! Desigur, acest lucru nu poate fi prevăzut în prealabil, dat fiind faptul că evoluția unui program la un moment dat este dependentă de datele concrete asupra cărora operează. Totuși, Belady [4] a descris un model de evidență statistică prin care se poate prevedea cu o oarecare probabilitate care este pagina care va fi solicitată cel mai târziu.

Dintre metodele mai "ortodoxe", de înlocuire, descrise printre altele în [22], [10], noi vom detalia trei:

- înlocuirea unei pagini care nu a fost recent utilizată (NRU - Not Recently Used);
- înlocuirea în ordinea încărcării paginilor (FIFO - First In First Out);
- înlocuirea paginii nesolicitate cel mai mult timp (LRU - Least Recently Used).

Evident, metodele de înlocuire prezentate mai sus sunt foarte simplu de descris. Pentru implementare trebuie avut în vedere faptul că întreținerea unei structuri de date care să permită decizia trebuie făcută *la fiecare acces la memorie*. În această situație, nu este permisă nici măcar întreținerea unei liste simplu înlănțuite! De cele mai multe ori aceste metode se implementează prin hard, făcându-se uneori compromisuri.

#### 10.4.4.1 Metoda NRU.

Fiecare pagină fizică are asociați doi biți, prin intermediul cărora se va decide pagina de evacuat. Bitul **R**, numit bit de *referire*, primește valoarea 0 la încărcarea paginii. La fiecare referire a paginii, acest bit este pus pe 1. Periodic (de obicei la 20 milisecunde), bitul este pus iarăși pe 0. Bitul **M**, numit bit de *modificare*, primește valoarea 0 la încărcarea paginii. El este modificat numai la scrierea în pagină, când i se dă valoarea 1. Acești doi biți împart în fiecare moment paginile fizice în patru clase:

0. clasa 0: pagini nereferite și nemodificate;
1. clasa 1: pagini nereferite (în intervalul fixat), dar modificate de la încărcarea lor;
2. clasa 2: pagini referite dar nemodificate;
3. clasa 3: pagini referite și modificate.

Atunci când o pagină trebuie înlocuită, pagina "victimă" se caută mai întâi în clasa 0, apoi în clasa 1, apoi în clasa 2 și în sfârșit în clasa 3. Dacă pagina de înlocuit este în clasa 1 sau clasa 3, conținutul ei va fi salvat pe disc înaintea înlocuirii. Acest algoritm simplu, deși nu este optimal, s-a dovedit în practică a fi foarte eficient.

#### 10.4.4.2 Metoda FIFO.

Implementarea acestei metode este foarte simplă. Se crează și se întreține o listă a paginilor în ordinea încărcării lor. Această listă se actualizează *la fiecare nouă încărcare de pagină* (nu la fiecare acces la memorie!). Atunci când se cere înlocuirea este substituită prima (cea mai veche) pagină din listă. Bitul **M** de modificare indică dacă pagina trebuie sau nu salvată înaintea înlocuirii.

O primă îmbunătățire ar fi combinarea algoritmilor NRU și FIFO, în sensul că se aplică întâi NRU, iar în cadrul aceleiași clase se aplică FIFO.

O altă îmbunătățire este cunoscută sub numele de *metoda celei de-a doua șanse*. La această metodă, se testează bitul **R** de referință la pagina cea mai veche. Dacă acesta este 0, atunci pagina este înlocuită imediat. Dacă este 1, atunci este pus pe 0 și pagina este pusă ultima în listă, ca și cum ar fi intrat recent în memorie. Apoi căutarea se reia cu o nouă listă. Evident, șansa de scăpare a unei pagini victimă este să existe o pagină mai "tânără" ca ea și care să nu fi fost referită.

Si acum o curiozitate! Bunul simț ne spune că șansa ca o pagină să fie înlocuită scade pe măsură ce numărul de pagini fizice crește. Si totuși nu este așa! În [4] [22] este dat un contraexemplu, cunoscut sub numele de *anomia lui Belady*, pe care îl ilustrăm în fig. 10.21a și 10.21b. Este vorba de un program care are 5 pagini virtuale, pe care le solicită în ordinea:

0 1 2 3 0 1 4 0 1 2 3 4

În fig. 10.21a este ilustrată evoluția când există 3 pagini fizice, iar în fig. 10.21b când există 4 pagini fizice.

a) Trei pagini fizice:

Solicitare pagina: 0 1 2 3 0 1 4 0 1 2 3 4

Pagina recenta : 0 1 2 3 0 1 4 4 4 2 3 3

. 0 1 2 3 0 1 1 1 4 2 2

Pagina mai veche : . . 0 1 2 3 0 0 0 1 4 4

Inlocuire pagina : I I I I I I I I I = 9

b) Patru pagini fizice:

Solicitare pagina: 0 1 2 3 0 1 4 0 1 2 3 4

Pagina recenta : 0 1 2 3 3 3 4 0 1 2 3 4

. 0 1 2 2 2 3 4 0 1 2 3

. . 0 1 1 1 2 3 4 0 1 2

Pagina mai veche : . . . 0 0 0 1 2 3 4 0 1

Inlocuire pagina : I I I I I I I I I = 10

**Figura 10.21 Anomalia lui Belady**

#### 10.4.4.3 Metoda LRU.

LRU (pagina mai puțin folosită în ultimul timp) este un algoritm bun de înlocuire. El are la bază următoarea observație, reieșită (tot) din principiul vecinătății. O pagină care a fost solicitată mult de către ultimele instrucțiuni, va fi probabil solicitată mult și în continuare. Invers, o pagină solicitată puțin (sau deloc), va rămâne probabil tot așa pentru câteva instrucțiuni.

Problema este cum să se țină evidența utilizărilor? Se *exclude* din start întreținerea unei liste înlănțuite care să fie modificată la fiecare acces la memorie. Prețul plătit este mult prea mare. Iată două posibile rezolvări.

*Numărătorul de accese* se implementează hard. Există un registru numit *contor* reprezentat (de regulă) pe 64 de biți. La fiecare acces, valoarea lui este mărită cu o unitate. În tabela de pagini, există câte un spațiu rezervat pentru a memora valoarea contorului. În momentul accesului la o pagină, valoarea contorului este memorată în acest spațiu rezervat din tabela de pagini. Atunci când se impune o înlocuire, este înlocuită pagina care a reținut cea mai mică valoare a contorului.

*Matricea de referințe.* Pentru un sistem de calcul care are  $n$  pagini fizice, se utilizează o matrice binară de  $n \times n$ . La pornire, toate elementele au valoarea 0. Atunci când se face referire la o pagină  $k$ , linia  $k$  a matricei este înlocuită peste tot cu 1, după care coloana  $k$  este înlocuită peste tot cu 0. În fiecare moment, numărul de cifre 1 de pe o linie oarecare  $l$  arată de câte ori a fost referită pagina  $l$  după încărcare.

Iată, spre exemplu, în fig. 10.22, cum arată evoluția matricei de referințe într-un sistem de calcul care are patru pagini fizice, solicitate în ordinea:

0 1 2 3 2 1 0 3 2 3:

0 1 2 3 2 1 0 3 2 3

```

0111 0011 0001 0000 0000 0000 0111 0110 0100 0100
0000 1011 1001 1000 1000 1011 0011 0010 0000 0000
0000 0000 1101 1100 1101 1001 0001 0000 1101 1100
0000 0000 0000 1110 1100 1000 0000 1110 1100 1110

```

**Figura 10.22 Evoluția unei matrice de referințe**

Implementarea mecanismului matricei de referințe se face destul de ușor prin hard, în orice caz mult mai ușor decât implementarea mecanismului cu numărător de referințe.

Desigur, politicile de înlocuire descrise mai sus pot fi simulate foarte bine prin soft. Din păcate, eficiența mecanismelor scade drastic. Acesta este motivul pentru care unele implementări ale memoriei virtuale (nu dăm aici nume) au eșuat lamentabil. Nu-i nimic, s-a câștigat experiență și asta nu este lucru puțin.

#### 10.4.5 Cum funcționează o memorie cache?

Apariția memoriei *cache* a fost dictată de necesitatea creșterii performanțelor sistemului de calcul. Memoria cache conține copii ale unor blocuri din memoria operativă. Când procesorul încearcă citirea unui cuvânt de memorie, se verifică dacă acesta există în memoria cache. Dacă există, atunci el este livrat procesorului. Dacă nu, atunci el este căutat în memoria operativă, este adus în memoria cache împreună cu blocul din care face parte, după care este livrat procesorului. Datorită vitezei mult mai mari de acces la memoria cache, randamentul general al sistemului crește.

Aici apar o serie de probleme: cât de mare este o memorie cache? care este dimensiunea optimă a blocului de memorie destinat schimbului cu această memorie? În lucrarea [42] se dau răspunsuri la aceste întrebări. Tot acolo se fac convenite evaluări ale raportului cost/performanță pentru dimensionarea unei memorii cache, precum și care este probabilitatea (hit ratio) ca o adresă cerută de procesor să fie găsită în memoria cache.

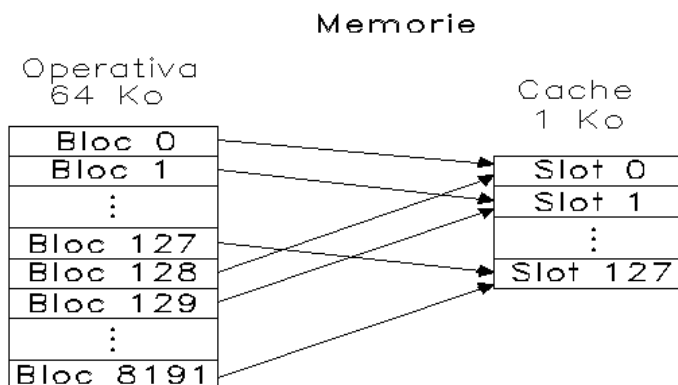
Noi ne rezumăm la ilustrarea modului în care se face corespondența dintre memoria operativă și memoria cache, precum și politicile de schimb dintre cele două tipuri de memorie.

Memoria cache este împărțită în mai multe părți egale, numite *sloturi* (poziții). Un slot are dimensiunea unui bloc de memorie, care este în mod obligatoriu o putere a lui 2. Fiecare slot conține în plus câțiva biți (numiți generic *tag* = etichetă) care indică blocul de memorie operativă depus în slotul respectiv. Dimensiunea unui tag depinde de mecanismul de schimb dintre cele două memorii, dar noi nu ne vom ocupa aici de el. De exemplu, în fig. 10.23 și 10.24 am ales o memorie operativă de 64 Ko și o memorie cache de 1 Ko. Dimensiunea unui slot, identică cu dimensiunea unui bloc de memorie, am ales-o de 8 octeți (64 de biți).

Se cunosc [42] mai multe metode de *proiectare a spațiului memoriei operative pe memoria cache*. Cea mai simplă metodă este *proiectarea directă*. Dacă  $C$  indică numărul total de sloturi din memoria cache,  $A$  este o adresă oarecare din memoria operativă, atunci numărul  $S$  al slotului în care se proiectează adresa  $A$  este:

$$S = A \bmod C$$

În fig. 10.23 este ilustrată această corespondență.



**Figura 10.23 Proiectare directă pe memoria cache**

Principalul ei avantaj este simplitatea. În schimb, există un mare dezavantaj, generat de faptul că fiecare bloc are o poziție fixă în memoria cache. Dacă, spre exemplu, se cer accese alternative la două blocuri care ocupă același slot, atunci are loc un trafic intens între cele două memorii, fapt care face să scadă mult randamentul general al sistemului.

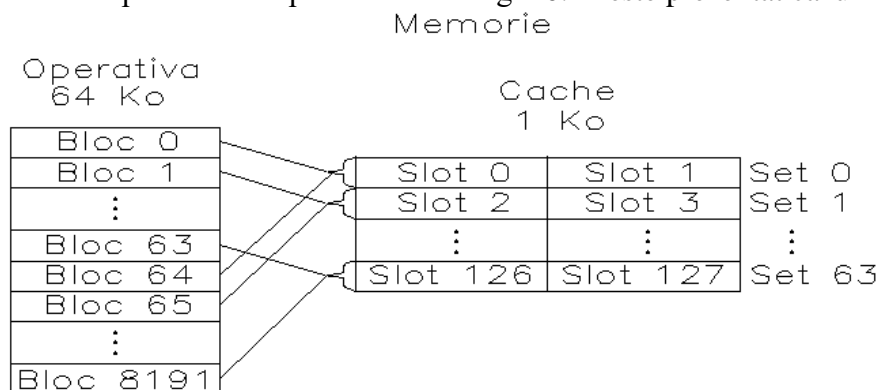
O a doua metodă poartă numele de *proiectare asociativă*, apărută pentru eliminarea dezavantajului de mai sus. Această metodă plasează un bloc de memorie într-un slot oarecare liber.

Gestiunea sloturilor memoriei cache în acest caz ridică aceleași probleme ca și alocarea paginată a memoriei. Principala problemă care se ridică aici este cea legată de *politica de înlocuire*. În paragraful 2.4.3 am tratat această problemă, descriind algoritmi NRU, FIFO și LRU. Tot ceea ce s-a spus acolo este valabil, fără modificări și la înlocuirea unui slot atunci când un bloc solicitat din memoria operativă nu mai are loc în memoria cache.

O a treia metodă poartă numele de *proiectarea set-asociativă*, și ea este o combinație a precedentelor două metode. Ideea ei este următoarea. Memoria cache este împărțită în  $I$  seturi, un set fiind compus din  $J$  sloturi. Avem deci relația  $C = I \cdot J$ . La cererea de memorie de la adresa  $A$ , se calculează numărul  $K$  al setului în care va intra blocul, astfel:

$$K = A \bmod I$$

Având fixat numărul setului, blocul va ocupa unul dintre sloturile acestui set. Alegerea slotului este de această dată o problemă de planificare. În fig. 10.24 este prezentat cazul  $I=64$  și  $J=2$ .



**Figura 10.24 Proiectarea set - asociativă**

Metoda set-asociativă este mult folosită în practică. Compromisul proiectării directe și a celei asociative face să fie aplicată o politică de înlocuire numai atunci când într-un slot există deja *J* blocuri, ceea ce se întâmplă foarte rar. În rest se aplică proiectarea directă, care este mult mai simplă.