

GESTIUNEA MEMORIEI

Gestiunea resurselor memoriei este un aspect complex al unui sistem de operare. Iată câțiva pași parcurși înspre o organizare eficientă și performantă.

Cuprins

1. **Elemente de baza in gestiunea memoriei** *Neculoiu Paul*
 - 1.1. Monoprogramarea (mono/single tasking) fara utilizarea de swap sau paginare;
 - 1.2. Multiprogramarea cu partitii fixe;
 - 1.3. Modelarea multiprogramarii;
 - 1.4. Analiza performantelor sistemelor cu multiprogramare;
 - 1.5. Protectie si relocare;
 - 1.6. Rezumat

2. **Swapping** *Raducanu Vlad*
 - 2.1. Introducere;
 - 2.2. Gestiunea memoriei cu harti de biti;
 - 2.3. Gestiunea memoriei cu liste inlantuite;
 - 2.4. Swapping in cazul Linux;
 - 2.5. Swapping in cazul Windows 2000
 - 2.6. Rezumat

3. **Memoria virtuală** *Musat Liana*
 - 3.1. Paginare ;
 - 3.2. Tabele pagină;
 - 3.2.1. Tabele de pagină cu nivel multiplu;
 - 3.2.2. Structura unei intrări în tabel;
 - 3.3. TLB-uri – Translation Lookaside Buffers;
 - 3.4. Tabele de pagină inversate.

4. **Algoritmi de înlocuire a paginilor** *Datcu Octaviana*
 - 4.1. Introducere
 - 4.2. Strategii de înlocuire a paginilor – algoritmi
 - 4.2.1. Modelul înlocuirii optime a paginilor (OPRA)
 - 4.2.2. Not Recently Used (Neutilizată recent)
 - 4.2.3. First In, First Out (primul intrat, primul ieșit)
 - 4.2.4. Least Recently Used (Cel mai recent utilizat)
 - 4.2.5. Implementarea cu stivă a LRU
 - 4.2.6. Algoritmi de aproximare LRU
 - 4.2.6.1. Second Chance (A Doua Șansă)
 - 4.2.6.2. Clock (Ceasul)
 - 4.2.6.3. Not frequently used (Neutilizată frecvent)
 - 4.2.9. Aging (Îmbătrânirea)
 - 4.2.7. Random
 - 4.3. Impactul dimensiunii paginii asupra performanțelor memoriei virtuale
 - 4.4. Comportamentul programului în cazul paginării
 - 4.5. Working Set (Setul de lucru)
 - 4.6. Fenomenul “thrashing”
 - 4.7. PFF (Frecvența de eroare de pagină)
 - 4.8. Alți algoritmi
 - 4.8.1. “Low inter-reference Recency Set Replacement Policy” LIRS
 - 4.8.2. “Enhanced second chance/clock”
 - 4.8.3. “Page Buffering”
 - 4.9. Concluzii asupra performanței principalilor algoritmi
 - 4.10. Caracteristici ale paginării în Windows XP/ Linux
 - 4.11. Referințe
-
5. **Modelarea algoritmilor de înlocuire a paginilor** *Velican Valentin*
 - 5.1. Definiție. Scurt “istoric”.
 - 5.2. Înlocuirea locală vs. Înlocuirea globală
 - 5.3. Algoritmul optim de înlocuire a paginilor. Imposibilitatea realizării sale.
 - 5.4. Algoritmi de înlocuire. Tipuri. Discuție.
 - 5.5. Concluzii
 - 5.6. Rezumat
 - 5.7. Referințe
-
6. **Probleme de proiectare pentru sisteme de paginare** *Stoian Andrei*
 - 6.1. Alocare globală vs alocare locală;
 - 6.2. Controlul alocării paginilor;
 - 6.3. Dimensiunea paginilor;
 - 6.4. Spațiul datelor și spațiul instrucțiunilor;
 - 6.5. Pagini partajate;
 - 6.6. Eliberarea (curățarea) paginilor;
 - 6.7. Interfața memoriei virtuale;

7. **Probleme de implementare** *Nistor Adrian*

- 7.1. Implicatiile Sistemului de Operare in paginare;
- 7.2. Gestionarea erorilor de paginare;
- 7.3. Backupul instructiunilor;
- 7.4. Blocarea paginilor in memorie;
- 7.5. Termenul de "Backing Store";
- 7.6. Separarea regulilor si a mecanismului;

8. **Gestionarea memoriei în UNIX** *Ledeanu Mihai Silviu*

- 8.1. Concepte fundamentale
- 8.2. Implementarea Gestionării memoriei în Unix
 - 8.2.1. Swapping
 - 8.2.2. Paginarea
 - 8.2.3. Algoritmul de înlocuire a paginilor
- 8.3. Rezumat

1 Elemente de baza in gestiunea memoriei *Neculoiu Paul*

“640K ought to be enough for anybody.” – atribuit lui Bill Gates (1981)

Adevarat sau nu citatul in 1981 acest concept nu mai este de ceva timp de actualitate. Nevoile de memorie ale aplicatiilor software din ziua de azi au depasit de mult acea valoare. Desi tehnica a evoluat o data cu necesitatea acestora de memorie totusi exista inca anumite limitari tehnice sau financiare care retin utilizatorul din a atinge utopia memoriei non-volatile infinit de mari si infinit de rapida.

Din acest motiv se pune un mare accent in ziua de azi pe gestionarea cat mai eficient a memoriei limitate aflate la dispozitia utilizatorului.

1.1. Monoprogramarea (mono/single tasking) fara utilizarea de swap sau paginare

Desi tehnica nu mai este demult utilizata in statiile din ziua de azi, inca se mai poate gasi in unele sisteme dedicate sau unitati de calcul portabile (palmtop) si printre primele statii aparute desi tehnicile de implementare pentru fiecare desi acestea difera de la caz la caz.

Primele mainframeuri aveau sistemul de operare la baza memoriei RAM, restul revenindu-i memoriei de program, primele PCuri aparute fiind echipate si cu o memorie ROM rezervata pentru drivere si intrare/iesire cunoscuta sub numele de BIOS (basic input/output system), indtrodus pentru prima data pe scara larga de catre IBM si la acea data era una dintre putinele solutii care presupunea un nivel de abstractizare independent de sistemul de operare si care avea in acelasi timp si calitatea de a fi suficient de bine documentat. Acesta deservea sistemul cu o serie de functii de vaza realiza pasii necesari incarcarii sistemului de operare in RAM si totdata punea la dispozitie unele optiuni de configurare.

Necesitatea de flexibilitate a zilei de azi dar si a programarii defectuase au dus la inlocuirea memoriei ROM ca mediu de stocare in favoare unor memorii EEPROM (sau flash-ROM) aceasta fiind mai usor de actualizat.

Alternativ sistemul de operare se poate pastra intr-o unitate ROM separata de memoria RAM de program, sistem indelung utilizat in palmtopuri si unitatile de calcul embeded. Sistemul de operare era initial “incrustat” si nu putea fi inlocuit decat prin inlocuirea fizica a ROM-ului. In ziua de azi, cu evolutia memoriilor configurabile non-volatile (flash-ROM) este posibila actualizarea continutului ROM-ului respectiv fara inlocuirea fizica a acestuia. In principiu exista doua tehnici de baza pentru a schimba continutul unui flash-ROM in palmtopuri. Prima este referita ca umbrire (“shadowing”) care implica stocarea componentelor actualizabile ale ROM-ului in zona RAM, aceasta avea totusi anumite probleme, componentele respective ocupand RAM-ul respectiv in primul rand, (unele fiind chiar prea mari pentru a incapea in acesta), in al doilea rand daca unitatea ar fi fost vreodata pornita “la rece” (situatie in care RAM-ul se goleste), actualizarile respective s-ar fi pierdut. A doua tehnica utilizata implica reinscriptionarea intregii imagini, utilizand loaderul de boot (bootstrap). Acesta intruieste initial unitatea portabila sa incarce anumite componente in ROM cand acesta este pornit (cald sau rece). Un dezavantaj cu aceasta tehnica este aceea ca majoritatea acestor sisteme au spatiu limitat de memorie si nu pot stoca intreaga imagine in timpul reinscriptarii.

Sistemul astfel organizat unitatea de calcul nu poate procesa decat un program intr-un anumit moment de timp, acesta copiindu-l de pe unitatile de stocare in memorie si executandu-l pana la finalizare, moment in care procesarea se opreste, unitatea de calcul asteptand urmatoarea comanda. La primirea acesteia incarca noul program in memorie suprascriindu-l pe cel dintai si executandu-l.

Modul monoprogramarii nu mai este folosit in ziua de azi decat in unele unitati embeded, pentru care costul de productie mai redus si/sau lipsa nevoiei pentru aplicatia respectiva a unui nivel superior de gestionare a justificat renuntarea la modurile mai avansate.

1. 2. Multiprogramarea cu partitii fixe.

Necesarul zilei de azi a dus la adptarea unor tehnici ce pot rula procese multiple simultan. Avantajul acestui sistem este ca procesorul poate prelucra continuu, putand prelua un program in timp ce alt program asteapta finalizarea procedurilor de intrare/iesire, crescand, astfel, eficienta sistemului de calcul fata de modul "clasic" de procesare fiecarui program in parte. Tehnica a fost intotdeauna utilizata in serverele de retea pentru a rula procese multiple (pentru clienti diferiti ce-i drept) simultan, dar in ziua de azi si unitatile client in sine au aceasta capacitate, monoprogramarea disparand de mult din aceasta zona de prelucrare.

Multiprogramarea este tehnica de exploatare a sistemelor care permite existenta simultana in memoria interna a mai multor programe care se executa concomitent in partitii fixe de memoriecu conditia ca acestea sa nu utilizeze in acelasi timp simultan aceeaasi resursa. Executia in multiprogramare a lucrarilor se face pe loturi, fiecare lot de lucrari avand afectata o partitie fixa din memria interna. O partitie este o zoan continua de memorie, de o lungime si adresa fixa, partitii partajate la inceput, de pilda, la pornirea sistemului.

In cadrul fiecarui lot, lucrarile sunt executate secvential, fiind lansate automat in executie. Sub controlul sistemului de operare, unitatea de prelucrare comuta de la o partitie la alta, pentru a realiza executia simultana a proceselor, comutarea facandu-se in momentul in care unitatea nu este utilizata de procesul respectiv (ex. Asteapta terminarea operatiilor de intrare/iesire), comutarea intre procese facandu-se in urma unor evenimente interne proceselor din executie. Fiecare partitie are asociata o prioritate de executie, resursele sistemului de calcul fiind alocate proceselor conform solicitarilor acestora, si in functie de disponibilitatea reurselor. In cazul solicitarii unei resurse care nu este disponibila, procesul respectiv intra in asteptare, pana la eliberarea acesteia. Oridnea de alocare a resurselor intre procesele care solicita aceeaasi resursa fiind determinata de prioritatea de partitiei (procesele cu prioritatea mai mare au acces la aceeaasi resursa inaintea celor cu prioritatea mai mica). Sistemul de calcul dispune de un sistem de intreruperi prin intermediul caruia se semnaleaza aparitia unui eveniment care poate fi cauza comutarii intre procese.

Accesul programelor la partitii pentru organizarea loturilor devine o problema in sensul ordinii acestora, dimensiunea partitiilor fiind limitata de dimensiunea totala a memoriei, aceasta putand fi divizata in portiuni inegale pentru stocarea programelor. Un algoritm utilizat este asezarea programului la coada pentru cea mai mica partitie care o poate contine. Dezavantajul unei astfel de tehnici consta in momentul in care aplicatiile mici, desi existand spatiu suficient intr-o partitie de dimensiune mai mare pt ea, asteapta loturi multiple pentru a intra in partitia ei corespunzatoare.

O alta strategie ar fi organizarea intrarii pe rand a proceselor astfel in cat acestea sa intre in prima partitie de dimensiuni suficient de mari care devine disponibila. Din nefericire, in cazul sarcinilor mici care intra in partiile mari, spatiul neutilizat se pierde, acest lucru fiind nedorit. Alternativ se poate cauta intreaga lista de intrare cand o partitie devine libera sa se aleaga cea mai mare care incapa, desi in acest mod cele mai mici risca sa nu mai apuce sa intre.

1.3. Modelarea Multiprogramarii

Prin multiprogramare se poate obtine o utilizare crescuta a procesorului. Pe scurt, daca, in medie, un proces este procesat numai 20% din timpul pe care il petrece in memorie, cu cinci procese simultane in memorie, procesorul ar trebui sa fie utilizat in permanenta. Acest model este nerealist de optimist, totusi, deoarece presupune ca toate cele cinci procese nu vor astepta niciodata aceeaasi interfata de intrare/iesire.

Un model mai bun il reprezinta privirea utilizarii procesorului dintr-un punct de vedere probabilistic. Presupunand ca un proces utilizeaza o fractiune p din timp asteptand terminarea procesului de intrare/iesire, cu n procese in memorie simultan probabilitatea ca toate cele n procese sa astepte pentru I/O (caz in care procesorul sta) este p^n . Utilizarea procesului apoi fiind $1-p^n$. Deci, daca procesul si-ar petrece 80% din timp asteptand I/O, cel putin zece procese trebuie sa fie in memorie simultan pentru a aduce risipa de procesor sub 10%, procesele interactive asteptand de exemplu ca utilizatorul sa introduca date la un terminal, 80% risipa de procesor nefiind un lucru intocmai iesit din comun.

Pentru acuratete, modelul probabilistic descris este numai o aproximare. Acesta presupune implicit ca cele n procese sunt independente, fiind suficient de acceptabil pentru un sistem cu cinci procese in memorie sa aiba trei ruland si doua in asteptare. Dar cu un singur procesor nu putem rula trei procese simpultan drept urmare un proces devenit pregatit cat timp procesorul este ocupat va trebui sa astepte. Drept urmare procesele nu sunt independente.

De asemenea, presupunand ca o masina de calcul are, spre exemplu, 32MB de memorie, sistemul de operare ocupand 16MB si fiecare program ocupand 4MB. Aceasta ar permite patru programe sa fie in memorie in acelasi timp. Cu o medie de asteptare de 80%, avem o utilizare a procesorului (ignorand necesitatile sistemului de operare) de aproximativ 60%. Adaugarea a inca 16MB de memorie permite sistemului sa treaca de la 4 programe la 8, crescand astfel utilizarea procesorului la 83%. Adaugarea a inca 16MB nu ar creste utilizarea procesorului decat de la 83% la 93%, punand astfel sub semnul intrebării, utilitatea reala adaugarii repsectivei extinderi de memorie.

1.4. Analiza performantelor sistemelor multiprogramate

Modelul discutat poate de asemenea fi utilizat pentru a analiza prelucrării pe loturi. Nu toate procesele utilizeaza procesorul pentru acelasi interval de timp. Cu o medie de asteptare de 80% un proces poate petrece de 5 ori mai mult timp in memorie fata de cat timp, utilizeaza, practic procesorul, chiar si fara concurenta pentru accesul la acesta. In cazul sosirii unei a doua sarcini, utilizarea procesorului creste datorita gradului ridicat de multiprogramare, fiecare proces avand acces jumătate de timp la procesor, pierderea de timp de acces la procesor a primei sarcini fiind redusa. In cazul sosirii unei a 3a sau a 4a sarcini timpul de acces al fiecareia scade iar utilizarea procesorului creste.

1.5. Relocarea si Protectia

Utilizarea multiprogramarii ridica doua probleme esentiale, relocarea si protectia. Sarcini diferite ruleaza de la adrese diferite. Cand programele sunt combinate sa inceapa intr-un singur spatiu de adresa, linkerul trebuie sa stie de la ce adresa va incepe programul in memorie.

Daca prima instructiune ar fi un salt la o procedura la adres absoluta 100, in cazul in care programul e incarcat in prima partitie (la adresa 100K de ex), instructiunea respectiva va sari la adresa 100 care se afla in interiorul sistemului de operare. Ceea ce este necesar este ca aceasta sa sara la 100K+100. Daca programul ar fi incarcat in partitia 2 ar trebui sa sara la 200K+100. Problema adresarii relative la locul de inceput al programului este cunoscuta ca problema relocarii.

O solutie posibila consta in modificarea instructiunilor in sine din program in timp ce acesta este incarcat in memorie. Programele incarcate intr-o partitie li se adauga adresa de inceput a acelei partitii. Performarea relocarii in timpul incarcarii in acest mod, linkerul trebuia sa atribuie programului binar o lista continand care din liniile din program sunt adrese si care sunt comenzi, constante sau alte obiecte care nu trebuie relocate.

Relocarea in timpul incarcarii nu rezolva problema protectiei. Un program malitios poate oricand construi o noua instructiune la care sa sara.

Deoarece programele in acest sistem utilizeaza adrese absolute de memorie in loc de adrese relative la registru nu exista metode de a opri un program din a construi o instructiune care citeste sau scrie orice cuvand din memorie. In sistemele multiuser, nu este de dorit a lasa procesele sa scrie si sa citeasca bucati de memorie apartinand altor utilizatori.

Solutia data de IBM a fost impartirea memoriei in blocuri de 2KB si asignarea acestora a unui cod de protectie de 4 biti pentru fiecare bloc. Acesta bloca orice tentativa a vreunui proces de accesare a memoriei a carui cod de protectie diferea de al lui. De vreme ce numai sistemul de operare putea schimba codurile de protectie procesele utilizatorilor erau prevenite din a interfera cu ale altora si cu sistemul de operare in sine.

O solutie alternativa atat la relocare cat si la problema protectiei o reprezinta echiparea sistemului cu doi registri hardware speciali denumiti base(baza) si limit(limita). Cand unui proces ii vina randul, registrul de baza este incarcat cu adresa de start a partitiei, iar registrul de limita cu lungimea acesteia. Fiecare adresa de memorie generata automat l se adauga continutul registrului de baza inainte de a fi trimisa in memorie.

Astfel, instructiunea in sine nu mai este modificata. Adresele sunt de asemenea verificate sa nu depaseasca limita, astfel incat sa nu isi paraseasca partitia curenta. Registrele sunt protejate la scriere din partea utilizatorilor.

Dezavantajul acestei scheme este nevoia de a aduna si compara fiecare referinta de memorie. Comparatia se face rapid dar adunarile sunt incete datorita propagarii bitului de carry decat daca nu sunt utilizate circuite speciale de adunare.

CDC 6600 a fost primul supercomputer din lume si utiliza aceasta tehnica. Procesoarele intel 8088 utilizate pentru primele IBM PC foloseau o versiune mai slaba a acesteia, doar registre de baza dar fara cele de limita. Putine calculatoare o mai folosesc in ziua de azi.

1.6. Rezumat

monoprogramarea (mono/single tasking) fara utilizarea de swap sau paginare

Tehnica monoprogamarii inca se mai gaseste in uz in unele sisteme dedicata sau unitati de calcul portabile nemaifiind utilizata in statii de cateva generatii.

Informatiile erau stocate in functie de necesitati fie in RAM-ul temporar sau in memorii mai nevolatile ROM, necesitatile curente inlocuind in timp pe acestea din urma cu variante mai performante.

Sistemul in schimb nu poate procesa decat un program intr-un anumit moment de timp, acesta nemaifiind utilizat in ziua de azi decat in unele aplicatii restranse.

Multiprogramarea cu partitii fixe.

Evolutia de la monoprogamare, multiprogramarea aduce avantajul ca poate prelucra multiple procese simultan, procesorul putand lucra continuu putand prelua un program in timp ce altul asteapta finalizarea procedurilor de intrare/iesire ducand la o eficienta mai mare.

Aceasta permite existenta simultana in memoria interna a mai multor programe care se eexecuta concomitent in partitii fixe de memorie cu conditia ca acestea sa nu utilizeze in acelasi timp simultan aceeasi resursa.

Modelarea Multiprogamarii.

Prin multiprogamare se poate obtine o utilizare crescuta a procesorului. Se pot gasi multiple metode de modelare a acesteia pentru determinarea performantelor fiecare cu partile lui slabe si tari.

Este nevoie a se tine cont de faptul ca procesele nu ruleaza simultan in mod real ci intra succesiv pe bucati in procesor, "mimand" astfel procesul de rulare simultana.

Analiza performantelor sistemelor multiprogramate

Modelul discutat poate de asemenea fi utilizat pentru a analiza prelucrarii pe loturi. Nu toate procesele utilizeaza procesorul pentru acelasi interval de timp. In cazul sosirii unei a doua sarcini, utilizarea procesorului creste datorita gradului ridicat de multiprogamare, fiecare proces avand acces jumătate de timp la procesor, pierderea de timp de acces la procesor a primei sarcini fiind redusa.

Relocarea si Protectia

Utilizarea multiprogamarii ridica doua probleme esentiale, relocarea si protectia. Sarcini diferite ruleaza de la adrese diferite. Cand programele sunt combinate sa inceapa intr-un singur spatiu de adresa, linkerul trebuie sa stie de la ce adresa va incepa programul in memorie.

2.1. Introducere

Exista mai multe metode de gestiune a memoriei, dintre acestea metodele cu partitii fixe pot fi utile in cazul unui sistem cu prelucrare pe loturi. Acest tip de metode poate fi eficient in cazul in care memoria este suficient de mare ca sa tina suficient de multe procese astfel incat procesorul sa fie ocupat cat mai mult timp.

In sisteme multiutilizator cu partajare a timpului, sau in computere care au de indeplinit sarcini consumatoare de memorie apare problema ca memoria principala nu este suficienta. Solutia in acest caz este ca o parte din procese sa fie stocate pe harddisk si aduse in memoria principala atunci cand este nevoie. Se folosesc doua strategii generale in acest caz: cea mai simpla, numita swapping, care consta in mutarea proceselor cu totul intre memorie si hard, si cealalta, memoria virtuala, care permite stocarea unor parti din proces pe hard si altora in memoria principala.

Avantajul la swapping este practic ca partiile memoriei sunt de marime variabila, in functie de process si astfel memoria poate fi utilizata mai eficient, principalul dezavantaj insa este ca se complica algoritmul de alocare a memoriei.

Dupa folosirea swappingului un timp e posibil sa se creeze mai multe zone de memorie libera care sa fie prea mici pentru a putea incarca un proces intreg in ele, pentru a rezolva aceasta problema se poate compactata memoria astfel incat sa avem o singura zona de memorie libera foarte mare si compacta. Insa aceasta actiune este o mare consumatoare de timp si nu este in general utilizata.

Spatiul alocat pentru procese noi, sau aduse de pe harddisk trebuie sa fie mai mare decat procesul in sine, deoarece procesele au tendinta sa creasca in timpul rularii, de asemenea este util sa se aloce acest spatiu in plus intre o zona de stiva (care creste in jos) si o zona de date (care creste in sus) astfel incat utilizarea acestui spatiu sa fie cat mai eficienta.

2.2. Gestiunea memoriei cu harti de biti (bitmaps)

O metoda de gestiune a memoriei in cazul swappingului este cu ajutorul hartilor de biti (bitmaps). Metoda presupune creerea unei harti in care fiecare bit corespunde unei zone de memorie de dimensiune fixa, bitul fiind 0 pentru o zona libera si 1 pentru o zona ocupata. Problema principala in acest caz este alegerea dimensiunii zonei de memorie corespunzatoare unui bit, deoarece alegerea unei dimensiuni prea mari duce la pierderea de spatiu, iar alegerea unei dimensiuni prea mici duce la marirea inutila a hartii memoriei. In cazul in care zonele de alocare au dimensiune prea mare in momentul in care o parte dintr-o zona este utilizata, in harta memoriei intreaga zona este marcata ca "folosita" si astfel se pierde restul de memorie din zona respectiva.

O alta problema cu aceasta metoda este ca atunci cand trebuie adus un proces in memorie trebuie cautata o secventa de biti 0 (corespunzatoare unei secvente de zone libere in memorie) suficient de mare cat sa poate tine procesul respectiv, cautarea unei astfel de secvente intr-o harta de biti este foarte consumator de timp.

2.3. Gestiunea memoriei cu liste inlantuite

Utilizarea listelor inlantuite pentru gestiunea memoriei presupune creerea unei liste in care fiecare element se refera la o zona de memorie astfel: zona poate fi ocupata sau libera, se retine adresa de inceput a zonei si adresa de sfarsit a zonei si un pointer catre elementele vecine. Lista de obicei este ordonata dupa adrese astfel incat atunci cand trebuie modificata aceasta sa se faca usor : atunci cand se dezaloca memorie dintr-o zona elementele corespunzatoare zonelor vecine libere impreuna cu elementul respectiv se combina formand un nou element cu adresa de inceput a primului si adresa de final a ultimului si marcat ca liber.

Pentru a aduce un proces din memorie trebuie cautata o zona de memorie adecvata, pentru aceasta exista mai multi algoritmi cum ar fi first fit (care alege prima zona suficient de mare pentru a tine procesul respectiv, dupa care la urmatorul proces cautarea se face din nou de la inceput), next fit (alege prima zona suficient de mare, dupa care la urmatorul proces cautarea se face incepand cu ultima zona unde s-a ajuns), best-fit (se cauta prin toata lista pana se gaseste cea mai mica zona libera capabila sa acomodeze procesul respectiv), worst fit (se cauta prin toata lista si se alege cea mai mare zona libera). Ideea la worst fit este ca in urma celorlalti algoritmi ajunge la un moment dat sa se creeze niste zone libere prea mici pentru a fi utile si in situatia lui worst fit spatiul liber ramas va fi cat mai mare posibil. In practica insa se dovedeste ca next fit e chiar ceva mai slab decat first fit, ca best fit e chiar mai slab deoarece dureaza f. Mult si creeaza si zonele mici de mem. libera si ca nici worst-fit nu este atat de bun.

O posibila imbunatatire este sa fie tinute 2 liste separate una pentru zonele alocate, alta pentru zonele libere, astfel cautarea se face doar in lista de zone libere si se imbunatatesc performanta algoritmilor. Folosind aceasta metoda nici nu mai e nevoie pentru o structura separata pentru lista, ci ar putea pur si simplu fi implementata direct in memorie, astfel la inceputul fiecarei zone sa fie un mic element care sa descrie zona, inclusiv cu pointer catre urmatoarea zona libera si astfel se poate renunta la informatia de liber/ocupat folosind doar 2 cuvinte pentru un element din lista, in loc de 3.

O alta imbunatatire se numeste quick-fit si consta in a tine o lista separata cu zonele libere de marimi uzuale (zone care sunt mai probabil sa fie cerute). Pentru procese ce necesita zone de marimi mai ciudate se pot incarca in zone un pic mai mari, sau se poate tine o lista separata pentru zonele de marimi ciudate. Oricum pentru acest tip de swapping apar din nou problemele de fragmentare.

2.4. Swapping in cazul Linux

Memoria in cazul sistemului de operare Linux este organizata intr-un spatiu virtual de memorie pentru fiecare proces consistand din trei segmente. Segmentul text, in care este tinut "codul" procesului, acest segment este de obicei "read-only", un segment de date in care sunt tinute variabilele si datele necesare programului, acest segment este impartit in doua o parte pentru datele neinitializate si o parte pentru datele initializate, si un segment de stiva care creste in jos (incepe de la adrese mari si creste catre 0), programele nu pot gestiona explicit marimea acestui segment si in momentul in care se depaseste marimea acestui segment apare o eroare hard si segmentul este automat marit cu o pagina. Segmentul de date poate creste si se poate micșora si exista o functie de sistem cu ajutorul careia programele pot aloca sau dealoca memorie (brk).

Swappingul in Linux este comandat de nivelul superior al unui programator pe 2 niveluri, anume swapperul. Mutarea din memorie spre harddisk este initiata atunci cand nu mai e suficienta memorie libera dintr-unul din mai multe motive, de exemplu este apelata functia brk pentru a extinde un segment de date, sau este apelata functia fork si e nevoie de memorie pentru procesul fiu sau o stiva creste mai mult decat spatiul alocat ei. Pentru mutarea inversa, anume trebuie adus in memorie un proces care a stat prea mult pe harddisk gasirea unui proces care sa fie trimis spre harddisk pentru a face loc in memorie este gestionata tot de swapper si anume sunt intai analizate procese care asteapta ceva (de exemplu input de la utilizator), daca sunt gasite mai multe este ales cel al carui timp de rezidenta in memorie si al carui prioritate sunt mai mari. Swapperul verifica lista de procese care nu sunt in memorie o data la cateva secunde pentru a decide daca trebuie adus vreunul in memorie. Pentru a preveni suprasolicitarea hardului niciun proces nu poate fi scos din memorie daca nu a stat macar 2 secunde.

2.5. Swapping in cazul Windows 2000

In cazul Windows gestionarea memoriei este mai complicata si swappingul in sine este strans legat de mecanisme complexe de inlocuire a paginilor. Ideea de baza este ca paginile sunt tinute fie in niste seturi de lucru, adica elemente necesare in acel moment pentru procesare, fie in niste liste de pagini cu diferite tipuri de pagini de genul pagini care au fost modificate si nu au fost inca scrise pe disc, pagini care au fost folosite dar nu au fost modificate fata de varianta de pe disc, pagini libere (pagini care contin date, dar nu mai sunt necesare) si pagini cu zero (pagini libere care sunt completate cu 0 peste tot) de asemenea mai exista si o lista cu locatii defecte din RAM astfel incat sa nu se incerce scrierea acolo. Exista mai multe threaduri care muta pagini intre harddisk si memoria RAM, printre care un thread de swapping care din 4 in 4 secunde cauta procese care au fost inactive o anumita perioada de timp si le

muta pe una din listele de pagini modificate, sau de pagini in standby(pagini nemodificate).

2.6.Rezumat

Swappingul este o metoda de gestiune a memoriei ce presupune transferul unor procese din memoria ram pe harddisk pentru a elibera memorie pentru alte procese. Spre deosebire de memoria virtuala care permite stocarea doar unor parti din proces, la swapping se stocheaza intreg procesul.

In urma utilizarii acestei strategii e posibil ca in memorie sa apara zone goale de memorie foarte mici ce practic devin inutilizabile, astfel este util sa se foloseasca un algoritm eficient pentru swapping.

Exista doua metode de a implementa swappingul si anume cu hartile de biti si cu liste inlantuite. Hartile de biti ocupa mai mult spatiu in plus decat listele si algoritmi de cautare a unor spatii libere suficient de mari sunt ceva mai greoi. Pentru gestionarea memoriei cu ajutorul listelor inlantuite exista mai multi algoritmi cu diferite avantaje si dezavantaje.

In cazul Linux-ului swappingul este comandat de un singur nivel dintr-un sistem de gestiune a memoriei pe mai multe nivele, iar algoritmul implementat evita destul de usor thrashing-ul harddiskului (adica utilizarea excesiva cauzata de un transfer nejustificat de des al proceselor intre harddisk si memorie)

In cazul Windows situatia este mai complicata iar swappingul este comandat de mai multe threaduri si in stransa legatura cu algoritmi de inlocuire a paginilor. Astfel sunt folosite diferite liste pentru a selecta paginile ce trebuie transferate spre hard si eventual reimprospatate, threadul care face swapping scaneaza memoria pentru procese ce au fost inactive mai mult de o perioada determinata de timp si trimite paginile din care este alcatuit pe una din listele respective.

Bibliografie:

Andrew S. Tanenbaum - Modern operating systems

Andrew S. Tanenbaum - Operating Systems Design and Implementation

3. Memoria virtuala *Muşat Liana*

Memoria virtuală reprezintă o metodă de organizare a memoriei prin intermediul căreia programatorul "vede" un spațiu virtual de adresare foarte mare care este mapat în memoria fizic disponibilă, fără ca programatorul să "simtă" aceasta.

Aceasta „soluție” a aparut din nevoia unui spațiu mare de memorie, pentru programe prea mari pentru a intra în memoria disponibilă. Soluția folosită în mod uzual a fost împărțirea programului respectiv în bucăți numite overlay. Overlay 0 începea să ruleze primul. Atunci când acesta termină, chemă un alt overlay. Unele sisteme bazate pe overlay erau destul de complexe permițând mai multe overlay-uri în memorie, în același timp. Overlay-urile erau ținute pe disc și erau schimbate între ele în memorie, în mod dinamic, conform cerințelor. Schimbarea efectivă a overlay-urilor și operația de împărțire a programului în bucăți este făcută de sistem, cu ajutorul memoriei virtuale.

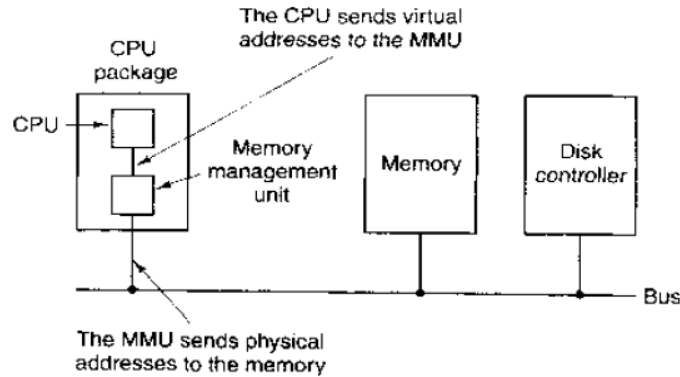
Ideea de bază a memoriei virtuale este aceea că mărimea totală a programului, a datelor, poate depăși mărimea memoriei fizice disponibile pentru acesta. Programul menține părțile ce sunt folosite în memorie iar restul părților sunt ținute pe disc. De exemplu, un program de 16 MB poate rula pe un sistem cu 4MB alegând cu grijă care 4MB trebuie să fie menținuți în memorie în orice moment, iar părțile de care este nevoie se vor schimba între disc și memorie.

Prin mecanismele de memorie virtuală se mărește probabilitatea ca informația ce se dorește a fi accesată de către CPU din spațiul virtual (disc), să se afle în memoria principală, reducându-se astfel în mod semnificativ timpul de acces de la 10-15 ms (timp acces discuri), la 50-80 ns (timp acces DRAM-uri) în tehnologiile actuale.

3.1. Paginare

Majoritatea sistemelor cu memorie virtuală folosesc o tehnică numită paginare. De obicei, spațiul virtual de adresare este împărțit în entități de capacitate fixă (4-64 Ko actualmente), numite pagini. Unitățile corespondente din memoria fizică sunt numite cadre de pagină. Paginile și cadrele de pagină sunt mereu de aceeași mărime. În general, prin mecanismele de memorie virtuală, memoria principală conține paginile cel mai recent accesate de către un program, ea fiind pe post de "cache" între CPU și discul hard.

La orice computer există un set de adrese de memorie ce pot fi produse de programe. Aceste adrese generate de program sunt numite adrese virtuale și formează spațiul adresei virtuale. La computerele ce nu au memorie virtuală, adresa virtuală este pusă direct pe bus-ul memorie și face ca, cuvântul de memorie virtuală să fie scris sau citit. Atunci când se folosește memoria virtuală, adresele virtuale nu ajung direct pe bus-ul de memorie. În schimb, acestea ajung la un Unitate de Management al Memoriei care marchează adresele virtuale în adrese de memorie fizică, după cum este prezentat și în figura.



Transformarea adresei virtuale emisă de către CPU într-o adresă fizică (existentă în spațiul MP) se numește **mapare** sau **translatare**. Memoria virtuală oferă o funcție de **relocare** a programelor (adreselor de program), pentru că adresele virtuale utilizate de un anumit program sunt mapate spre adrese fizice diferite, înainte ca ele să fie folosite pentru accesarea memoriei. Această mapare permite aceluiași program să fie încărcat în memoria principală, modificările de adrese realizându-se automat prin mapare (fără memorie virtuală un program depinde de obicei în execuția sa de adresa de memorie unde este încărcat de către sistemul de operare).

Transferul dintre RAM și disc se realizează în unități de mărimea unei pagini. Atunci când unitatea încearcă să acceseze o adresă 0, folosind instrucțiunea:

```
MOV REG 0
```

adresa virtuală 0 este trimisă la Unitate de Management al Memoriei. Acesta observa că adresa intră în pagina 0 (adică de la 0 la 4095). Memoria nu știe nimic despre Unitate de Management al Memoriei și vede doar o cerere de scriere sau de citire, cerere pe care o onorează. Astfel Unitate de Management al Memoriei marchează toate adresele virtuale între 0 și 4095. Avem posibilitate de a marca 16 pagini virtuale pe oricare 8 cadre de pagina setând harta Unitate de Management al Memoriei în mod corespunzător nu rezolvă problema prin care adresa virtuală este mai mare decât memoria fizică. Din moment ce dispunem doar de 8 cadre numai 8 din paginile virtuale sunt marcate pe adresa fizică. Celelalte nu sunt marcate. Astfel un bit prezent/absent urmărește care pagini sunt fizic prezente în memorie.

Atunci când programul încearcă să folosească o pagină nemarcată Unitate de Management al Memoriei observă că pagina nu este marcată și face ca CPU-ul să execute o operație. Această operație se numește **eroare de pagină**. Sistemul de operare ia un cadru de pagină puțin uzat și scrie conținutul acestuia înapoi pe disc. După aceea ia pagina la care s-a făcut referință și o aplică cadrului de pagină, schimbă marcarea și repornește instrucțiunea prinsă.

Adresa virtuală este împărțită într-un număr de pagini, numărul de pagină este folosit ca un index în tabelul paginii, obținându-se numărul cadrului conform paginii virtuale. De exemplu, adresa virtuală, 8196 (001000000000100 în cod binar), o adresă de 16 biți de intrare care este împărțită într-un număr de pagină de 4 biți, și un rest de 12 biți. Cu 4 biți pentru numărul paginii, avem 16 pagini și cu 12 biți pentru offset putem adresa cei 4096 de biți în cadrul unei singure pagini. Dacă bit-ul prezent/absent este 0, se declanșează o restricționare la sistemul de operare. Dacă bit-ul este 1, numărul cadrului ce se găsește în tabelul din pagină este copiat în comanda superioară de 3 biți a registrului de ieșire, împreună cu offset-ul de 12 biți, acesta din urmă fiind copiat nemodificat din adresa virtuală de intrare. Împreună se formează o adresă fizică de 15 biți. După aceea registrul de ieșire este pus pe bus-ul de memorie sub forma adresei memoriei fizice.

3.2. Tabele pagină

În cazul cel mai simplu, marcarea adreselor virtuale în adresele fizice se realizează conform procedurii descrise mai sus. Adresa virtuală este împărțită într-o număr de pagină virtuală (biți de ordine mare) și într-un offset (biți de ordine inferioară). De exemplu, la o adresă de 16 și o mărime a paginii de 4KB cei 4 biți superiori pot specifica una dintre cele 16 adrese virtuale și cei 12 biți pot specifica offset-ul (0 la 4095) în cadrul adresei selectate. Însă se poate realiza și o împărțire cu 3 sau 5 sau cu un alt număr de biți. Împărțiri diferite implică mărimi diferite de pagină.

Numărul paginii virtuale este folosit drept index în tabelul paginii pentru a găsi intrarea acelei paginii virtuale. Prin intermediul intrării din tabel (dacă aceasta există) se găsește numărul cadrului paginii. Numărul cadrului este atașat la partea superioară a offset-ului, înlocuind numărul paginii virtuale, pentru a forma o adresă fizică ce poate fi trimisă memoriei.

Scopul tabelului din pagină este acela de marca paginile virtuale în cadre. Vorbind din punct de vedere matematic tabelul este o funcție, având numărul paginii virtuale drept argument iar numărul cadrului fizic este rezultatul. Folosind rezultatul acestei funcții câmpul paginii virtuale este o adresă virtuală ce poate fi înlocuită de un cadru de câmp de pagină, astfel formându-se o adresă în memoria fizică.

trebuie să se țină cont de două probleme importante:

1. tabelul din pagină poate fi foarte mare
2. marcarea trebuie să se realizeze rapid.

Primul punct rezultă din faptul că, calculatoarele moderne folosesc adrese virtuale de cel puțin 32 de biți. Luând în considerare o pagină cu mărimea de 4 Kb, un spațiu de adresă de 32 biți are 1 milion de pagini, și un spațiu de adresă de 64 de biți are mai mult decât ați putea calcula. Cu un milion de pagini în spațiul pentru adresa virtuală tabelul paginii trebuie să aibă 1 milion de intrări. Și țineți cont de faptul că fiecare proces are nevoie de propriul tabel de pagină (deoarece are propriul spațiu pentru adresa virtuală).

Al doilea punct este o consecință a faptului că marcarea de la virtual la fizic trebuie să fie făcută la fiecare referință a memoriei. O instrucțiune tipică conține un cuvânt de instrucțiune și deseori un operand. În consecință este nevoie să faceți 1, 2 sau chiar mai multe referințe la tabelul paginii, pentru fiecare instrucțiune în parte. Dacă o instrucțiune durează, să spunem, 4 ns, căutarea în tabelul paginii trebuie să se realizeze în mai puțin de 1 ns pentru a evita stagnarea.

Cel mai simplu design (cel puțin conceptual vorbind) este de a avea un singur tabel de pagină ce este alcătuit dintr-o gamă de registre de componente rapide, cu câte o intrare pentru fiecare pagină virtuală, intrare indexată de numărul paginii virtuale. Atunci când un proces este inițializat, sistemul de operare încarcă regiștrii cu tabelul paginii procesului, dintr-o copie ce este ținută în memoria principală. În timpul execuție procesului nu mai sunt necesare referințe pentru tabel. Avantajele acestei metode sunt acelea că este o metodă directă și nu necesită accesul la memorie în timpul procesului de marcare. Dezavantajul constă în faptul că poate fi scump (în cazul în care tabelul este mare). Necesitatea de a încărca întreg conținutul tabelului duce la diminuarea performanțelor.

La cealaltă extremitate se află posibilitatea conform căreia tabelul din pagină este în întregime în memoria principală. În această situație hardware-ul are nevoie de un singur registru ce indică începutul tabelului paginii. Acest design permite schimbarea identificării memoriei la un comutator de context, reîncărcând un singur registru. Desigur, acesta are dezavantajul de a necesita una sau mai multe referințe de memorie pentru a citi intrările din tabel, în timpul execuției fiecărei instrucțiuni. Din acest motiv, această abordare este rar folosită în forma ei cea mai pură, dar mai jos vom studia câteva variații care au performanțe mult mai bune.

3.2.1. Tabele de pagină cu nivel multiplu

Pentru a evita problema ce ține de stocarea continuă a unor tabele imense, multe calculatoare folosesc un tabel cu nivel multiplu. Secretul metodei de tabel de pagină cu nivel multiplu este de a evita ținerea tuturor tabelurilor în memorie în orice moment. În special, acelea de care nu este nevoie nu ar trebui să fie ținute în memorie.

Pentru un tabel cu doua nivele, cu 1024 de intrări. Fiecare din aceste 1024 de intrări reprezintă 4M deoarece întreg spațiul adresei virtuale de 4 Gb este împărțit în bucăți de 1024 de bytes. Intrarea localizată cu ajutorul indexării în partea de sus a tabelului produce adresa sau numărul cadrului paginii pentru un tabel de pagină de nivelul doi. Intrarea 0 din partea de sus a tabelului de pagini indică tabelul de pagină al programului următor, intrarea 1 indica tabelul de pagină ce conține datele, și intrarea 1023 indică tabelul de pagină pentru majoritatea informației. Celelalte intrări nu sunt folosite.

3.2.2. Structura unei intrări în tabel

Aspectul unei intrări în tabel este dependent de mașină, dar tipul de informații prezent este în mare același de la mașină la mașină.

Cel mai important câmp este *Page frame number*. Scopul mapării este de a localiza această valoare. Lângă acesta avem *Present/absent bit*. Dacă acest bit este 1 intrarea este validă și poate fi folosită. Dacă intrarea este 0 pagina virtuală corespondentă intrării nu este în memorie. Accesarea unei intrări din tabel ce are acest bit 0 duce la apariția unei erori de pagină. Biții *Protection* spun ce fel de acces este permis. În forma cea mai simplă câmpul conține 1 bit, cu 0 pentru citire/scriere și 1 bit doar pentru citire. O aranjare mai sofisticată constă în a avea 3 biți, câte unul pentru permiterea citirii, scrierii și executării paginii.

Biții *Modified* și *Referenced* țin cont de folosirea paginii. Când se scrie pe o pagină, hardware-ul declanșează automat bitul *Modified*. Acest bit este de valoare atunci când sistemul de operare decide să revindică un cadru. Dacă pagina din acesta a fost modificată (de exemplu este „murdară”) aceasta trebuie să fie scrisă pe disc. Dacă aceasta nu a fost modificată (de exemplu este „curată”) poate pur și simplu să fie abandonată din moment ce copia de pe disc este încă valabilă. Bit-ul este uneori numit **bit-ul murdar** din moment ce acesta reflectă starea în care se află pagina.

Bit-ul *Referenced* este setat atunci când o pagină este referită, fie pentru citire sau scriere. Scopul acestuia este de a ajuta sistemul de operare în alegerea unei pagini ce va fi evacuată atunci când apare o eroare de pagină. Paginile care nu sunt folosite sunt candidați mai buni decât paginile care sunt, și acest bit joacă un rol important în câțiva dintre algoritmii de înlocuire a paginii, algoritmi ce vor fi studiați mai târziu în această lucrare.

Ultimul bit permite dezactivarea cache-ului pentru această pagină. Această caracteristică este importantă pentru paginile care se mapează pe regiștri în loc de memorie. Dacă sistemul de operare are un timp restrâns de acționare așteptând răspunsul unui dispozitiv I/O la o comandă ce tocmai a fost dată, atunci este esențial ca hardware-ul să ia cuvântul de la dispozitiv și să nu folosească copia care se află în cache. Cu ajutorul acestui bit cache-ul poate fi oprit. Sistemele care un spațiu I/O separat și nu folosesc I/O mapate pe memorie nu au nevoie de acest bit.

Adresa discului ce este folosită pentru a reține pagina atunci când aceasta nu este în memorie nu face parte din tabel. Motivul este simplu. Tabelul conține numai acea informație de care hardware-ul are nevoie pentru a transforma o adresă virtuală într-o adresă fizică. Informația de care sistemul de operare nu are nevoie pentru a manevra erorile de pagină este

ținută în tabele de program, în interiorul sistemului de operare. Hardware-ul nu are nevoie de ea.

3.3.TLB-uri – Translation Lookaside Buffers

În majoritatea schemelor de paginare tabelele sunt ținute în memorie din cauza mărimii lor. Din punct de vedere teoretic, acest tip de design are un impact enorm asupra performanțelor. Considerăm, de exemplu, o instrucțiune ce copiază un registru la altul. În lipsa paginării această instrucțiune face o singură referință la memorie pentru a lua instrucțiunea. Cu paginarea, va fi nevoie de referințe de memorie suplimentare pentru a accesa tabelul. Din moment ce viteza de execuție este limitată de rata cu care CPU-ul poate lua instrucțiuni și date din memorie, faptul că este necesar să se facă referințe la tabel pentru două pagini duce la reducerea performanței cu 2/3. În aceste condiții nimeni nu ar folosi acest sistem.

Soluția are la bază observația conform căreia majoritatea programelor au tendința de a face un număr mare de referințe la un număr mic de pagini și nu invers. Astfel, numai o fracțiune mică din intrările din tabel sunt citite în mod constant; celelalte sunt rar folosite.

Soluția cu care s-a venit a fost echiparea calculatoarelor cu un dispozitiv hardware mic pentru maparea adreselor virtuale în adrese fizice fără a trece prin tabel. Dispozitivul se numește **TLB** sau uneori o **memorie asociativă**, este de obicei în interiorul MMU și este alcătuit dintr-un număr mic de intrări, opt în acest exemplu, dar rareori sunt mai mult de 64. Fiecare intrare conține informații cu privire la o pagină inclusiv numărul paginii virtuale, un bit ce este setat când pagina este modificată, codul de protecție (permisiile de citire/scriere/executare), și cadrul de paginii fizice în care pagina este situată. Aceste câmpuri au o corespondență unu la unu cu câmpurile din tabelul de pagină. Un alt bit indică dacă o intrare este validă sau nu. Atunci când o adresă virtuală este prezentată la MMU pentru translație, hardware-ul verifică mai întâi să vadă dacă numărul paginii virtuale este presetat în TLB comparându-l simultan cu toate intrările (în paralel). Dacă se găsește o pereche validă și accesul nu încalcă biții de protecție cadrul paginii este luat din TLB, fără a mai merge la tabel. Dacă numărul paginii virtual este prezent în TLB dar instrucțiunea încearcă să scrie pe o pagină ce nu poate fi scrisă, o eroare de protecție este generată, la fel cum s-ar fi întâmplat și din tabel.

Intrările TLB sunt încărcate în mod explicit de către sistemul de operare. Când un TLB dă o eroare, în loc ca MMU să se ducă la tabele pentru a găsi și extrage referința de pagină necesară acesta generează o eroare TLB și pune problema în seama sistemului de operare. Sistemul trebuie să găsească pagina, să extragă o intrare din TLB, să introducă una nouă și să repornească instrucțiunea ce a dat eroare. Dacă TLB este destul de mare (să spunem 64 de intrări) pentru a reduce rata de ratare, managementul programului corespondent TLB se dovedește a fi eficient. Principalul câștig aici este un MMU mult mai simplu, care eliberează suprafață semnificativă pe CPU, pentru cache-uri și alte caracteristici ce ajută la îmbunătățirea performanței. Mai multe strategii au fost dezvoltate pentru a îmbunătăți performanța mașinilor ce realizează managementul TLB prin intermediul programului. O abordare acoperă atât reducerea numărului de ratări ale TLB cât și reducerea costului unei ratări TLB, atunci când aceasta apare.

Modul normal de procesare al unei ratări TLB fie că este în program sau în hardware, este să se meargă la tabel și să se realizeze operațiile de indexare pentru a localiza pagina la care se face referință.

3.4.Table de pagină inversate

Tabelele tradiționale din tipul celor descrise până acum au nevoie de o intrare pentru fiecare pagină virtuală din moment ce sunt indexate conform numărului paginii virtuale.

Tabelele de pagină inversate economisesc cantități mari de spațiu, urmărind ce proces este localizat în cadrul paginii, cel puțin atunci când spațiul adresei virtual este mult mai mare decât memoria fizică, acestea au o parte negativă: translația virtual-fizic devine mult mai dificilă. Când procesul nu face referire la pagina virtuală, hardware-ul nu va mai putea găsi pagina fizică folosind ca index în tabel.

TLB poate reține toate paginile care sunt folosite des, translația poate avea loc la fel de rapid ca în cazul tabelelor normale. Însă, la o ratare TLB, tabelul de pagină inversat trebuie să fie căutat în program. Un mod fezabil de a realiza această căutare este de a avea un tabel mărunțit, în adresa virtuală. Dacă tabelul mărunțit are la fel de multe sloturi ca și numărul de pagini fizice al mașinii, lanțul mediu va fi de o intrare, fapt ce mărește viteza de mapare. Odată ce numărul cadrului paginii a fost găsit, noua pereche (virtual, fizic) este introdusă în TLB.

Tabelele inversate sunt în prezent folosite pe câteva stații de lucru IBM și HP și vor deveni din ce în ce mai des întâlnite pe măsură ce mașinile pe 64 de biți sunt mai des folosite.

Bibliografie

- 1) Patterson D., Hennessy J.- Computer Organization and Design:
The Hardware - Software Interface
- 2) Tanenbaum – Modern operating systems
- 3) Bensoussan, A. & Clingen, C. T. (May 1972),
[The Multics Virtual Memory: Concepts and Design](#)

4. Algoritmi de înlocuire a paginilor *Datcu Octaviana*

4.1. Introducere

Atunci când are loc un *defect de pagină** este necesar ca sistemul de operare să elibereze un cadru de pagină prin expulzarea unei pagini care are cea mai mică probabilitate să fie solicitată în lucru, în viitorul apropiat. Algoritmii de înlocuire a paginilor sunt cei care fac posibilă o selecție cât mai bună a paginii “victimă”.

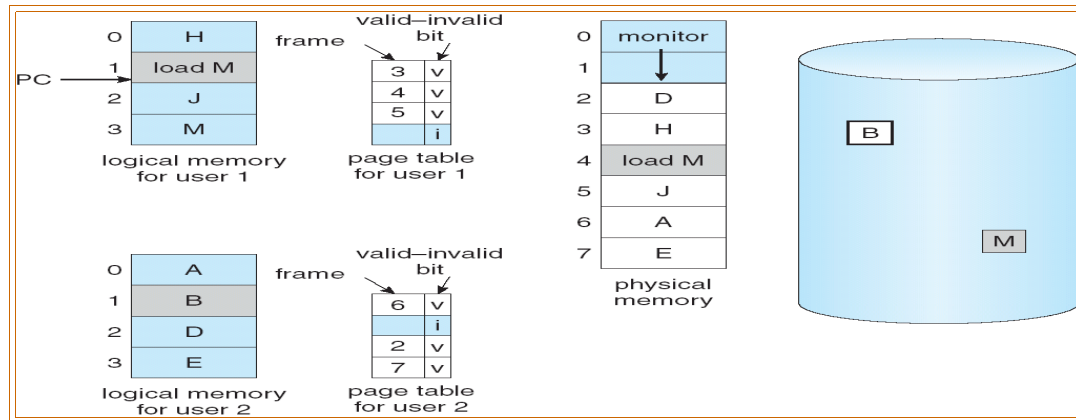


Figura 1 Este necesară înlocuirea paginilor

Aceasta are loc într-un sistem de operare ce utilizează paginarea pentru memoria virtuală.

Paginarea survine atunci când o pagină “liberă” nu poate fi folosită pentru ca alocarea să fie făcută cu succes, fie pentru că nu mai există pagini libere, fie din cauză că numărul acestora este mai mic decât un anumit prag.

Structurile de date și mecanismele utilizate pentru alocarea și dealocarea paginilor sunt critice pentru menținerea eficienței subsistemului format de memoria virtuală.

Performanța acestuia este mult mai bună dacă este aleasă pentru înlocuire o pagină care nu este intens folosită.

În continuare vor fi prezentați algoritmii cel mai adesea utilizați pentru a realiza înlocuirea paginilor.

Au fost efectuate studii referitoare la performanțele acestor algoritmi și au fost propuse variante de îmbunătățire a lor. Vom prezenta, deci, și câteva dintre aceste variante.

*defect de pagină** = procesorul dorește să acceseze o adresă virtuală care nu se află în memoria principală.

4.2. Strategii de înlocuire a paginilor – algoritmi

4.2.1. Modelul înlocuirii optime a paginilor (OPRA)

Modelul înlocuirii optime a paginilor a fost descris de către L.A. Belady.

Principiul acestui algoritm este: “înlocuiește pagina care nu va fi folosită perioada cea mai îndelungată în viitor”. Acest deziderat îi conferă statutul de cel mai bun algoritm de înlocuire a paginilor (din punct de vedere teoretic). Cu implementarea însă lucrurile stau diferit, căci sistemul de operare nu poate privi în viitor, deci este imposibil de realizat practic.

O altă importantă calitate a acestui algoritm este aceea că el oferă suportul de evaluare a celorlalți algoritmi de înlocuire a paginilor.

Iată, pe scurt, descrierea OPRA :

1. are loc o eroare de pagină ;
2. un set de pagini există în memorie;
3. etichetează fiecare pagină cu numărul de instrucțiuni ce vor fi executate înainte ca această pagină să fie folosită din nou, în viitor;
4. înlocuiește pagina ce are eticheta de valoare cea mai mare.

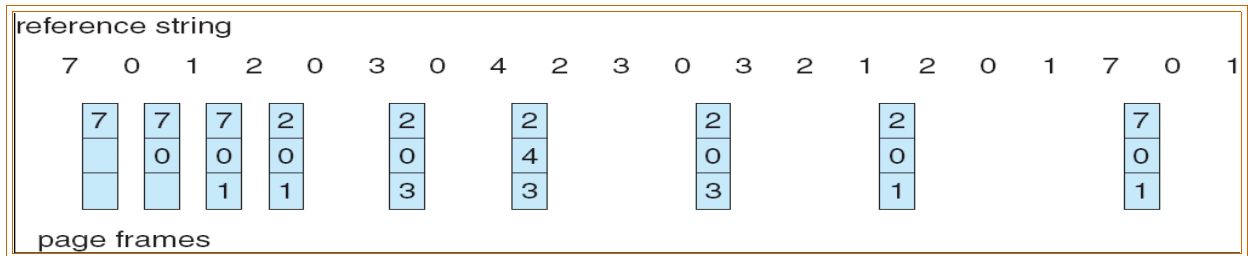


Figura 2 Exemplu

Surse bibliografice: [1],[2],[4],[5],[6],[8],[9],[10],[11],[12],[14]

4.2.2. Not Recently Used (Neutilizată recent)

NRU utilizează biți de stare, R și M, asociați fiecărei pagini, cu următoarea semnificație:

R: pagină la care s-a făcut referire (citită sau scrisă);

M: pagină modificată (scrisă).

În funcție de valoarea acestor biți, 0 sau 1, algoritmul face distincția între patru clase posibile:

		R	M
Clasă 1	Pagină nereferită, nemodificată	0	0
Clasă 2	Pagină nereferită, modificată	0	1
Clasă 3	Pagină referită, nemodificată	1	0
Clasă 4	Pagină referită, modificată	1	1

Figura 3 Tabel ce face distincția între clase (de pagini)

La prima vedere, paginile de clasă 1 nu pot fi întâlnite (modificată, dar totuși nereferită ?!). Este însă posibilă existența lor, atunci când o pagină de clasă 3 (referită, dar nemodificată) are, datorită unei întreruperi (clock interrupt), bitul R modificat la 0. Întreruperile

nu modifică și bitul M, deoarece informația conținută de acesta este utilă, necesară chiar, pentru a ști dacă pagina respectivă trebuie, sau nu, să fie rescrisă pe disk. Deci, iată o pagină (0,0), de clasă 1.

Algoritmul înlocuiește, în mod aleator, o pagină din cea mai joasă (ca număr) clasă nevidă; este, deci, înlocuită o pagină modificată, care nu a fost referită în cel puțin un tact de ceas (20 ms, tipic), mai curând decât o pagină nemodificată utilizată intens.

NRU este accesibil ca nivel de înțelegere, implementare și are o performanță satisfăcătoare.

Surse bibliografice: [1],[4],[10],[12],[14]

4.2.3. First In, First Out (primul intrat, primul ieșit)

Menține o listă înlănțuită a tuturor paginilor, în ordinea în care acestea apar în memorie. Principiul First In First Out implică faptul că pagina de la începutul listei va fi cea înlocuită.

FIFO are marele dezavantaj de a permite îndepărtarea unei pagini necesare. Din această cauză, el este rar utilizat în formă "pură"

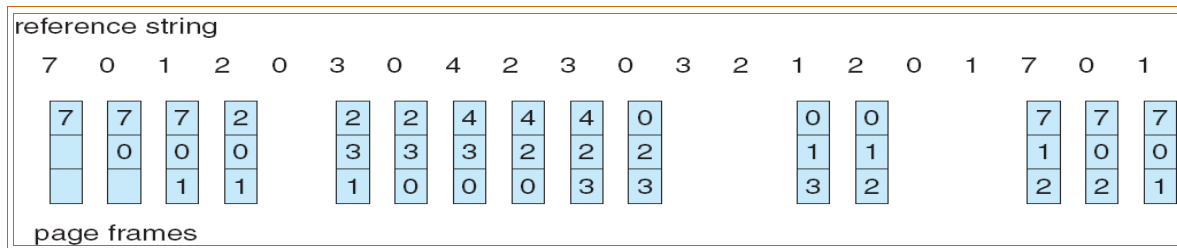


Figura 4 Exemplu

Anomalia lui Belady:

Anomalia lui Belady demonstrează (1969) că este posibil să existe mai multe erori de pagină atunci când cadrele de pagină cresc în timpul utilizării algoritmului FIFO.

Procesorul poate încărca un număr limitat de pagini la un moment dat. El pretinde un cadru pentru fiecare pagină ce o poate încărca.

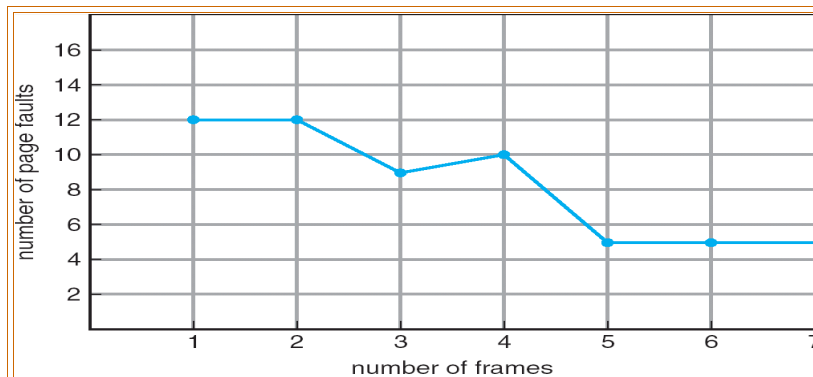


Figura 5 Ilustrarea anomaliai lui Belady

Surse bibliografice: [1],[2],[4],[5],[6],[7],[8],[10],[11],[12],[14]

4.2.4. Least Recently Used (Cel mai recent utilizat)

Paginile intens folosite în ultimele câteva instrucțiuni vor fi probabil intens utilizate din nou în următoarele câteva instrucțiuni. Paginile care nu au fost folosite timp îndelungat, au probabilitate mare de a rămâne neutilizate pentru o perioadă mare, în viitor. Această idee se apropie promițător de aceea a algoritmului lui Belady.

LRU folosește, deci, informația referinței, pentru a lua o decizie în privința înlocuirii unei pagini într-o mai bună “cunoștință de cauză”. Urmând ideea de predicție a viitorului sugerată de algoritmul optimal, “ghicește” viitorul pe baza experienței anterioare.

Implementarea LRU nu este ieftină, pentru că pentru aceasta este necesar a menține o listă înlănțuită a tuturor paginilor din memorie, cu cel mai recent folosită pagină în prim plan, cea care a fost utilizată cel mai puțin recent, rămânând pe o ultimă poziție. Lista trebuie să fie reactualizată la fiecare referire a memorie. De aceea, se consumă timp mult prea mare (se caută pagina în listă, este găsită și ștersă din memorie, și rescrisă la începutul listei).

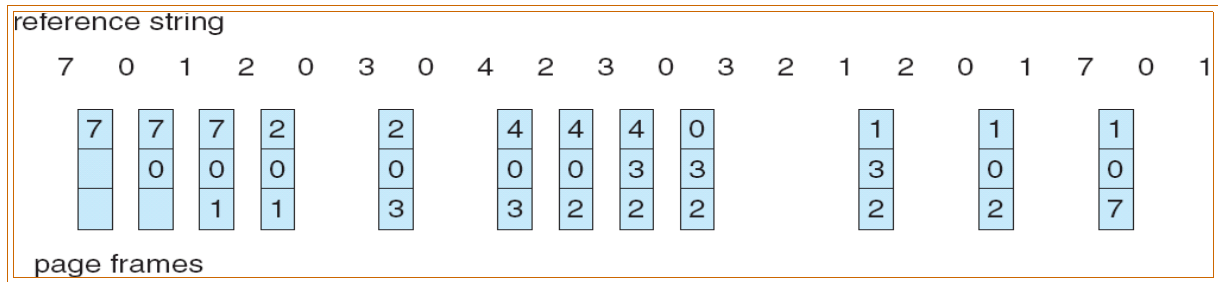


Figura 6 Exemplu

Este necesară, deci, o aproximare a acestui algoritm, pentru a-i conferi o îmbunătățire.

Surse bibliografice: [1], [2], [3],[4],[5],[6],[7],[8],[10],[11],[12],[14]

4.2.5. Implementarea cu stivă a LRU

Este păstrată o stivă a numerelor paginilor. Pagina referită este mutată în capul listei. Dezavantajul acestei operații constă în utilizarea a 6 pointeri pentru a fi modificată. Avantajul metodei este că nu se face căutare pentru înlocuire.

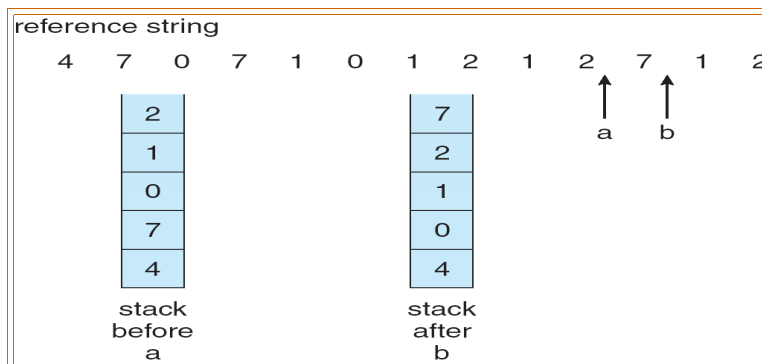


Figura 7 Exemplu

4.2.6. Algoritmi de aproximare LRU

4.2.6.1. Second Chance (A Doua Șansă)

Bitul de referire:

1. asociem fiecărei pagini un bit, R, inițializat cu 0;
2. când pagina este referită, bitul R devine 1;
3. înlocuiește pe cea cu R=0, dacă există.

A doua șansa:

1. este nevoie de bit de referire;
2. este o înlocuire cu ordine;
3. dacă pagina de înlocuit (în ordinea ceasului) are R = 1 atunci:
 - 1) setează R = 0;
 - 2) păstrează pagina în memorie;
 - 3) înlocuiește următoarea pagină (în ordinea acelor de ceasornic);
 - 4) se repetă algoritmul de la 3) pentru fiecare pagină.

Second Chance (A Doua Șansă) este o formă modificată a FIFO, îmbunătățită. În aceeași manieră ca și FIFO, punctează la capătul cozii. Dar, față de FIFO, verifică, în plus, dacă bitul său R este setat. În cazul în care nu este setat, pagina este "scoasă din joc". Pentru R setat, acesta este resetat, pagina este inserată la sfârșitul cozii, și procedul este repetat.

A Doua Șansă dă fiecărei pagini o "a doua oportunitate" - o pagină care a fost referită este cel mai probabil în folosință, și nu ar trebui îndepărtată, spre deosebire de o nouă pagină care nu a fost referită.

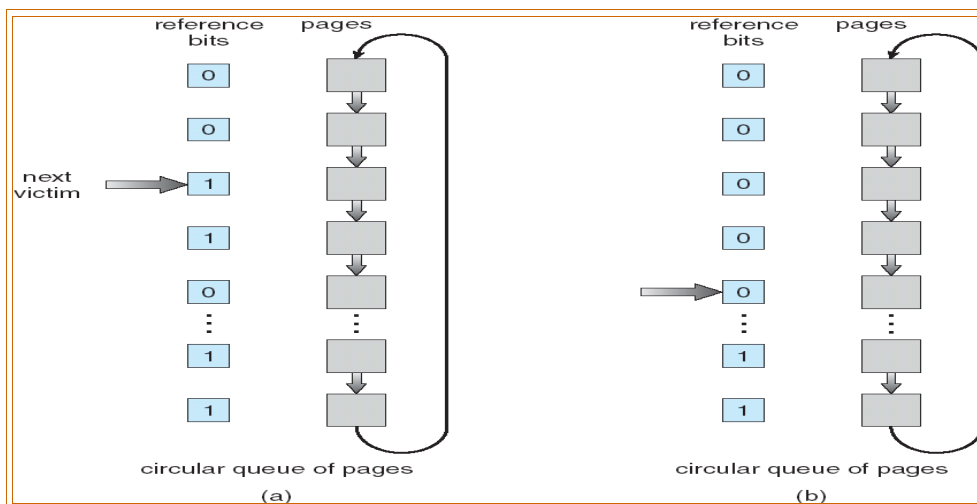


Figura 8 Exemplu

Surse bibliografice: [1],[2],[3],[5],[6]

4.2.6.2. Clock (Ceasul)

Algoritmul A Doua Șansă este rezonabil, dar ineficient întrucât rotește paginile în listă în mod constant. O abordare mai bună constă în păstrarea tuturor cadrelor într-o listă circulară în forma unui ceas.

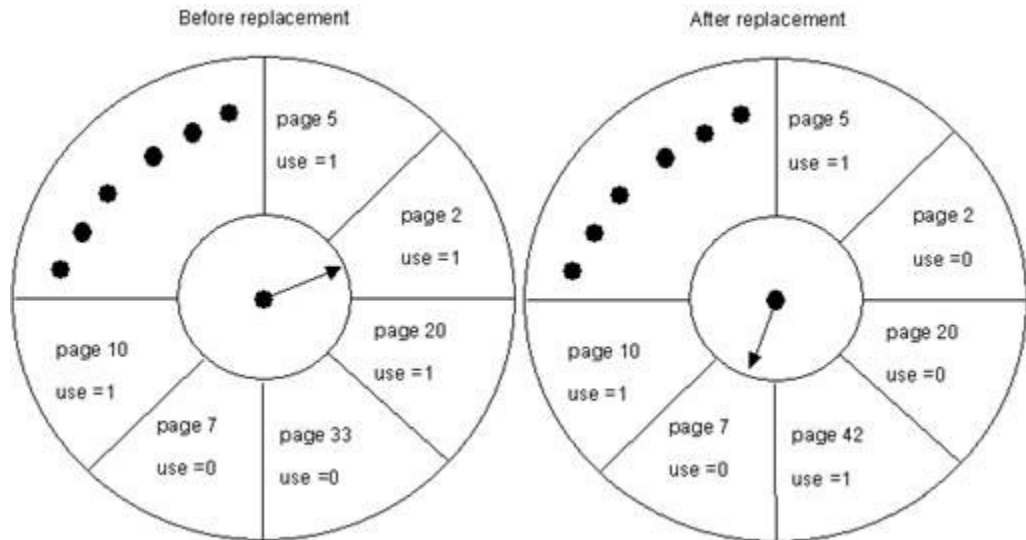


Figura 9 Exemplu

Pagina utilizată în urmă cu cel mai mult timp (prima din coadă) este analizată. Dacă bitul ei $R=0$, pagina este înlocuită, o nouă pagină este inserată în ceas în locul său. Dacă $R=1$, el este resetat la 0 și se trece la următoarea pagină în ordinea impusă de FIFO (pagina curentă este păstrată în memorie).

Algoritmul diferă de cel numit A doua Șansă doar prin implementare.

Surse bibliografice: [3],[10],[11],[14]

4.2.6.3. Not frequently used (Neutilizată frecvent)

NFU generează mai puține defecte de pagină decât LRU atunci când tabelul de pagini conține valori nule de pointeri.

Presupune existența unui numărator, fiecare pagină având un contor propriu, inițializat cu 0. La fiecare interval de ceas, pentru toate paginile care au fost referite în acel interval număratorul le va fi incrementat cu 1, aceste numărătoare reținând frecvența cu care o anumită pagină a fost folosită. Astfel, pagina cu număratorul având valoarea ce mai mică poate fi îndepărtată atunci când este necesar.

Dezavantajul acestei strategii este că ține cont de frecvența de utilizare, fără a reține și timpul de utilizare. De aceea, o pagină care a fost intens utilizată în timpul primului pas, va fi favorizată față de una care este utilizată mai puțin, în cel de-al doilea pas, deoarece are o frecvență mai mare. Acest lucru conduce la o performanță scăzută.

Surse bibliografice: [2],[10],[11],[14]

Există un algoritm similar, dar mai bun:

4.2.6.4. Aging (Îmbătrânirea)

Modificare care îi este adusă lui NFU, transformându-l în "Aging" este aceea că, noul algoritm ține cont nu numai de frecvența de utilizare a paginilor, ci și de timpul cât durează această utilizare.

Numărătorul nu mai este pur și simplu incrementat cu 1, ci, anterior acesta este deplasat la dreapta (deci, împărțit la 2), la stânga a ceea ce s-a obținut fiind adăugat bitul de referință.

Paginile referite cel mai recent, chiar dacă mai puțin referite, vor avea o prioritate mai mare decât paginile referite mai frecvent, dar în trecut.

Când este necesar să fie îndepărtată o pagină, atunci va fi aleasă aceea care are cea mai mică valoare în numărător.

Iată, deci, că, "Aging" poate oferi o performanță apropiată celei optimale, contra unui cost mediu.

Surse bibliografice: [2]

4.2.7. Random

Algoritmul de înlocuire "Random" alege în mod aleator, așa cum sugerează și denumirea sa, o pagină pentru a fi înlocuită. Aceasta elimină inconvenientul costului necesar memorării referirii la pagini.

Este mai bun decât FIFO. Pentru referiri ale memoriei în buclă este mai bun chiar și decât LRU, deși acesta își găsește o mai bună aplicabilitate în practică.

Surse bibliografice: [4],[7]

4.3. Impactul dimensiunii paginii asupra performanțelor memoriei virtuale

Dimensiunea paginilor este de dorit a fi aleasă astfel încât să se reducă fragmentarea internă.

$$\text{overhead} = \left(\frac{s \cdot e}{p} \right) + \left(\frac{p}{2} \right)$$

Spațiul tabelului Fragmentarea internă

s = dimensiunea medie a unui proces [bytes];

p = dimensiunea paginii [bytes];

e = dimensiunea intrării de pagină;

Dimensiunea optimă se determină prin calculul derivatei expresiei în funcție de p și anularea acesteia.

$$-(s \cdot e)/(p \cdot p) + 1/2 = 0 \rightarrow p = (2 \cdot s \cdot e)^{1/2}$$

Soluționarea unei erori de pagină:

1. sistemul de operare consultă un alt tabel:
 - dacă adresa nu este validă: sfârșitul procesului;
 - dacă adresa este validă: tabelul indică adresa paginii pe disc.
2. se alocă un frame liber;

3. se încarcă pagina referită în memorie;
 4. se reactualizează tabelul paginilor (bitul de invalidare = 1);
 5. se pornește reexecuția instrucțiunii
- Această modalitate se numește a face *paginarea la cerere*.

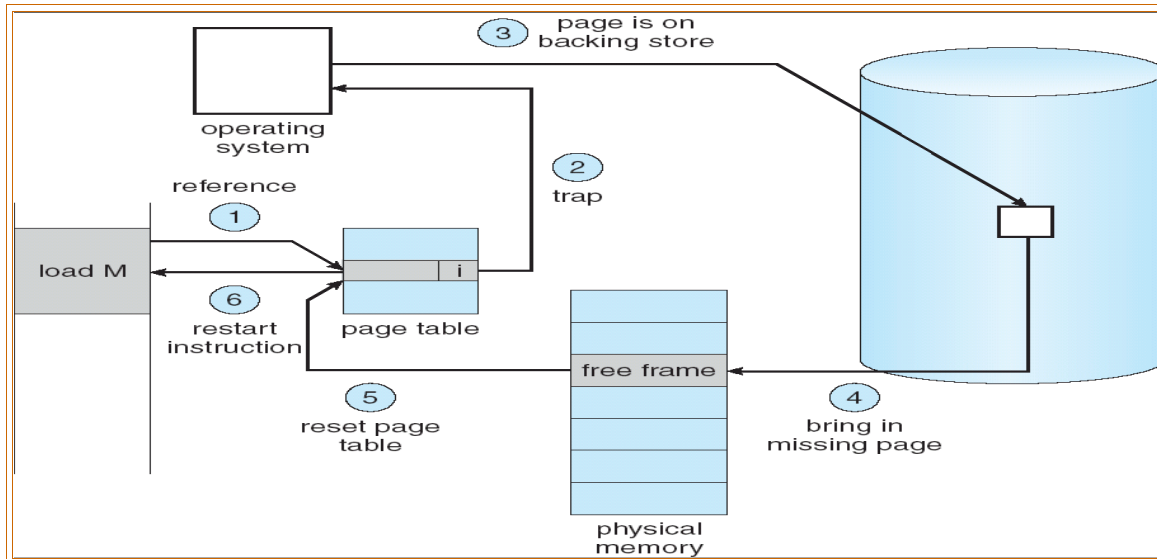


Figura 9 Soluționarea unei erori de pagină

Surse bibliografice: [1],[3],[6]

4.4.Comportamentul programului în cazul paginării

Sistemul de operare determină dimensiunea datelor și pe cea a instrucțiunilor. Alocă și inițializează tabelul de pagini. Apoi alocă aria de swap pe disc pentru a stoca tabelul de pagini când procesul este swapped out. Este memorat tabelul de pagini și aria de swap în tabelul procesului. Posibila este și preîncărcarea de pagini.

La începerea execuției programului, este resetat MMU (Memory Management Unit), încărcat TLB-ul și copiat tabelul de pagini sau adresele de start/sfârșit în registrele hardware.

Încheierea procesului presupune eliberarea tabelului de pagini, a cadrelor de pagini și spațiului de swap de pe disc.

În cazul *defectului de pagină*, S.O. determină care adresă virtuală a cauzat eroarea de pagină și află ce tabel este necesar pentru a o localiza pe disc. Apoi expulzează o pagină pentru a face loc în memoria principală (dacă e nevoie). Întoarce numărătorul de program pentru a indica la instrucțiunea eronată astfel încât aceasta să poată fi executată din nou (de această dată, cu pagina necesară aflată în memorie).

Aici își găsesc utilitatea algoritmiile prezentate la punctul 4.2.

Surse bibliografice: [3],[6]

4.5. Working Set (Setul de lucru)

Localitatea programelor este o proprietate observată experimental. În fiecare interval de timp, un process face referire doar la un subsansamblu al paginilor sale.

Formalizare:

- procesul ocupă n pagini : 1, ..., n-1
- execuția procesului generează o secvență de referințe la aceste pagini
- $w_{\Delta(t)}(t) =$ paginile referite în intervalul $[t-\Delta t, t]$

Principiul localității : $w_{\Delta(t)}(t) = w_{\Delta(t-\Delta t)}(t+\Delta t)$

page reference table

... 2 6 1 5 7 7 7 5 1 6 2 3 4 1 2 3 4 4 4 3 4 3 4 4 4 1 3 2 3 4 4 4 3 4 4 4 ...

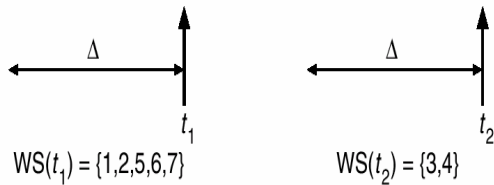


Figura 10 Setul de lucru

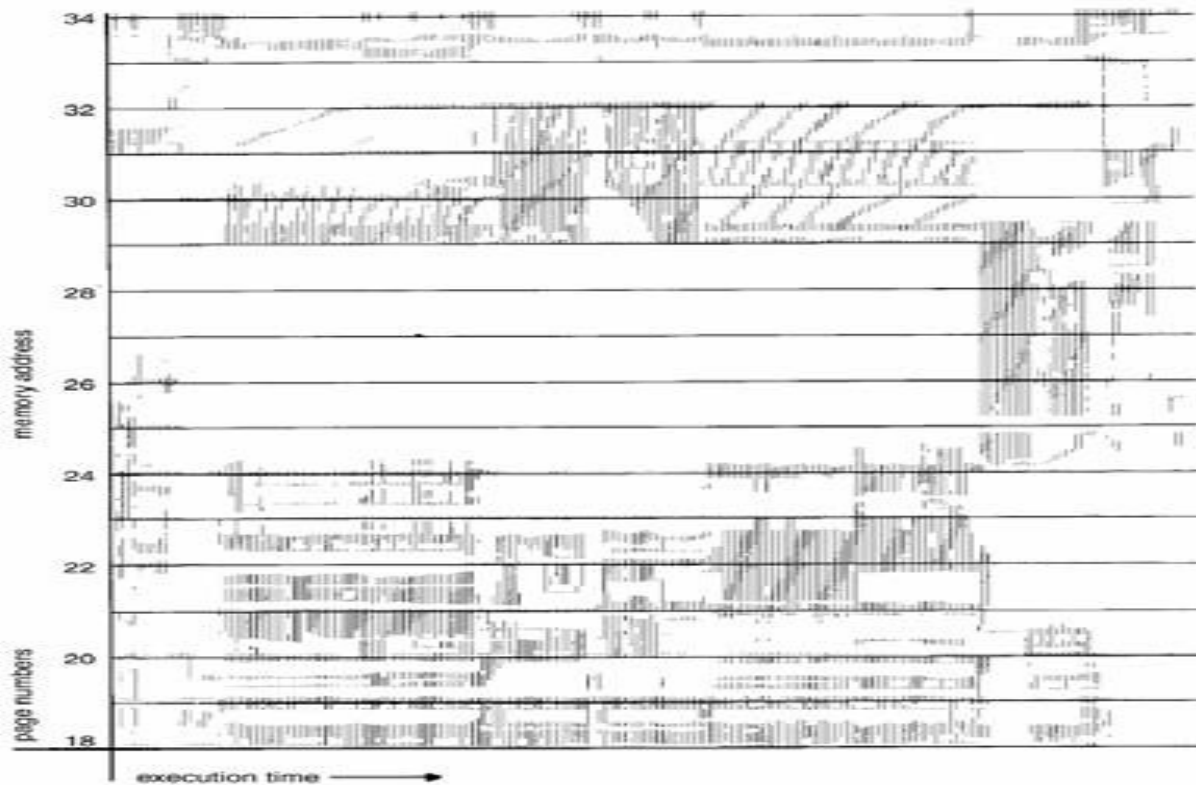


Figura 11 Localitatea

1. *Localitatea temporală* : o adresă care tocmai a fost referită are mari șanse să fie referită în imediata apropiere temporală.
Exemplu : bucle, variabile locale ale unei proceduri, proceduri etc.
2. *Localitatea spațială*: dacă o adresă a a fost referită, există șanse mari ca o adresă vecină lui a (adică în aceeași pagină) să fie referită în curând.
Exemplu: execuție secvențială a codului, parcurgerea tablourilor etc.

Deși costul paginării este ridicat, dacă frecvența producerii acesteia este suficient de mică, paginarea este acceptabilă.

De regulă, procesele implică ambele tipuri de localitate în timpul execuției lor, dând paginării practicabilitate.

Din perspectiva sistemului de operare, când memoria este plină, pagini sunt expulzate, și încărcate de pe disc atunci când sunt referite din nou. Referirea la pagini expulzate cauzează erori în TLB.

Cererea de pagini are loc și atunci când un proces este lansat pentru prima dată; procesul are un tabel al paginilor nou, cu toți biții valizi 0 și nici o pagină în memorie. Când își începe execuția, instrucțiunile produc defecte în cod și pagini de date. Erorile încetează atunci când întregul cod și paginile de date necesare se află în memorie. Înărcarea este necesară doar pentru codul și paginile active ale procesului.

Spații fixe / spații variabile

Într-un sistem cu multiprogramare este nevoie să alocăm memorie proceselor concurente. Dar cum să determinăm ce cantitate de memorie să dăm fiecărui proces?

Algoritmi de spațiu fix:

- fiecărui proces îi este dat un anumit număr de pagini limitat ca să le poată folosi;
- când procesul atinge această limită începe să înlocuiască din propriile pagini;
- înlocuirea este locală;
- unele procese pot să funcționeze foarte bine, în detrimentul altora.

Algoritmi de spațiu variabil:

- setul de pagini al procesului crește și se reduce dinamic;
- înlocuirea este globală;
- unul dintre procese poate să deterioreze funcționarea celorlalte.

Setul de lucru al unui proces este utilizat pentru a modela localitatea dinamică a folosirii memoriei în cazul său. Reprezintă numărul adecvat de pagini pentru ca un proces să funcționeze eficient, deci, acel număr de pagini pe care procesul le folosește active, la un moment dat, care trebuie să fie încărcat în memorie la acel moment pentru ca procesul să fie executat.

Colocvial, WS este setul de pagini pe care procesul le folosește în chiar acest moment. Mai formal, WS este setul tuturor paginilor pe care procesul le-a referit în ultimele T (Δt) secunde.

Modul în care sistemul de operare alege T -ul este următorul: la un defect de pagină corespund $10 \text{ ms} = 2$ milioane de instrucțiuni. Dar T trebuie să fie mult mai mare ca 2 milioane de instrucțiuni. Dacă T este prea mic, sunt îndepărtate pagini aflate încă în uz, la cealaltă extremă (T prea mare) aflându-se utilizarea ineficientă a memoriei (reducerea inutilă a nivelului de multiprogramare).

Paginile setului de lucru sunt înlănțuite (listă circulară). Fiecare pagină a WS are un bit de referință setat la 1 pe durata tuturor referințelor la pagina respectivă. Căutarea paginii care va fi adusă presupune parcurgerea listei, pornind de la poziția curentă. Dacă pagina are $R=1$, bitul este resetat la 0. Dacă $R=0$, pagina este adusă în memorie.

Numărul de pagini referite în intervalul $[t, t-T]$ este dimensiunea WS-ului. Setul de lucru se schimbă cu localitatea programului. Pe durata unei localități slabe sunt referite mai multe pagini, dimensiunea WS-ului fiind mai mare.

Problemele care se pun sunt: 1) cum determinăm T ; 2) cum știm care este momentul în care se schimbă WS-ul. Din această cauză WS nu este folosit în practică ca un algoritm de înlocuire a paginilor. Este folosit ca o abstractizare.

Surse bibliografice: [1],[2],[4],[5],[6],[7],[8]

4.6. Fenomenul “thrashing”

Thrashing-ul are loc atunci când cea mai mare parte a timpului sistemul de operare se ocupă de paginarea datelor pe/de pe disc, în loc să progreseze dând prioritate lucrului util. Sistemul este astfel suprasolicitat.

Nu știm care pagini sunt cele care ar trebui să fie în memorie pentru a reduce defectele. Este posibil să nu existe suficientă memorie fizică pentru toate procesele din sistem, pur și simplu. Posibilele soluții sunt swapping-ul sau cumpărarea a mai multă memorie.

Algoritmii de înlocuire a paginilor evită thrashing-ul .

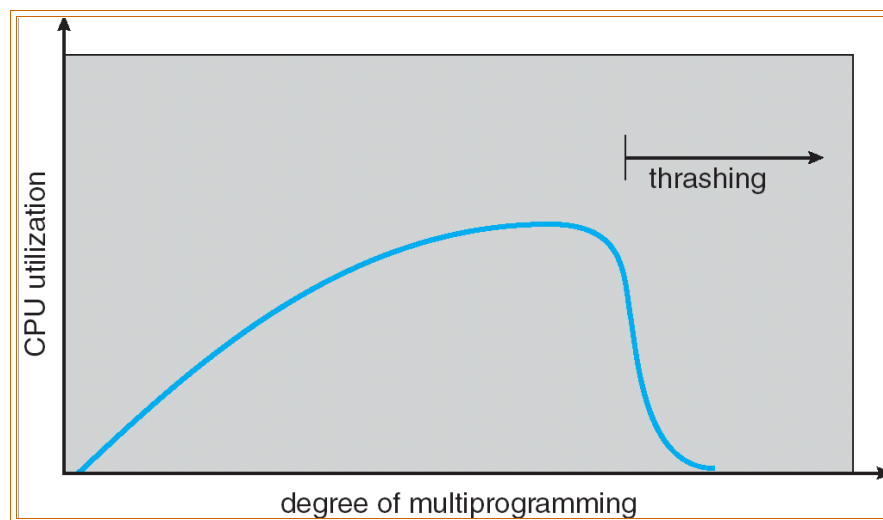


Figura 12 Thrashing

Surse bibliografice: [2],[5],[6]

4.7. PFF (Frecvența de eroare de pagină)

Este un algoritm de spațiu variabil care folosește o abordare mult mai “ad-hoc”.

Monitorizează frecvența de eroare pentru fiecare proces. Dacă aceasta este peste un prag înalt ales, i se dă procesului mai multă memorie, astfel încât să fie mai puține defecte, dar nu funcționează întotdeauna (Ex.: FIFO, anomalia lui Belady). Dacă frecvența de eroare este sub un alt prag ales, jos, se ia memorie de la proces. Ar fi normal ca el să producă mai multe defecte, dar nu este regulă.

Este greu să se utilizeze PFF pentru a distinge între schimbările de localitate și schimbările în dimensiunea WS-ului.

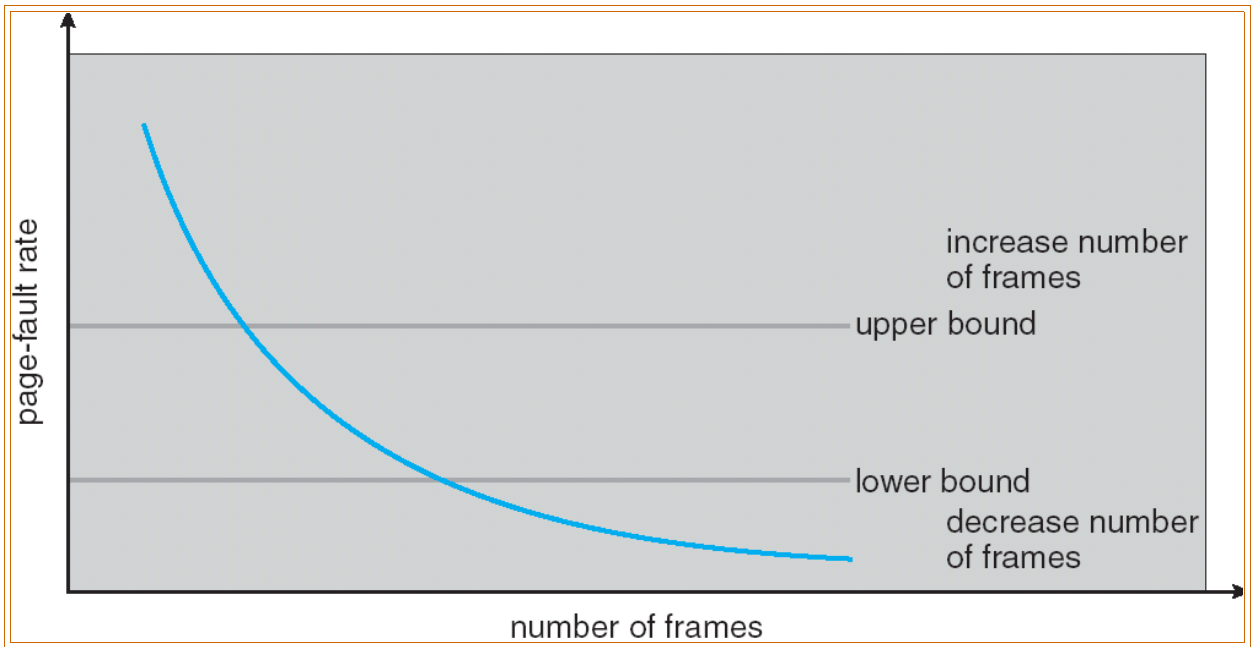


Figura 13 Frecvența de defect de pagina în funcție de numărul de cadre.

Surse bibliografice: [5],[6]

4.8. Alți algoritmi

4.8.1. “Low inter-reference Recency Set Replacement Policy” LIRS

Deoarece de a pretinde să acoperim toate aspectele acestui algoritm, vom prezenta modul în care acesta funcționează:

Blocurile de referință sunt împărțite în două clase: HIR (High Inter-reference Recency) și LIR (Low Inter-reference Recency). Fiecare bloc deține informații referitoare la istoria paginilor conținute, în cache, având intrările în înregistrările din cache ca HIR de non-rezidență.

Cache-ul, cu dimensiunea L (blocuri) este împărțit, la rândul său, într-o parte majoră (de dimensiune L_{lirs}) și una minoră (de dimensiune L_{hirs}).

$$L = L_{lirs} + L_{hirs}$$

												“Recency”	IRR
E									x			0	Inf
D		x						x				2	3
C				x								4	Inf
B			x		X							3	1
A	x					x		x				1	1
blocuri/timp virtual	1	2	3	4	5	6	7	8	9	10			

Figura 14 Cum este ales blocul “victimă” și cum sunt interschimbate stările LIR/HIR

x prezent într-o casuță a tabelului înseamnă că rândul corespunzător este referit la momentul virtual din coloana corespunzătoare. Recency și IIR reprezintă valorile ce corespund momentului virtual 10, pentru fiecare bloc.

Astfel, presupunând $L_{lirs} = 2$ și $L_{hirs} = 1$, la momentul 10, setul HIR = { C,D,E}, setul LIR = { A, B}. Singurul bloc rezident este E.

Surse bibliografice: [13]

4.8.2. “Enhanced second chance/clock”

Algoritmul folosește clasele (R,M), definite în tabelul din **Figura 3**. Înlocuiește prima pagină pe care o găsește în clasa nevidă cu cea mai mică prioritate.

Un dezavantaj este acela ca poate fi necesar să se scaneze coada de mai multe ori până la înlocuire.

Dacă “victima” este “murdară” (M=1), atunci se consumă timp cât aceasta este swappată. Este, deci, mai convenabil să se aleagă o victimă “curată” (M=0).

Astfel este determinată prioritatea claselor:

- clasa 1 (R=0,M=0): pagina ce îi aparține este un bun candidat pentru înlocuire;
- clasa 2 (R=0, M=1): vechile pagine trebuie să fie scrise; nu este un candidat la fel de bun ca anterioara;
- clasa 3: (R=1,M=0): pagina este recent referită, ceea ce o face un candidat prost;
- clasa 4: (R=1,M=1) : în nici un caz un bun candidat, din cauză că trebuie, de asemenea, să fie scrisă.

Surse bibliografice: [3],[5]

4.8.3. “Page Buffering”

Până acum a fost presupusă aplicarea algoritmilor ori de câte ori este nevoie ca o pagină să fie adusă în memorie.

Majoritatea algoritmilor discutați anterior sunt prea costisitori pentru a fi utilizați la fiecare defect de pagină.

“Page Buffering” reține o “piscină” a paginilor libere. Algoritmul de înlocuire este utilizat abia atunci când aceasta devine prea neîncăpătoare (“low water mark”), eliberând sufficient spațiu pentru o nouă aprovizionare cu pagini (“high water mark”).

Când apare un defect de pagină, se alege un cadru din lista liberă. Cadrele din această listă conțin încă informația precedentă, și pot fi salvate dacă pagina virtuală este referită anterior realocării.

Surse bibliografice: [1]

4.9. Concluzii asupra performanței principalilor algoritmi

Atunci când are loc un defect de pagină este necesar ca sistemul de operare să elibereze un cadru de pagină prin expulzarea unei pagini care are cea mai mică probabilitate să fie solicitată în lucru, în viitorul apropiat. Algoritmii de înlocuire a paginilor sunt cei care fac posibilă o selecție cât mai bună a paginii “victimă”.

Iată cum “se descurcă” algoritmi de înlocuire a paginilor întâlniți frecvent și prezentați în cele anterior spuse:

1. OPRA este un algoritm optimal, ce conferă (teoretic) numărul minim de defecte de pagină. Utilizat ca etalon pentru ceilalți algoritmi.
2. FIFO înlocuiește pagina încărcată în cel mai îndepărtat moment de timp. Poate îndepărta pagini importante (necesare)
3. LRU înlocuiește pagina referită la momentul cel mai îndepărtat în viitor. Excelent ca performanță, dar dificil de implementat.
4. CLOCK înlocuiește pagina cu vechimea cea mai mare.
5. WORKING SET păstrează în memorie setul de pagini cu frecvență de eroare minimă. Costisitor.
6. PFF mărește/micșorează setul de pagini în funcție de frecvența de eroare.
7. NRU prea radical.
8. A Doua Șansă ameliorare considerabilă a FIFO.
9. CEASUL (clock) realist.
10. NFU aproximare grosieră a LRU.
11. AGING aproximare eficientă a LRU.

4.10. Caracteristici ale paginării în Windows XP/ Linux

Windows XP

- înlocuire a paginilor locală;
- FIFO pe process;
- Sunt luate pagini de la procese ce folosesc mai multe pagini decât dimensiunea minimă a WS-ului lor, și oferite altor procese care suferă de lipsa paginilor disponibile;
- Procesele sunt lansate cu WS = 50 de pagini, implicit;
- Sistemul monitorizează frecvența de eroare (defect) și ajustează WS-ul în concordanță cu aceasta;
- La apariția unui defect de pagină, grupuri de pagini din jurul paginii lipsă sunt aduse în memorie.

Linux

- înlocuire a paginilor globală;
- utilizează algoritmul clock pentru înlocuire;
- paginile îmbătrânesc cu fiecare pas parcurs de “mâna” indicatoare a clock-ului;
- paginile ce nu au fost folosite vreme îndelungată vor avea în final valoarea 0.
- sistemul este în plină dezvoltare (noutăți apar încontinuu).

Surse bibliografice: [5]

4.11. Referințe

1. LABORATOIRE DE SYSTEMES REPARTIS, Systèmes d'exploitation, Implantation de programmes en mémoire, Mémoire virtuelle, Applications, Gestion de la mémoire, Alain Sandoz, Semestre été 2007;
2. OS @ SEP501 (Spring 2007) -- Jin-Soo Kim (jinsoo@cs.kaist.ac.kr);
3. Department of Computer Science, UMass Amherst, Andrew H. Fagg, CMPSCI 377: Operating Systems, Lecture 22;
4. CS 241 Fall 2007 System Programming, Memory Paging and Replacement, Lawrence Angrave;
5. CSCI 4401 Principles of Operating Systems I, Virtual Memory, Vassil Roussev vassil@cs.uno.edu;
6. Operating Systems Concepts, Silberschatz, Galvin and Gagne, 2005;
7. Sistemi a microprocessore, Memoria Virtuale, Corso di Laurea in Ingegneria dell'Informazione, Università degli Studi di Firenze ;
8. Virtual Memory: Page Replacement, Victoria University of Wellington, New Zealand;
9. A study of replacement algorithms for a virtual-storage computer by L. A. Belady;
10. Tanenbaum, Andrew S. - Operating Systems Design and Implementation, 2Ed;
11. Jackson libri - Tanenbaum - Moderni Sistemi Operativi;
12. CSE 120, Principles of Operating Systems, Fall 2004, Lecture 11: Page Replacement, Geoffrey M. Voelker;
13. LIRS: An Efficient Low Interference Recency Set Replacement Policy to Improve Buffer Cache Performance ;
14. Andrew Tanenbaum, Modern Operating Systems.

Notă: Internetul a fost cel cu ajutorul căruia au fost acumulate materialele bibliografice.

5. Modelarea algoritmilor de înlocuire a paginilor *Velican Valentin*

5.1. Definiție. Scurt “istoric”.

Într-un sistem de operare ce utilizează paginarea pentru organizarea memoriei virtuale, **algoritmii de înlocuire a paginilor** decid care dintre pagini trebuie înlocuită (swap out) când o nouă pagină din memorie trebuie alocată.

Din punct de vedere istoric acest subiect a fost foarte dezbătut în perioada anilor '60; '70 rezultând, prin studierea sa, algoritmi LRU foarte sofisticăți. Însă, din momentul respectiv și până în prezent, câteva ipoteze “tradiționale”, luate în calcul la dezvoltarea acestor algoritmi, au devenit un real obstacol pentru performanță având drept finalitate o nouă intensificare a cercetărilor. Spre exemplu pot fi reamintite:

- dimensiunea unității de stocare ce a crescut cu câteva ordine de mărime; în acest fel algoritmii ce necesitau o verificare periodică a fiecărei locație de memorie devin din ce în ce mai puțin practici.
- ierarhizarea memoriei a crescut; astfel, “costul” unui “cache miss” raportat de procesor a devenit din ce în ce mai ridicat.

- răspândirea programării orientate pe obiect, ce utilizează un număr ridicat de funcții precum și a unor tipuri de date sofisticate (e.g. structuri arborescente) aduce cu sine o utilizare “haotică” a memoriei.

Cerințele de la algoritmi s-au schimbat datorită diferențelor dintre sistemele de operare; în particular majoritatea sistemelor de operare moderne au reunit memoria virtuală cu fișierele cache de sistem. Astfel algoritmul este nevoit să aleagă o pagină atât din paginile adreselor virtuale cât și din cel al fișierelor cache. Cele din urmă au anumite particularități, precum cerințelor de scriere într-o anumită ordine. În plus algoritmul trebuie să ia în calcul și cerințele de memorie impuse de alte sub-sisteme ale kernel-ului (ce la rândul lor alocă memorie). Ca rezultat înlocuirea paginilor în kernel-uri moderne (Linux, FreeBSD, Solaris) tinde să lucreze la nivel mai scăzut față de un subsistem al memoriei virtuale.

5.2.Înlocuirea locală vs. Înlocuirea globală

Algoritmii de înlocuire pot să fie *locali* sau *globali*. Un algoritm de înlocuire “local” alege pentru schimb o pagină aparținând aceluiași proces în timp ce un algoritm “global” este liber să aleagă o pagină de oriunde din memorie.

Înlocuirea locală a paginilor implică existența unui mod de partiționare a memoriei care determină exact câte pagini vor fi atribuite unui anumit proces sau grup de procese. Cele mai cunoscute moduri de partiționare sunt algoritmii de “partiționare fixă” și “set echilibrat”. Avantajul înlocuirii la nivel local este independența pe care o are un proces în gestionare, nedepinzând de alte structuri de date globale.

5.3.Algoritmul optim de înlocuire a paginilor. Imposibilitatea realizării sale.

Algoritmul optim de înlocuire a paginilor (cunoscut și sub denumirea de OPT sau algoritmul “clairvoyant”) lucrează în felul următor: când o pagină trebuie schimbată, sistemul de operare înlocuiește pagina ce va fi utilizată cel mai târziu în viitor.

Acesta nu poate fi implementat în realitate datorită imposibilității de a calcula cu exactitate când vor fi utilizate paginile, decât în foarte puține cazuri. Totuși algoritmele actuale oferă performanțe aproape de nivelul optim – la prima rulare a programului, sistemul de operare reține toate paginile accesate și folosește aceste informații pentru organizarea lor ulterioară (la următoarele trimiteri în execuție a programului). Este de remarcat faptul că eficiența maximă nu este obținută la prima lansare în execuție și depinde totodată și de tiparul referințelor către memorie (ce trebuie să rămână asemănător la lansări în execuție ulterioare).

5.4.Algoritmi de înlocuire. Tipuri. Discuție.

- [Not Recently Used](#) -

Pentru ca un sistem de operare să rețină o statistică utilă a modului în care paginile sunt utilizate, fiecare pagină are asociați doi biți de status: R – setat atunci când pagina este citită/scriasă; M – atunci când pagina este scrisă. Biții sunt reținuți în intrările din tabela de paginare. Este important de reținut că fiecare acces la memorie este însoțit de un update al celor doi biți, fiind deci esențial ca setarea să fie făcută prin intermediul hardware-ului. Odată ce un bit a fost setat, el rămâne astfel până ce sistemul de operare îl resetează la zero.

Totuși dacă hardware-ul nu dispune de acești biți, ei pot fi simulați în felul următor : Când un proces este început ,intrările sale din tabela de paginare sunt marcate ca inexistente

în memorie. În momentul în care se referențiază o pagină, este returnat, deci, un “page fault”. Sistemul de operare setează deci bitul R (în tabelele sale interne) și schimbă intrarea din tabelă astfel încât să indice pagina corect, cu modul READ ONLY. Dacă o operație de scriere apare ulterior atunci același mecanism este reluat, permițând sistemului de operare să seteze bitul M iar pagina va fi accesată cu READ/WRITE.

Biții M și R pot fi utilizați pentru crearea unui algoritm în felul următor:

Când un proces este început, ambii biți ai fiecărei pagini sunt setați la zero de către sistemul de operare. Periodic bitul R este resetat, distingând astfel paginile care nu au fost referențiate (recent) față de cele utilizate. La apariția unui page fault, sistemul de operare verifică toate paginile și le împarte în patru categorii:

- clasa 0: nereferențiate, nemodificate
- clasa 1: nereferențiate, modificate
- clasa 2: referențiate, nemodificate
- clasa 3: referențiate, modificate

Deși la o primă vedere clasa 1 pare imposibil de realizat, totuși există posibilitatea ca unei pagini din clasa 3 să i se modifice bitul R la trecerea unei perioade de timp.

Algoritmul NRU scoate în mod aleatoriu o pagină din cea mai joasă clasă (din punct de vedere al valorii numerice) nevidă. Implicit, în algoritmul de față se consideră a fi mai bine să fie schimbată o pagină care nu a fost referențiată în perioada de timp decât o pagină liberă dar care este frecvent utilizată. Avantajele NRU sunt ușurința în înțelegere precum și eficiența sa, chiar dacă nu optimă.

- First In First Out (FIFO) -

Ideea algoritmului FIFO este evidentă citind chiar numele său – sistemul de operare urmărește utilizarea paginilor cu ajutorul unei cozi; cele mai recente dispuse la sfârșit și cele mai puțin utilizate la început. Când este necesară înlocuirea unei pagini, se selectează pagina de la începutul cozii (cea mai puțin utilizată) pentru a face schimbul.

Deși FIFO este intuitiv și ușor, se comportă destul de slab în aplicații, astfel, fiind rar utilizat în forma sa clasică (descrisă mai sus). Algoritmul prezintă așa numita “Anomalia Belady”. Cu alte cuvinte, crescând numărul de pagini pe care procesorul le poate încărca la un moment dat (în același timp utilizând metoda FIFO) crește probabilitatea de apariție a page fault-urilor.

Este utilizat în cadrul sistemului de operare VAX / VMS.

- Second Chance –

Reprezintă un algoritm de înlocuire a paginilor ce are la bază algoritmul FIFO. Se comportă sensibil mai bine decât FIFO cu foarte puține pierderi în ceea ce privește viteza. Deosebirea față de algoritmul precedent este că deși se caută să se schimbe tot prima pagină aflată în coadă, totuși se face o verificare suplimentară asupra acesteia, și anume se verifică bitul dacă bitul R este setat. Dacă acesta este zero atunci pagina este schimbată, altfel ea este adusă la sfârșitul cozii (deci coada poate fi privită și ca o coadă circulară), bitul R este resetat și algoritmul este reluat. Astfel, orice pagină primește o a doua șansă. Dacă toți biții R sunt setați atunci la a doua întâlnire a primei pagini din listă, aceasta este schimbată, deoarece acum ea are bitul de referențiere zero.

Așa după cum îi sugerează numele, algoritmul acordă o a doua șansă fiecărei pagini: o pagină veche care a fost referențiată, este probabil în utilizare și nu ar trebui schimbată spre deosebire de o pagină nereferențiată.

- Clock -

Clock este o variantă de FIFO mult mai eficientă decât "Second Chance" deoarece nu presupune trecerea paginilor, în mod constant, la spatele listei deși îndeplinește, în mare măsură, aceleași funcții precum S-C. Clock păstrează în memorie o listă circulară de pagini, indicând care este pagina cea mai "veche" din listă. Dacă apare un page fault, atunci este verificat bitul R al paginii celei mai puțin utilizate (indicată), dacă acesta este găsit ca fiind setat atunci se incrementează indicatorul de pagină. Procedura se repetă până este găsită o pagină ce îndeplinește condițiile de schimb.

- Least Recently Used (LRU) -

Algoritmul LRU, deși similar în denumire cu NRU (not recently used), diferă prin faptul că reține o statistică a utilizării paginilor pe o perioadă scurtă de timp, în timp ce NRU verifică utilizarea paginilor pe perioada ultimului tact. LRU lucrează pe principiul conform căreia paginile care au fost intens utilizate în ultimile câteva instrucțiuni, au probabilitate mare de a fi utilizate și în următoarele câteva instrucțiuni. Deși în teorie este apropiat de nivelul de performanță optimă, totuși este greu de implementat. Există câteva variante ale algoritmului care încearcă să reducă din complexitatea sa sacrificând cât mai puțin din performanță.

Cea mai complexă metodă este reprezentată de "metoda listei înlănțuite", ce conține o astfel de listă ce reține toate paginile din memorie. La spatele listei este adăugată pagina cel mai puțin utilizată recent iar la începutul listei pagina cea mai frecvent utilizată. Dezavantajul constă în faptul că fiecare referințiere a unei pagini duce la o modificare a listei ajungând astfel la un proces consumator de timp.

O altă metodă, ce are nevoie și de suport hard este următoarea: se presupune ca hardware-ul deține un contor pe 64 de biți ce este incrementat la fiecare instrucțiune. De fiecare dată când o pagină este accesată, primește o valoare egală cu cea a contorului la momentul accesului. Când o pagină trebuie schimbată, sistemul de operare alege acea pagină care are valoare alocată cea mai mică. Totuși această metodă nu este realizabilă datorită hard-ului actual, care nu conține un astfel de contor.

Un avantaj important al algoritmului LRU îl reprezintă faptul că se pretează unei analize statistice complete. A fost demonstrat, spre exemplu, ca un LRU nu poate avea mai mult de N ori erori decât algoritmul optim; unde N este proporțional cu numărul de pagini. Pe de altă parte LRU are tendința să piardă din performanță în multe cazuri des întâlnite de accesare a paginilor. (spre exemplu lucrul cu o aplicație ce realizează o buclă peste N+1 pagini, în timp ce există N pagini reținute de LRU, aduce cu siguranță câte un page fault la fiecare acces). Multe modificări în LRU, urmăresc deci să prezică diferitele moduri în care sunt utilizate paginile, alegând pentru cazurile în care algoritmul nu mai prezintă randament alte modalități de înlocuire a paginilor.

Alte variante asupra LRU sunt așa numitul LRU-K, ce aduce îmbunătățiri în ceea ce privește localitatea în timp. Este cunoscut și sub denumirea de LRU-2

Algoritmul ARC. Extinde LRU și aduce îmbunătățiri datorită unei mai bune organizări a listei ce reține utilizarea paginilor. El păstrează o istorie a paginilor recent schimbate și o utilizează în alegerea proprietăților paginilor stocate pentru o mai bună decizie în eventualitatea unui swap.

- Algoritmul de alegere aleatoare (Random) -

Acest algoritm înlocuiește în mod aleator o pagină din memorie. Avantajul clar al acestui algoritm este simplitatea sa, timpul pierdut pentru alegerea paginii de schimb după o anumită organizare fiind nul. În general se comportă mai bine decât FIFO iar pentru

referențieri prin ciclări se comportă mai bine chiar decât LRU, deși în general acesta din urmă este mai performant. OS/390 (sistem de operare creat de IBM pentru utilizarea în mainframe-uri) folosește LRU iar în cazul scăderii performanței trece la utilizarea unui algoritm aleator.

- Not Frequently Used -

NFU generează mai puține page fault-uri decât LRU în cazul în care tabela de pagini conține pointeri nuli.

NFU are nevoie pentru funcționare de un counter, fiecare pagină în parte având desemnat un astfel de counter, inițial fiind setat cu valoarea zero. La fiecare clock, paginile care au fost referențiate în acest interval vor avea counterul incrementat cu 1. Se va putea urmări în acest mod frecvența utilizării paginilor. Pagina cu valoarea cea mai mică a counterului va fi schimbată.

Principala problemă cu acest algoritm este aceea că urmărește frecvența de utilizare fără a ține cont de timpul cât sunt utilizate. Aceasta va avea drept rezultat o scădere a performanței.

5.5. Concluzii

Concluziile care se pot trage sunt următoarele:

- algoritmul optimal de înlocuire a paginilor nu poate fi implementat datorită imposibilității de a prezice cu exactitate ce pagină urmează să fie utilizată de procesul care rulează.
- algoritmul NRU se distinge prin simplitate.
- algoritmul FIFO nu are eficiența dorită, poate schimba pagini necesare (nu verifică dacă pagina a fost referențiată). De asemenea prezintă anomalia Belady.
- Second Chance, deși mai performant decât FIFO rulează lent.
- Clock este o alternativă bună față de Second Chance deoarece păstrează performanțele sale și în același timp ameliorează timpul necesar rulării algoritmului.
- LRU prezintă performanțe excelente dar este în general dificil de implementat.
- NFU este o aproximare destul de "brutală" a LRU; prezintă performanțe superioare doar în anumite condiții.
- pentru atingerea unor performanțe cât mai apropiate de optim este bună utilizarea mai multor algoritmi și alegerea celui mai bun în funcție de condițiile de moment. (exemplu: OS/390)

5.6. Rezumat

- Din punct de vedere istoric acest subiect a fost foarte dezbătut în perioada anilor '60; '70.
- În prezent cerințele de la algoritmi s-au schimbat datorită diferențelor dintre sistemele de operare.
- Algoritmii de înlocuire pot să fie locali sau globali.
- Algoritmii optimi de înlocuire a paginilor nu pot fi implementați în realitate datorită imposibilității de a calcula cu exactitate când vor fi utilizate paginile
- Algoritmii NRU se distinge prin simplitate.

- Algoritmul FIFO nu are eficiența dorită, poate schimba pagini necesare (nu verifică dacă pagina a fost referențiată). De asemenea prezintă anomalia Belady.
- Second Chance, deși mai performant decât FIFO rulează lent.
- Clock este o alternativă bună față de Second Chance deoarece păstrează performanțele sale și în același timp ameliorează timpul necesar rulării algoritmului.
- LRU prezintă performanțe excelente dar este în general dificil de implementat.
- NRU este o aproximare destul de “brutală” a LRU; prezintă performanțe superioare doar în anumite condiții.
- Pentru atingerea unor performanțe cât mai apropiate de optim este bună utilizarea mai multor algoritmi și alegerea celui mai bun în funcție de condițiile de moment. (exemplu: OS/390)

5.7.Referinte

* Exemplu de implementare a unui algoritm FIFO în limbajul de programare C (găsirea paginii celei mai vechi) :

```
int fifo_alg() {
int save_reg,page,frame,earliest_time;
save_reg = FTBR;
earliest_time = INFINITY;
do {
FTBR++;
if (FTBR >= FTBLR) FTBR = 0;
page = FTBR->entries[FTBR].page;
if (page == NIL) return(FTBR);
if (PTBR->entries[page].first_time < earliest_time) {
earliest_time = PTBR->entries[page].first_time;
frame = FTBR;
}
} while (FTBR != save_reg);
FTBR = frame;
return(FTBR);
}
```

unde: - PTBR este un pointer către tabela de pagini.
- FTBR este un pointer către tabela de frameuri
- FTBR este un pointer către tabela de adrese
- FTBLR este numărul total de locații (din spațiul de adrese fizice) disponibile
- algoritmul returnează cea mai veche pagină acesată

*

- 1) www.informit.com/
- 2) <http://en.wikipedia.org/>
- 3) <http://gaia.ecs.csus.edu/~zhangd/oscal/PagingApplet.html> - script ce prezinta simularea algoritmilor de înlocuire a paginilor.
- 4) <http://www.sci.csu Hayward.edu/>

6. Probleme de proiectare pentru sisteme de paginare *Stoian Andrei*

6.1. Alocare globală vs alocare locală

Când un proces declanșează un page fault (pagina accesată nu se găsește în memoria principală), un algoritm de alocare locală a paginilor selectează pentru înlocuire o pagină care aparține acelui proces (sau un grup de procese care împart aceeași partiție de memorie). Algoritmul de alocare globală a paginilor poate selecta orice pagină din memorie pentru înlocuire.

Alocarea locală a paginilor presupune existența unei partiționări a memoriei care să determine câte pagini sunt alocate unui proces sau unui grup de procese.

În general, algoritmi de alocare globală sunt mai eficienți decât cei de alocare locală, mai ales când dimensiunea setului în lucru variază pe toată durata procesului. Dacă se folosește un algoritm pentru alocare locală, iar dimensiunea setului în lucru crește, va rezulta thrashing-ul. Dacă dimensiunea setului scade, algoritmi locali folosesc memoria în mod ineficient. Pentru un algoritm de alocare globală, sistemul trebuie să decidă mereu câte pagini să aloce pentru fiecare proces. O metodă este de a monitoriza dimensiunea setului în lucru cu ajutorul biților de vârstă, dar în acest caz nu se poate spune cu certitudine dacă nu va rezulta thrashing-ul. Dimensiunea setului se poate modifica în câteva microsecunde, pe când biții de vârstă durează câteva tacturi de ceas.

O altă metodă este existența unui algoritm care să aloce pagini proceselor în lucru. Se determină numărul proceselor în lucru și se alocă pentru fiecare un număr egal de pagini.

La folosirea unui algoritm de alocare globală, este posibil să alocăm pentru început un număr de pagini proporțional cu dimensiunea procesului, dar alocarea trebuie actualizată pe măsură ce procesele se execută. Astfel se folosește algoritmul frecvenței page fault-urilor (PFF – page fault frequency), care stabilește când să se mărească sau să scadă numărul de pagini alocate unui proces. Acest algoritm nu dă nici o informație asupra cărui pagini să fie înlocuite când apare un fault, ci doar controlează dimensiunea setului alocat.

6.2. Controlul alocării paginilor

Chiar dacă se folosesc cei mai buni algoritmi de alocare a paginilor, se poate întâmpla să rezulte thrashing-ul. De fapt, de fiecare dată când suma tuturor seturilor în lucru ale proceselor este mai mare decât capacitatea memoriei, rezulta thrashing-ul.

Acest lucru se poate observa atunci când algoritmul PFF indică existența unui proces care are nevoie de mai multă memorie, dar nici un proces nu are nevoie de mai puțină memorie. Situația în care se ajunge este destul de delicată: nu se poate aloca memorie în plus procesului care are nevoie, deoarece nu avem de unde. Trebuie să "scăpăm" temporar de anumite procese. Singura metodă pentru a realiza acest lucru, fără a termina forțat procese, este să le mutăm pe disc (swap). Astfel, paginile rămase libere de la procesele mutate vor fi acum alocate proceselor care au mai multă nevoie de memorie.

6.3. Dimensiunea paginilor

Dimensiunea paginilor este de obicei determinată de arhitectura procesorului.

Un sistem cu dimensiune mică a paginilor utilizează mai multe pagini, astfel că tabela de pagini (page table), care este utilizată pentru a face legătura între adresele virtuale și adresele fizice, va ocupa mai mult spațiu. De exemplu, dacă un spațiu de adrese virtuale de dimensiune 2^{32} este împărțit în pagini de dimensiune 4KB (2^{12}), numărul paginilor virtuale este de 2^{20} ($20=32-12$). Dacă dimensiunea paginii este de 32KB (2^{15}), atunci sunt necesare mai puține pagini, adică 2^{17} .

Pentru a face legătura între adresele fizice și adresele virtuale, procesoarele au nevoie de o tabela specială, numită Translation Lookaside Buffer, sau TLB, care este "consultată" la fiecare acces de memorie. TLB nu are o dimensiune fixă, iar când nu poate completa cerere de adresă (TLB miss), tabelele de pagini trebuie cautate manual (hardware or software), pentru a se găsi adresa căutată. Dacă dimensiunea paginilor este mare, atunci tabela TLB poate memora mai multe legături între adresele virtuale și adresele fizice, evitându-se astfel acele TLB misses, mari consumatoare de timp.

Rareori se întâmplă ca un proces să aibă nevoie de un număr fix (întreg) de pagini. Astfel, ultima pagină va fi parțial ocupată, restul de memorie neocupată din acea pagină neputând fi utilizată. Astfel, dimensiunea mare a paginilor va crește probabilitatea ca porțiuni mari din memorie să nu poată fi utilizate. Dimensiunea mică a paginilor asigură o mai bună egalitate între dimensiunea memoriei și dimensiunea totală a proceselor.

De exemplu, dacă dimensiunea paginii este de 1MB, iar un proces are nevoie de 1025KB, atunci pentru acest proces vor fi alocate 2 pagini, rezultând astfel 1023KB care nu vor putea fi utilizați.

La transferul de pe disk în memoria principală, mare parte a timpului necesar este datorat timpilor de "găsire" a informației – seek time. Din această cauză, transferuri mari și rare sunt mult mai eficiente decât transferuri mici și dese.

Pe baza a doi parametri – dimensiunea medie a unui proces și dimensiunea unei înregistrări în tabela de pagini – se poate calcula dimensiunea optimă a unei pagini:

$$\text{dim ens. optima} = \sqrt{2 \cdot \text{dim ens. medie} \cdot \text{dim ens. inregistrare}}$$

De exemplu, pentru procese de dimensiune medie 1MB, iar înregistrări în tabela de pagini de 8 bytes, dimensiunea optimă a paginilor este de 4KB.

6.4. Spațiul datelor și spațiul instrucțiunilor

Instrucțiunile și datele ocupă spații de adrese diferite în memorie. Astfel, există o adresă 0 în spațiul adreselor instrucțiunilor, care pointează către o locație unde este memorată o instrucțiune; la fel și pentru date.

La un calculator cu acest design, ambele spații de adrese (al instrucțiunilor și al datelor) sunt paginate independent una de cealaltă. Fiecare are tabela ei de pagini, cu legături independente între adrese virtuale și adrese fizice de date și instrucțiuni. Când este nevoie de o instrucțiune, procesorul știe ca trebuie să folosească spațiul instrucțiunilor și tabela de pagini a instrucțiunilor. La fel și pentru date.

6.5. Pagini partajate

Intr-un sistem multiprogramat, este ceva normal ca mai mulți utilizatori să ruleze același program în același timp. Este evident că este mai eficient să se utilizeze aceleași pagini, decât să fie mai multe copii ale acelor pagini în memorie.

Dacă există spații separate pentru date și instrucțiuni, este evident că pentru a partaja programele, procesele vor folosi aceeași tabelă de pagini pentru instrucțiuni, dar tabele de pagini diferite pentru date. Tabelele de pagini sunt structuri de date independente de tabela de procese. Fiecare proces are doi pointeri în tabela de procese: unul pentru tabela de date și unul pentru tabela de instrucțiuni.

Când mai multe procese folosesc aceleași instrucțiuni, apar probleme cu paginile pe care le folosesc. Dacă un proces se execută și este eliminat din memorie, toate conținuturile paginilor care i-au fost alocate sunt șterse, iar paginile vor fi alocate altor procese. Astfel, procesele care foloseau paginile primului proces vor genera multe fault-uri, pentru a încărca în memorie paginile necesare lor. Similar, când un procesul își termină execuția, este esențial să se știe care dintre paginile lui sunt folosite de alte programe, pentru a nu fi golite din greșeală. Se folosesc astfel structuri de date speciale care memorează care pagini sunt folosite de mai multe procese și care sunt aceste procese.

În cazul datelor, fiecărui proces îi este alocată o tabelă de pagini, fiecare pointând către același set de date. Nu este nevoie să se copieze setul de date, pentru fiecare proces. Paginile de date sunt Read Only pentru fiecare proces în parte.

Atâta timp cât procesele nu modifică datele, ci doar le citesc, nu este nici o problemă. Când un proces modifică o anumită dată din memorie, se creează o copie a paginii, astfel ca fiecare proces are acum pagina lui de date. Acestea sunt acum Read Write, astfel că viitoare modificări ale datelor nu vor produce nici o violare în sistem.

Această strategie mărește performanțele sistemului, deoarece paginile al căror conținut nu este modificat nu vor fi copiate.

6.6. Eliberarea (curățarea) paginilor

Când este nevoie de o nouă pagină, iar nu este nici o pagină goală în memorie spre a fi scrisă cu noua informație, una din paginile actuale trebuie mutată pe disc (swap).

Toate sistemele au un proces numit "paging daemon", care cercetează periodic memoria și menține un număr suficient de pagini goale spre a fi folosite la nevoie. El selectează și transferă pe disc, prin algoritmi specifici de alocare a paginilor, paginile care au fost modificate după ce au fost încărcate.

6.7. Interfața memoriei virtuale

Memoria virtuală este transparentă proceselor și utilizatorilor, este un spațiu mare de adrese virtuale, pe un calculator cu puțină memorie fizică.

În unele cazuri, utilizatorii au puțin control asupra memoriei și pot modifica modul de comportare al divizelor programe.

Cea mai utilizată tehnică de management a memoriei este DSM – Distributed Shared Memory. Scopul este acela de a permite mai multor procese într-o rețea să aibă acces la același set de pagini, printr-un singur spațiu de adrese. Când apare un fault, procedura de tratare a fault-urilor detectează în rețea ce calculator are pagina cerută și trimite un mesaj prin care pagina este eliminată din tabela de pagini și din memoria sursă, este trimisă prin rețea celui care a cerut-o, care o va indexa într-o tabelă de pagini și o va alocă programului care are nevoie.

Bibliografie:

- 1) http://en.wikipedia.org/wiki/Main_Page
- 2) "Virtual memory: Issues of implementation" Bruce Jacob, Trevor Mudge, 1998

- 3) www.patentstorm.us/
- 4) <http://www.windowstlibrary.com/>
- 5) <http://www.cs.mun.ca/~paul/cs3725/material/web/notes/>
- 6) <http://www.cs.jhu.edu/~yairamir/cs418/os6/>
- 7) http://www.ginfo.ro/8_4-5/
- 8) <http://funinf.cs.unibuc.ro/~gheorghe/curs/arhCalc/lec/I11four.pdf>
- 9) Modern Operating Systems 2Ed - Tanenbaum

7. Probleme de implementare *Nistor Adrian*

Abstract

Sistemul de operare are deosebita sarcina de organizare a memoriei virtuale pentru rezultate mai bune in ceea ce privește viteza de calcul și volumul de informații manipulate. In acest sens intervin o serie de probleme care trebuie sa aibă o rezolvare cât mai bună. În cele ce urmează se iau în calcul câteva aspecte importante în ceea ce privește problemele de implementare ale unui sistem de operare.

In primul rând sistemul de operare trebuie sa organizeze managementul paginării în fiecare din momentele în care se poate afla un proces: crearea sa, desfășurarea în timp, „page fault”-urile și momentul de încheiere a unui proces. Tot sistemul de operare are rolul de a trata mesajele de tip page-fault. În acest sens trebuie să aibă un mecanism corespunzător care să întrerupă procesul activ, să transfere controlul rutinei de tratare a PF-ului sistemului de operare și apoi să redea controlul procesului întrerupt. În rezolvarea page fault-ului, datorită nevoii de a re-executa instrucțiunea pentru care s-a adus în memorie pagina căutată (ce inițial nu era găsită), este nevoie de un sistem care sa rețină instrucțiunea curentă. Astfel backup-ul instrucțiunii este sistemul care rezolvă problema. Un alt aspect care trebuie tratat este problema ca un proces să aibă nevoie de o pagină care îi este eliminată din memorie de către sistemul de operare. Pentru ocolirea acestei situații sistemul de operare dispune de un sistem de blocare a paginilor în memorie, care reprezintă soluția salvatoare. Termenul de „Backing store” reprezintă procesul prin care sistemul de operare scrie pe disc datele unei pagini care este eliminată din memorie. Acest ansamblu de acțiuni, salvarea pe disc a unor pagini, respectiv extragerea de pe disc în memorie a paginilor căutate stă la baza managementului efectuat de către sistemul de operare. De asemenea tot sistemul de operare trebuie sa fie capabil sa gestioneze mesajele page fault și cu un paginator extern.

Sistemul de operare are deosebita sarcina de organizare a memoriei virtuale pentru rezultate mai bune in ceea ce privește viteza de calcul și volumul de informații manipulate. In acest sens intervin o serie de probleme care trebuie sa aibă o rezolvare cât mai bună.

7.1. Implicațiile Sistemului de Operare în paginare

Exista patru cazuri legate de paginare în care sistemul de operare trebuie sa ia în considerare un set anume de acțiuni în ceea ce privește desfășurarea unui anumit proces: crearea unui nou proces, desfășurarea în timp a procesului, „page fault”-urile și momentul de încheiere a unui proces.

Când un nou proces este creat într-un sistem de paginare sistemul de operare trebuie să determine mărimea programului și a datelor la momentul inițial și să creeze o tabelă de pagini pentru acel proces. Spațiul trebuie alocat în memorie pentru tabela de pagini și trebuie inițializat. Tabela de pagini trebuie să fie în memorie atunci când procesul rulează. În plus pe

hard disk trebuie rezervată memoria astfel încât atunci când o pagina este trecută pe disc să aibă unde să fie salvată. Această zonă trebuie inițializată cu textul programului și datele astfel încât atunci când un nou proces începe să primească page fault-uri să se poată aduce paginile cerute de pe disc.

Când un proces este programat pentru execuție MMU trebuie resetat pentru noul proces și TLB golit pentru a scăpa de eventualele urme ale procesului anterior. Tabela de pagini a procesului curent trebuie să fie cea curentă, acest lucru făcându-se de obicei prin copierea unui pointer în niște registre hardware corespunzătoare. Opțional paginile unora din procese pot fi aduse în memorie pentru a reduce numărul de page fault-uri inițial.

Când apare un mesaj de tip page fault, sistemul de operare trebuie să citească registrele hardware să determine care adresă virtuală a generat eroarea. Din această informație trebuie să calculeze ce pagină este necesară și să o identifice pe disc. Trebuie apoi să găsească un spațiu disponibil pentru pagina nouă înlăturând unele din vechile pagini folosite rar, și să aducă noua pagină în locația creată. Apoi trebuie să salveze contorul de program pentru a permite re-executarea instrucțiunii indicând către pagina căutată.

Când un proces se încheie, sistemul de operare trebuie să elibereze tabela sa de pagini, paginile sale și spațiul de pe disc ocupat de paginile salvate pe hard disk. Dacă unele pagini sunt partajate cu alte procese paginile din memorie sau de pe disc pot fi eliberate abia când ultimul proces se termină.

7.2.Gestionarea erorilor de tip PF (Page Fault)

Tratarea Page Fault necesită un mecanism de tratare a excepției care să întrerupă procesul activ, să transfere controlul rutinei de tratare a PF-ului sistemului de operare și apoi să redea controlul procesului întrerupt. PF va fi recunoscut pe parcursul ciclurilor de acces la memorie. Deoarece instrucțiunea care a cauzat PageFault *trebuie reluată*, trebuie salvat în stivă (automat) Contorul de Program aferent acesteia. Pentru asta, există un registru special, EPC (*Exception Program Counter*). Tabela de pagini necesită prezența unui bit „rezident” care ne arată dacă pagina respectivă este sau nu în memorie. Uneori se utilizează termenul „valid” pentru a indica prezența paginii în memorie. O pagină „invalidă” este astfel o pagină nerezidentă sau care are o adresă ilegală. Cu toate acestea, mai logic este să avem doi biți în tabela de pagini: unul care să indice faptul că pagina este validă iar al doilea să arate dacă pagina este rezidentă sau nu în memorie. Întotdeauna când apare un fenomen *page-fault*, se pune problema administrării (rezolvării) acestuia. În general, pașii care se urmează sunt

- Procesul necesită o pagină ce nu este rezidentă în memorie;
- Se verifică în tabela de pagini dacă referința de memorie este validă sau nu;
- Dacă referința este validă dar pagina nu este rezidentă, se încearcă obținerea acesteia din memoria secundară;
- Se caută și se alocă un cadru liber (o pagină de memorie fizică neutilizată până în prezent; uneori este necesară eliberarea unei pagini de memorie);
- Se planifică o operație de acces la disc pentru a se citi acea pagină din memoria secundară în cadrul nou alocat;
- După scrierea paginii în memorie se modifică tabela de pagini; pagina este acum rezidentă în memorie;
- Se repornește instrucțiunea ce a generat *page-fault*

7.3.Backup-ul instrucțiunilor

Când este cerută adresa unei pagini care nu există în memorie, instrucțiunea în cauză va genera un mesaj de tip page-fault către sistemul de operare. Acesta rezolvă situația creată și trebuie să re-execute instrucțiunea pentru a se finaliza cu succes. Problema apare datorită faptului ca de exemplu pentru o instrucțiune cu 2 operanzi sistemul de operare nu va ști ce anume a cauzat eroarea de tip page-fault: primul operand sau al doilea; Relativitatea este întărită și de nesiguranța valorii păstrate de contorul de program care trebuie atent modificat la re-executarea instrucțiunii după rezolvarea problemei ce generase inițial PF-ul. În acest scop, unele procesoare sunt construite special pentru a avea o soluție la această problemă. Metoda este prezența unor registre interne ascunse care au menirea de a păstra valoarea contorului de program înainte de executarea oricărei instrucțiuni. În acest fel sistemul de operare poate ști cu exactitate cum sa revină corect la „momentul” în care trebuie sa re-execute o anumită instrucțiune. Totuși daca procesoarele nu prezintă astfel de soluții hardware, sistemul de operare trebuie sa vină cu modificări suplimentare asupra procesului de executare a instrucțiunilor pentru a preveni potențiala ambiguitate în cazul unei erori de tip page-fault.

7.4. Blocarea paginilor in memorie

O alta problemă care poate apărea are legătură cu alte surse de informații/date, de exemplu porturile I/O. Pe scurt este vorba de șansa, sau neșansa, ca o pagină corespunzătoare unui proces de aducere a datelor dintr-un asemenea port, sa fie eliminată în favoarea unui alt proces mai puțin lent. Pentru înțelegere se poate lua un exemplu posibil, nu neapărat cu șanse foarte mari de apariție. Pentru un proces care se presupune că trebuie să citească date de la o componentă periferică legată printr-un port de tip I/O va exista un buffer cu anumite pagini in memorie. Având în vedere ca procesul se desfășoară mai lent, sistemul de operare va „profita de moment” să permită în acest timp unui alt proces sa intre în execuție. Presupunem că procesul al doilea generează un mesaj page-fault iar sistemul de operare trebuie să elimine din memorie o pagină pentru a aduce în loc informația căutată de proces. Exemplul de față fiind unul cu un anumit scop, luăm în considerare cazul în care SO alege pentru eliminare chiar paginile din bufferul în care primul proces inițial aducea date de pe portul de I/O. Acest lucru ar duce la situația în care unele procese ar modifica datele in paginile care nu le sunt asociate. Pentru aceasta, în special în cazul porturilor I/O, soluția găsită este ca paginile din memorie să poată fi blocate pentru ca sistemul de operare sa nu le elimine în cazul unui mesaj page-fault.

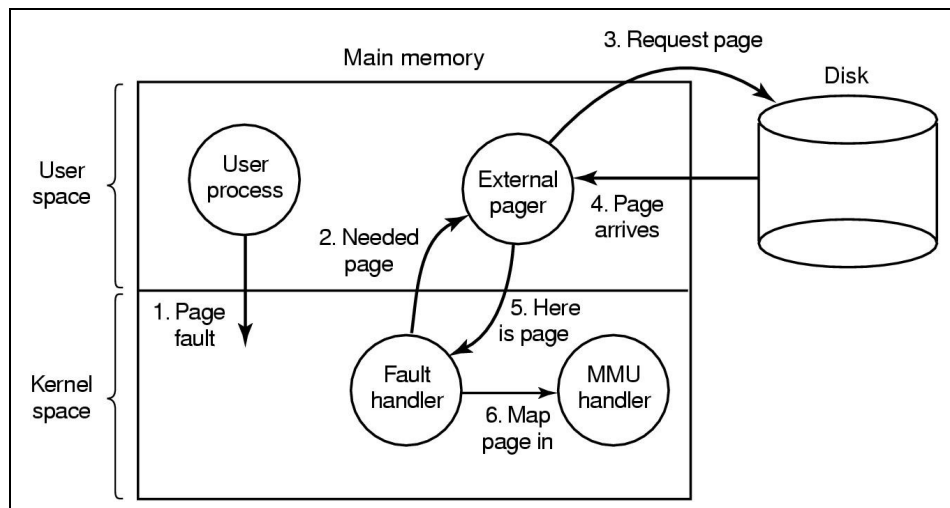
7.5. Termenul de "Backing Store"

Acest termen se referă la procesul prin care se scriu datele din memorie pe disc. Având în vedere ca SO-ul este cel care se ocupa de acest management, tot el trebuie sa aibă o bună metodă de a organiza pe disc datele corespunzătoare paginilor care „ies” din memorie. În acest sens există un spațiu delimitat, special destinat pentru paginile proceselor care vor fi în execuție, numit zonă „swap”. Sistemul de operare trebuie sa aleagă varianta optimă prin care să permită „trecerea” paginilor din memorie in zona de swap, ținând cont totodată și de faptul că un proces poate avea nevoie de o zonă mai mică sau mai mare de lucru, ba chiar își poate modifica această dimensiune pe parcursul desfășurării execuției. Din acest motiv s-au extras două situații prin care se poate rezolva această mică dilema. Ori să existe zone separate pe disc pentru program, date și stivă și sa se permită gruparea acestor zone; ori să nu se aloce nimic de la început, ci să se aloce spațiu pe disc pentru fiecare pagină nouă care trebuie să fie salvată in zona de swap. In acest al doilea caz mai este însă necesar și de o tabelă care sa rețină situația clară a paginilor salvate pe disc. Prima variantă,

paginarea în mod static a zonei de swap, presupune o corespondență clară de la început între paginile care trebuie scrise și spațiul alocat pentru așa ceva. Al doilea caz permite adaptarea în mod dinamic a spațiului în funcție de necesitățile de moment.

7.6. Separarea regulilor și a mecanismului

Pentru o buna organizare în cazul sistemelor complexe e nevoie de separarea mecanismului de reguli. Principiul poate fi aplicat și în cazul memoriei virtuale. Un exemplu simplu de separare este în figura de mai jos.

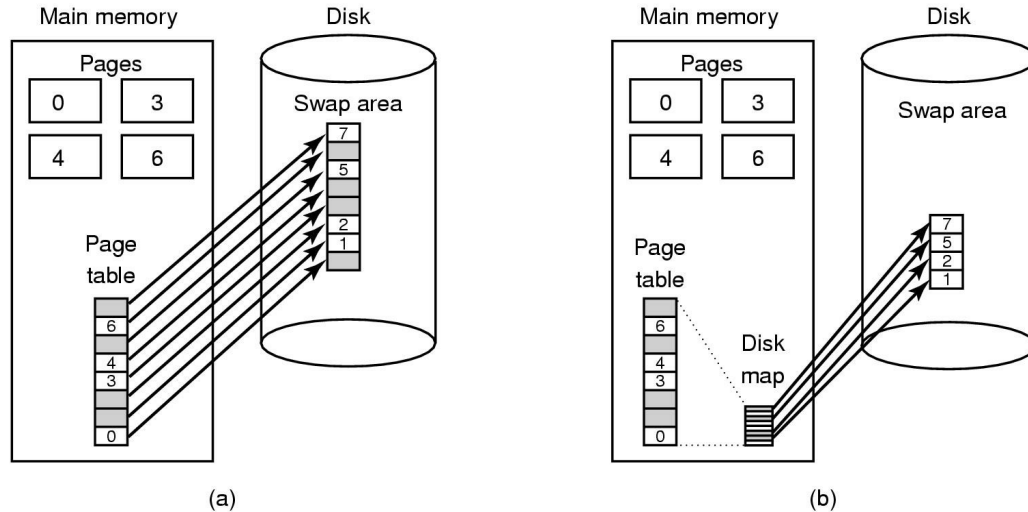


- 1 – O unitate de management a memoriei (MMU) la nivel scăzut
- 2 – O componentă de manevrare a mesajelor PageFault integrată în nucleu
- 3 – Un „paginator” extern ce rulează în spațiul de utilizator.

Toate detaliile despre cum ar trebui să funcționeze MMU sunt încapsulate în controlerul MMU, care are un cod dependent de mașină deci trebuie rescris pentru fiecare platformă pe care sistemul de operare este instalat. Controlerul mesajelor de tip "page fault" nu este dependent de mașină și conține marea parte a mecanismului de paginare.

Când un proces pornește, paginatorul extern este notificat pentru a pregăti harta paginilor de procese și să aloce pe hard disk spațiul necesar pentru "backing store". Cât timp procesul rulează paginatorul extern poate fi notificat în continuare.

Imediat ce procesul pornește, pot apărea mesaje tip page fault. Controlerul determină care pagina virtuală este necesară pentru folosire și trimite mesaj paginatorului extern "comunicând" situația. Paginatorul extern copiază apoi pagina în o porțiune din zona sa de adresa; transmite apoi controlerului de erori "page fault" unde se afla pagina. Controlerul demapează pagina din adresele paginatorului extern și întreabă controlerul MMU unde să pună în zona de adrese a utilizatorului la locul potrivit. Abia după acestea procesul poate fi restartat.



Avantajul principal al acestui tip de implementare este axarea codului mai mult pe module și flexibilitatea sporită. Dezavantajul, în schimb, este depășirea repetată a granițelor nucleului cât și a mesajelor transmise între diverse componente ale sistemului. În acest moment problema este controversată, dar pe măsură ce computerele devin din ce în ce mai rapide și software-ul devine mai complex, se va trece cu vederea peste acest sacrificiu în favoarea unui soft mai deschis și acceptabil pentru implementare.

Bibliografie:

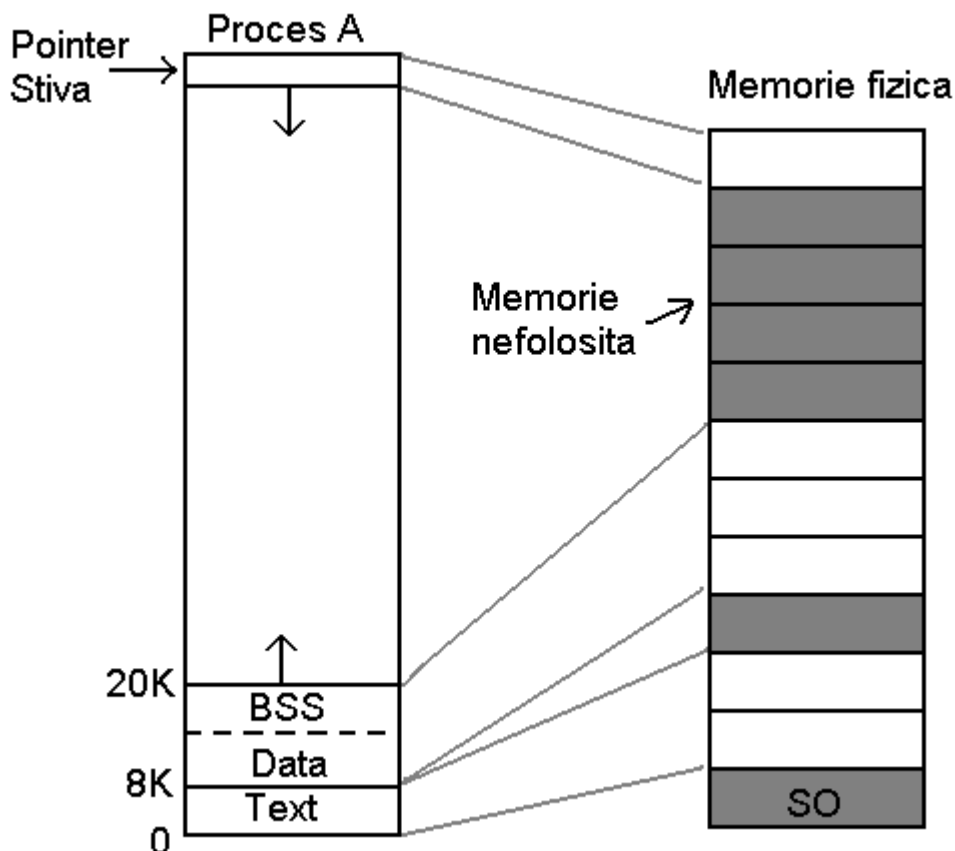
- 1) Modern Operating Systems 2Ed - Tanenbaum

8. Ledeanu Mihai Silviu

Spre deosebire de alte sisteme de operare, Unix-ul utilizează algoritmi foarte sofisticăți de gestionare a memoriei pentru a folosi foarte eficient resursele de memorie disponibile. De asemenea, modelul de memorie Unix încearcă să facă programele portabile și să facă posibilă implementarea Unix pe un număr mare de mașini diferite. Modelul de memorie s-a schimbat puțin odată cu trecerea timpului, el a functionat așa de bine încât nu s-a simțit nevoia multor îmbunătățiri ulterioare.

8.1. Concepte fundamentale

Fiecare proces Unix are un spațiu de adrese conținând trei segmente: text, data și stivă. Un exemplu de adrese de proces este Procesul A din figură:



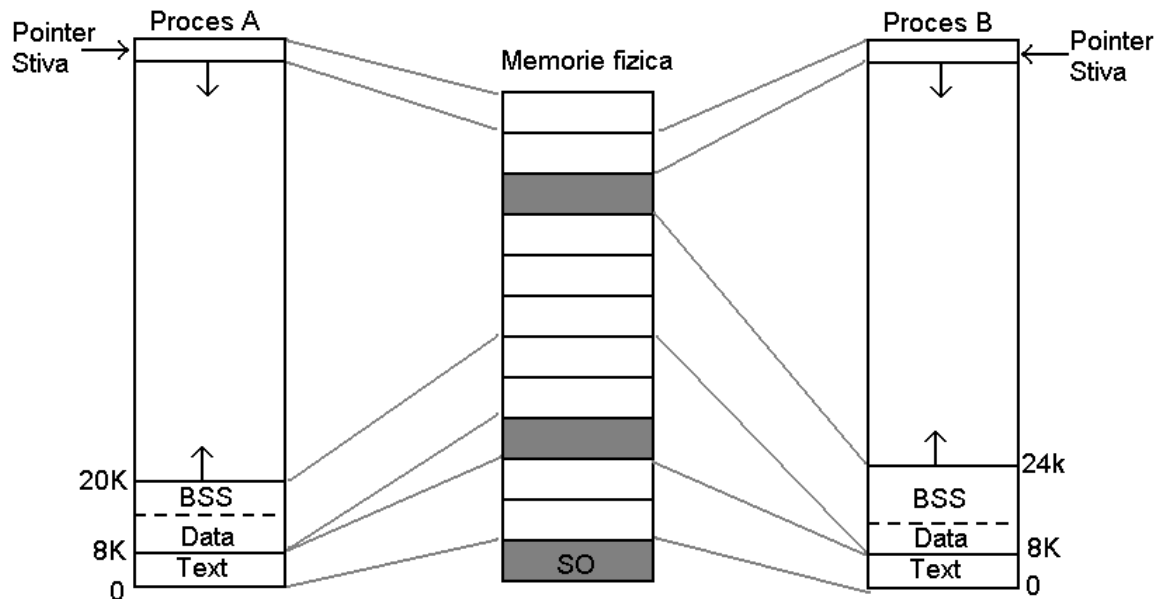
Segmentul de text conține instrucțiunile mașinii ce formează codul executabil al programului. Acesta e produs de compilator și asamblor translatând programul dintr-un limbaj în cod mașină. Acesta este de obicei read-only, programele care se automodifică au ieșit din circulație încă din anii 50, fiind prea greu de înțeles și de făcut debug.

Segmentul de date conține variabilele programului, șiruri, tabele și alte date folosite de program. Acesta e format din două părți: data inițializată și data neinițializată (numită BSS). Partea cu data inițializată conține variabile și constante ale compilatorului care au nevoie să fie inițializate când programul e pornit.

Spre deosebire de segmentul de text care nu se poate schimba, segmentul de date se schimbă, programele modificând tot timpul variabilele lor. Multe programe au nevoie să și aloce dinamic spațiu în timpul execuției. Unix permite segmentului de date să se mărească și micșoreze pe măsură ce memoria e alocată și dealocată. Există un system call, *brk*, ce permite unui program să seteze mărimea segmentului de dată, iar procedura *malloc*, din C/C++ e deseori folosită de programe pentru alocarea memoriei

Al treilea segment e cel de stivă. Aceasta pornește din varful spațiului de adrese virtuale și crește în jos către 0. Dacă stiva crește mai mult de partea de jos a segmentului de stivă, este generat un fault de hardware iar partea de jos a segmentului e coborâtă cu o pagină (adică este mărită cu o pagină). Programele nu gestionează mărimea segmentului de stivă. La începutul execuției unui program, stiva sa nu este goală, ci conține variabilele shell-ului și linia de comandă scrisă în shell pentru a rula programul. Astfel, programul are acces la argumentele sale.

Cand doi utilizatori rulează același program, ar fi ineficient să se pastreze două copii ale segmentului text. Astfel Unix suportă segmente de text partajate (împarțite). În figură vedem o posibilă împărțire a memoriei în care două procese (A și B) impart același segment de text.



Segmentele de dată și stivă nu sunt niciodată partajate, decăt după un *fork* (procedeu prin care un proces își crează o copie a sa, numita proces-copil) iar atunci numai acele pagini nu sunt modificate.

Pe unele calculatoare, hardware-ul suportă spații de adrese separate pentru program și pentru data. Unix-ul folosește aceasta proprietate, dacă este disponibilă. Astfel spațiile de adrese disponibile sunt dublate. De exemplu, pe un calculator cu adrese pe 16 biți, vor fi 2^{16} biți de adrese pentru program și 2^{16} pentru segmentul de dată și stivă.

Multe versiuni de Unix suportă memory-mapped files (fișiere cu memoria "mappată"). Astfel că e posibil să faci o hartă (map) al unui fișier într-o porțiune a spațiului de adrese a unui proces, așa că fișierul poate să fie scris și citit ca și cum ar fi un șir de octeți în memorie. Dacă mai multe procese mapează și accesează același fișier, aceeași memorie reală și pagini vor fi partajate tuturor proceselor. Aceasta dă voie mai multor programe să își partajeze data fără să fie nevoite să țina mai multe copii ale datelor în memorie.

8.2. Implementarea gestiunii de memorie în Unix

La începuturi, toate versiunile de Unix erau bazate pe swapping: când există mai multe procese ce pot fi ținute în memorie, unele dintre ele sunt înlocuite („swappate”) pe disc. Un proces este tot timpul swappat în întregime. Astfel ca un proces se afla ori în memorie, ori pe disc.

Odata cu apariția 3BSD (o distribuție Unix - Berkeley Software Distribution, versiunea 3BSD apare în 1979) universitatea Berkeley a adăugat și paginarea pentru a face față programelor tot mai mari ce erau scrise. Acum toate versiunile de Unix folosesc paginarea.

8.2.1. Swapping

Procesul în care unele pagini din memorie sunt mutate pe disc, și altele mutate în memorie, se numește swapping. Mutarea între memorie și disc este făcută de **swapper**. Mutarea (swappingul) din memorie către disc este făcută când kernelul ramana fara memorie libera din cauza unuia dintre evenimente:

- Un system call *fork* are nevoie de memorie pentru un proces-copil
- Un system call *brk* are nevoie sa isi extinda segemenul de data.
- O stiva creste și ramanae fara spațiu alocat.

De asemenea, când trebuia să fie aduc în memorie un proces ce a stat pe disc multa vreme, de obicei un alt proces trebuia să fie el eliminat pentru a face loc. Pentru a decide care proces să fie dus pe disc, swapperul caută procesele blocate, în asteptare (de exemplu pentru un periferic), fiind mai bine de eliminat un proces ce nu ruleaza decat unu care ruleaza. Din cele gasite este eliminat cel cu prioritatea mai mare. Dacă nici un proces nu este blocate, atunci se caută un alt proces după aceleasi criterii.

O data la cateva secunde swapperul verifica daca unul dintre procesele duse pe disc este gata de a reveni in memorie. Dacă există, este selectat cel care a ramas pe disc cel mai mult timp. Swapperul mai verifică și dacă swap-ul va fi unul ușor sau unul greu. Un swap ușor presupune că deja există destulă memori liberă pentru a fi adus procesul in memorie, un swap greu presupune că mai trebuie eliberată memorie pentru a fi adus procesul in memorie.

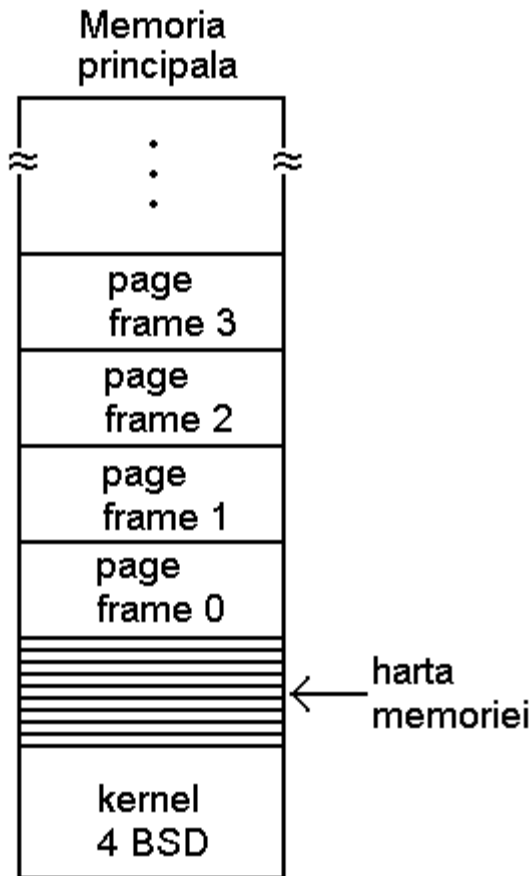
8.2.2. Paginarea

Odată cu apariția 3BSD de la Berkeley, paginarea a fost introdusă în aproape toate versiunile ulterioare de Unix.

Aici voi descrie principiul paginării distribuției 4BSD (apărută în noiembrie 1980) unde ideea este că un proces nu trebuie neapărat să fie adus în întregime în memorie ca să poată rula. Este nevoie numai de structura utilizatorului și de tabela paginii. Dacă acestea au fost aduse de swapper în memorie, atunci se spune că procesul este în memorie și poate rula. Paginile segmentelor de text, de date și de stivă sunt aduse în memorie dinamic, pe rând atunci când se face o referință către una din ele.

Paginarea este implementată parțial de kernel și parțial de un proces numit daemonul paginii. Daemon-ul paginilor este pornit periodic și verifică dacă are sarcini de făcut. Dacă el descoperă că numărul paginilor de pe lista de pagini libere este prea mic, atunci el începe să elibereze pagini.

Memoria principală, la 4BSD este organizată ca în figură, având trei părți: primele două părți, kernelul și core map (harta memoriei) sunt ținute permanent în memorie. Restul memoriei este divizată în cadre de pagini (page frame), fiecare putând conține o pagină de text, data sau stivă, o pagină de tabela, sau să fie liberă.



Harta memoriei (core map) conține informații despre conținutul cadrelor de pagini. Intrarea 0 din harta memoriei descrie cadrul de pagina 0, intrarea 1 descrie cadrul de pagină 1 și așa mai departe. Cu cadre de 1 KB și intrări în harta memoriei de 16 octeți, mai puțin de 2% din memorie este ocupată de harta memoriei.

Când un proces este pornit, este posibil să primească un page fault, pentru că una sau mai multe din paginile sale nu se găsesc în memorie. Dacă are loc un page fault, sistemul de operare ia primul cadru de pagini din lista de cadre libere, o elimină din listă și citește pagina de care este nevoie. Dacă lista de pagini libere este goală, procesul este suspendat până când daemon-ul de pagini eliberează o pagină

8.2.3. Algoritm de înlocuire a paginilor

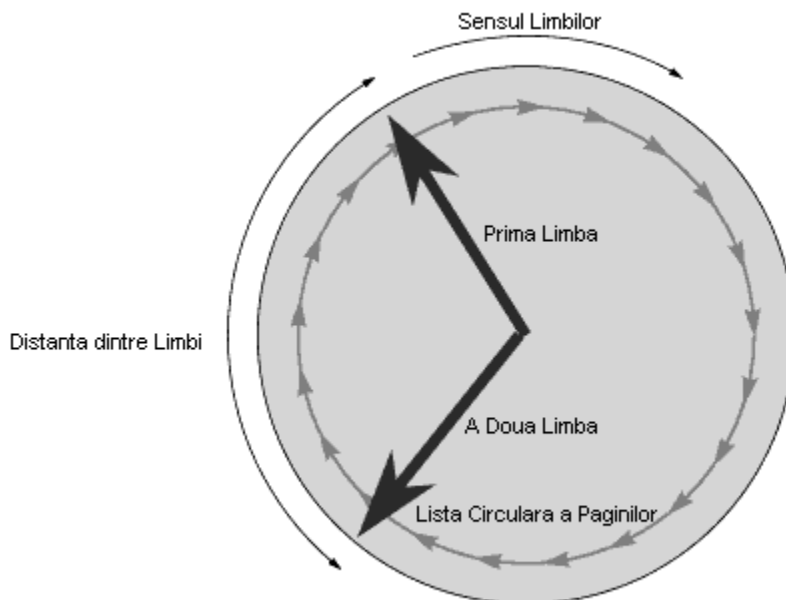
Algoritm de înlocuire a paginilor este executat de către daemon-ul de pagini. Acesta la fiecare 250 de milisecunde verifică dacă numărul de cadre de pagini libere este cel puțin egal cu un parametru al sistemului numit *lotsfree* (de obicei având valoarea unui sfert din memorie). Dacă sunt libere un număr insuficient de pagini, atunci daemonul începe să transfere pagini din memorie către disc, până ce se va ajunge la un număr de *lotsfree* de pagini libere. Dacă daemonul verifică că un număr mai mare de *lotsfree* de pagini sunt libere atunci se întoarce în modul de sleep. Dacă un calculator are multă memorie și puține procese active, atunci daemon-ul va sta cea mai mare parte a timpului în modul sleep.

Daemon-ul de pagini folosește o versiune modificată a algoritmului ceasului (clock algorithm). Algoritm de ceas scanează circular toate cadrele, asemenea unei limbi de ceas care se

primba de-a lungul tuturor orelor. La prima trecere, biții de folosire a cadrelor sunt setați zero. La a doua trecere, toate cadrele care nu au fost accesate de la ultima trecere a “limbii” sunt trecute în lista de cadre libere, după ce au fost mutate pe disc (mutarea pe disc are loc numai dacă acestea sunt “murdare”, adică dacă au fost modificate de când au fost aduse în memorie. Dacă nu au fost modificate, atunci nu are sens să fie remutate pe disc.)

Algoritmul acesta a fost folosit inițial de distribuțiile Unix ale universității Berkeley, dar mai tâziu s-a descoperit că pentru memorii mari, trecerea peste toate paginile dura prea mult. Așa că algoritmul a fost modificat în algoritmul ceasului cu două limbi (*two-handed clock algorithm*).

În acest algoritm daemon-ul ține doi pointeri către harta memoriei (cele două “limbi” ale ceasului. Prima setează biții de folosire zero, iar a doua verifică biții de folosire și mută paginile nefolosite pe disc. După ce amandouă limbile își fac verificările, fiecare avansează cu o poziție. Depinzând de deschiderea dintre cele două limbi, paginile vor avea mai mult sau mai puțin timp la dispoziție să nu fie accesate până sunt mutate pe disc.



Algoritmul Ceasului cu Doua Limbi

De fiecare dată când daemon-ul este pornit limbile avansează mai puțin ce un cerc întreg, numărul de avansări ale lor depinde de cât timp e nevoie pentru a ajunge ca numărul de pagini libere să ajungă la *lotsfree*. Dacă sistemul observă că rata de pagini este prea mare și numărul de pagini libere este tot simplu sub *lotsfree*, atunci sweeperul începe să înlăture procese din memorie.

Feluri de memorii:

- Principală – Memoria fizică RAM aflată pe placa de bază a calculatorului. Mai este numită și memorie reală. Aceasta nu include memoria video, cache-ul procesorului sau alte memorii periferice.

- Sistemul de fișiere – Memoria discurilor accesibilă prin pathname-uri (Pathname-urile sunt o secvență de simboluri și nume care identifică un fișier). Aceasta nu include spațiul de swap sau alte spații de stocare ce nu sunt adresabile prin pathname-uri. Nu include nici toate sistemele de fișiere de rețea.
- Spațiul de Swap – Memoria discurilor ce stochează data care nu e în memoria reală sau a file system-ului. Spațiul swap este cel mai eficient când e pe un disc separat sau partiție separată, dar câteodată este doar un fișier mare din sistemul de fișiere.

Memoria Sistemului de operare folosește:

- Kernel – Spațiul de memorie privat al Sistemului de operare. Acesta se afla tot timpul în memoria principală.
- Cache – Memoria care e folosită pentru a ține elemente ale sistemului de fișiere și alte operații de intrare/ieșire. A nu fi confundat cu cache-ul procesorului sau al hard disk-urilor, ce nu fac parte din memoria principală.
- Memorie Virtuală – Totalitatea spațiului de memorie adresabilă a tuturor proceselor ce rulează pe mașină. Locația fizică se poate afla pe oricare dintre cele trei tipuri de memorie.

Memoria Proceselor folosește:

- Data – Memoria alocată și folosită de program.
- Stiva – Stiva de execuție a programului (gestionată de sistemul de operare).
- Memoria mappată – Conținutul fișierului adresabil în spațiul de memorie al procesului.

8.3.Rezumat:

1. Concepte fundamentale

Modelul de memorie la Unix a suferit puține schimbări odată cu trecerea timpului.

Fiecare proces Unix are un spațiu de adrese ce conține trei segmente: segmentul de text, segmentul de date și segmentul de stivă.

Segmentul de text conține instrucțiunile mașinii ce formează codul executabil al programului. Este read-only.

Segmentul de date conține variabilele programului. La rândul său acesta conține datele inițializate și datele neinițializate. Acest segment se poate schimba.

Segmentul de stivă, la pornirea programului conține parametrii acestuia. Stiva crește în jos.

Segmentele de text pot fi partajate.

2. Implementarea gestiunii de memorie în Unix

2.1 Swappingul

La început, distribuțiile de Unix foloseau numai swappingul. Acesta presupune mutarea (swappingul) unui întreg program în memorie sau pe disc.

Swappingul este făcut de un scheduler numit swapper.

2.2 Paginarea

Odată cu trecerea timpului, în sistemele Unix a fost introdusă și paginarea, deoarece procesele creșteau din ce în ce mai mult și ocupau tot mai mult spațiu din memorie.

Ideea paginării este ca un proces nu trebuie adus cu totul în memorie pentru a putea fi rulat.

Paginarea este efectuată parțial de kernel și parțial de daemonul de pagini.

2.3. Algoritmi de înlocuire a paginilor

Algoritmul de înlocuire a paginilor este executat de daemon-ul de pagini care are grijă să fie cel puțin loto free pagini libere în memorie.

La început s-a folosit algoritmul ceasului, un pointer verifica lista hărții memoriei circular, iar dacă la a doua verificare pagina nu fusese activată atunci ea era scoasă din memorie. O pagină este rescrisă pe disc numai dacă a fost schimbată între timp, de când a fost adusă în memorie.

Pentru că parcurgerea întregii memoriei durează mult, s-a introdus algoritmul ceasului cu două limbi. Prima limbă setează pagina ca nefolosită, iar a doua verifică dacă între timp a fost folosită, dacă nu atunci este scoasă din memorie.

Bibliografie:

- 1) **Andrew Tanenbaum, Operating Systems, Design and Implementation**
- 2) <http://en.wikipedia.org>
- 3) <http://www.cpushack.net>
- 4) <http://www.freebsd.org>
- 5) notite curs SO cu St.St