

Capitolul 1 INTRODUCERE

Un sistem de operare este un program care gestioneaza hardware-ul calculatorului. Deasemenea asigura o baza pt. aplicatii si este ca un intermediar intre user si hardware-ul calculatorului.

1.1 Ce e un SO?

Este o parte importanta din aproape fiecare calculator. Un sistem poate fi impartit in 4 componente: Hardware, SO, aplicatiile shi utilizatori.

Scopul unui SO:

- Sa execute programele si sa usureze rezolvarea problemelor utilizatorilor.
- Sa faca sistemul mai usor de utilizat.
- Folosirea harware-ului unui calculator intr-ul mod eficient.

1. Hardware – Oferă resursele de calcul (CPU, memory, I/O devices).
2. SO – controleaza si coordoneaza utilizarea hardware-ului intre aplicatiile de la fiecare user.
3. Programele aplicatie – defineste modul in care resursele sistemului sunt folosite pentru a rzovla problemele de calcul ale utilizatorilor (compilayoare, baze de date, jocuri).
4. User (oameni, alte calc).

Schema pag 4 fig 1.1

1.1.1 User view(Cum este vazut de user)

Cum este vazut calculatorul de catre user variaza cu interfata folosita.

User sta in fata PC -> un asemenea sistem este facut ca un user sa detina monopolul asupra resurselor. In acest caz, SO este proiectat sa poate fi folosit cat mai usor(ease of use), cu putina grija asupra perfomantei si nici una fata de utilizarea resurselor.

Sistem de tip mainframe(sau minicomputer) -> SO este dezvoltat pt a maximiza utilizarea resurselor.

Statie conectate la retea -> SO este realizat sa faca un compromis intreutilizare individuala si utilizarea resurselor.

1.1.2 System View(cum este vazut SO de sistem)

Din punctul de vedere al calculatorului SO este un program care este cel mai apropiat hardware-ului. Putem vedea SO ca un alocator de resurse.

O alta vedere al unui SO pune accentul pe necesitatea de a controla dispozitivele de I/O si programele user. Aici SO este un program de control care previna erorile in executia programelor si utilizarea incorecta a calculatorului.

O definitie mult mai uzuala este ca un SO este programul care ruleaza tot timpul pe calculator -> kernel.

1.1.3 Scopul sistemului(System Goal)

Un SO este mai usor de definit daca spunem ce face decat ce este. Scopul principal al unor SO este sa ofere **comoditate** pt utilizator.(este cazul PC-urilor mici)

Alte SO au scopul principal de a avea operatii **eficiente** cu sistemul calcu-lui.(cazul sistemelor mari, multiuser).

In trecut eficienta era mai improtanta decat comoditatea. Cu timpul au fost adaugate intefete grafice -> comoditate.

1.2 Sisteme mainframe

Au fost primele calculatoare.

1.2.1 Sisteme Batch

User-ul nu interactioneaza direct cu calculatoru. User-ul pregatea un job si il introducea (submitted) in calculator. Job-ul era in forma cartelelor perfoarate; dupa un timp aparea rezultatul(output). Output-ul consta in rezultatul programului, continutul final al memoriei si continutul registrelor. SO era simplu, scopul lui era sa transfere controlul de la un job la urmatorul. SO era tot timpul rezident in memorie.

Pentru a mari rata de procesare: Operatorii grupau joburi cu nevoi asemanatoare si le rula ca un grup.

Fig 1.2 pg 8 memoria unui sistem batch simplu

Introducerea discurilor a permis pastrarea job-urilor pe disc. Avand acces direct la mai multe job-uri, SO putea face job scheduling, putea gestiona resurslele mai eficient.

1.2.2 Sisteme multimrogramate

Cel mai improtant aspect la job scheduling-ului este abilitatea de multiprogram.

Multiprogramarea mareste utilizarea CPU prin organizarea job-urilor astfel incat CPU sa albat tot timpul unul de rulat. SO pastreaza mai mult joc-uril simultan in memorie. Cand un job trebuie sa astepte SO comuta la un alt job.

Multiprogramarea este primul caz in care SO trebuie sa ia decizii pt useri. Toate joburile sunt pastrate in job pool(piscina cu joburi :))).Aici toate taskurile rezidente pe disc asteapta sa le fie alocata memorie principala. Daca mai multe joburi sunt gata sa fie adusa in mem, si nu este loc pt toate atunci SO trebuie sa aleaga pe care le aduce. Responsabil de aceasta este **job scheduling**. Daca mai multe joburi sunt gata sa ruleze **CPU scheduling** decide care sa ruleze.

Ce terbuie sa ofere SO?

- n Rutine I/O oferite de sistem
- n Managementul memoriei – sistemul alocă memoria mai multor joburi
- n CPU scheduling – siistemul trebuie sa aleaga ce job sa ruleza.

1.2.3 Sisteme time-sharing

Time sharing este o extensie logică a multiprogramării. CPU execută mai multe joburi prin comutare de la unul la altul, dar comutarea se face așa repede încât user-ul poate interacționa cu fiecare program în timp ce rulează. Un sistem interactiv facilitează comunicarea directă între user și sistem. Timpul de răspuns este scurt.

Un SO time-sharing permite mai multor utilizatori să folosească calculatorul simultan. Doar un mic timp CPU este necesar fiecărui utilizator. Ca să se obțină timpi de răspuns rezonabili, joburile trebuie să fie introduse în evacuare din Mem. principală de pe și pe disk. (memorie virtuală).

Este furnizată comunicarea între user și sistem, când SO termină de executat o comandă, caută următoarea comandă de la un alt user.

1.3 Sisteme desktop

Sunt dedicate unui singur utilizator.

Dispozitivele I/O sunt keyboard, mouse, display, printer. Principalul scop este maximizarea confortului și sensibilitatea utilizatorului.

1.4 Sisteme multiprocesor(paralele)

Este un sistem cu mai multe CPU care comunică. Procesoarele împart bus-ul, clock-ul și câteodată memoria. Comunicarea are loc de obicei prin memoria partajată.

Este un sistem cuplat foarte strâns.

Avantaje:

- Prin numărul mare de procesoare, se speră creșterii vitezei de procesare.
- Cost mai mic decât mai multe sisteme uni-procesor prin partajarea resurselor.
- Fiabilitate crescută: Defectarea unui procesor nu va opri tot sistemul.

Multiprocesare simetrica(SMP):

- fiecare procesor rulează o copie identică a SO
- Mai multe procese pot rula deodată fără a deteriora performanța
- Cele mai multe SO moderne folosesc multiprocesare simetrică

Multiprocesare asimetrica(SMP):

- Fiecare procesor are un atribuit un anumit task
- Un procesor controlează sistemul; el schedulează și distribuie lucrul către procesoarele slave
- Este regăsită în sisteme foarte mari.

1.5 Sisteme distribuite

Sistemele distribuite depind de rețea pentru o funcționalitate corectă. Prin comunicare Sist. distribuite pot împărtăși taskurile și a furniza un set bogat de opțiuni pt useri.

Un sist distribuit distribuie calculul între procesoare distincte.

Este un sistem cu o relație de legătură mai ușoară decât la sisteme paralele. Fiecare procesor are memoria sa proprie. Procesoarele comunică unele cu altele prin linii de comunicație (bus-uri, linii de telefon).

Avantaje:

- Impartirea resurselor
- Viteza de calcul mare
- Siguranță
- Comunicare

Are nevoie de infrastructura unei rețele, rețea locală LAN sau întinsă WAN. Poate fi un sistem peer-to-peer sau un sistem de tip client-server.

1.5.1 Sistem client-server

Sistemele centralizate de azi au rolul de sisteme server și satisfac cererile de la clienți.

- Sisteme Compute-server (server calcul) – oferă o interfață către care clienții pot trimite cereri de a executa un proces. Ca răspuns ele execută procesul și trimit clientului rezultatul.
- Sisteme File-server – oferă o interfață unde clienții pot crea, modifica, citi și șterge fișiere.

1.5.2 Sisteme peer-to-peer

1.6 Sisteme clustered

Aceste sisteme adună mai multe CPU laolaltă. Calculatoarele clusterate împart spațiul de stocare și sunt strâns legate printr-o rețea LAN. Fiecare mașină poate monitoriza o altă mașină. Dacă mașina monitorizată cedează, atunci cea care monitorizează poate să preia memoria sa și să reia execuția aplicațiilor care rula.

Clustering asimetric: o mașină se află în modul de standby, și monitorizează serverul activ. Dacă acesta se defectează, mașina din standby devine server.

Clustering simetric: Două sau mai multe mașini rulează aplicații și ele se monitorizează una pe alta.

1.7 Sisteme Real-Time

Sunt folosite când sunt puse restricții de timp asupra operațiilor unui procesor. Are restricții de timp bine definite, care trebuie respectate, altfel sistemul se oprește.

Sistemele în timp real pot fi hard sau soft.

Sisteme în timp real hard

Este garantat ca un task critic să fie terminat la timp.

- Dispozitivul de stocare secundar este limitat sau lipsește. Datele sunt stocate în memorie temporară sau în memorie read-only (ROM)
- Se generează conflicte cu sistemele time-sharing și nu sunt suportate de SO de uz general.

Sisteme în timp real soft

Task-urile critice au prioritate fata de celelelalte, aceasta prioritate este pastrata pana cand se termina task-ul.

- Utilizare limitata in controlul industrial si robotica.
- Folositor in aplicatii care necesita facilitati avansate ale SO.

1.8 Sisteme Handheld

Aici intra PDA-uri, celulare. Cei care dezvoltă SO pt dispozitive handheld trebuie sa ia in considerare:

- memoria limitata -> gestiune a memoriei foarte buna
- viteza mica a procesoarelor -> SO care sa nu supraincarce procesorul
- display-ul foarte mic

1.9 Migrarea trasaturilor

fig 1.6 pag 20

Trasaturi mai demult disponibile doar pe mainframe-uri au fost adoptate si la microcomputere.

1.10 Medii de calcul(Computing Enviroment)

1.10.1 Computing traditional

Mediu tipic pt aplicatii office. Tehnologia web impinge limitele mediului traditional.

1.10.2 Computing bazat pe Web

S-a marit importanta relelelor si viteza acestora.

1.10.3 Embedded Computing

Ruleaza SO embedded real time.

Capitolul 3 STRUCTURA UNUI SO

3.1 Componentele sistemului

Putem crea un sistem complex doar prin partitionarea lui in mai multe parti.

Componentele sistemului:

- n Gestionarea proceselor
- n Gestionarea memoriei principale
- n Gestionarea fisierelor
- n Gestionarea sistemelor I/O
- n Gestionarea memoriei secundare

- n Retelistica
- n Protectia sistemului
- n Interpretor de comenzi

3.1.1 Gestiunea Proceselor

Un proces este un program care se afla in executie. Un proces are nevoie de anumite resurse: timp CPU, memorie, fisiere, dispozitive I/O pentru a-si incheia task-ul. Executia unui proces este secventiala.

SO este responsabil de urmatoarele cand este vorba de gestiunea proceselor:

- Crearea si stergerea proceselor user si sistem
- Suspendarea si reluarea proceselor
- Asigura mecanisme pentru:
 - sincronizare
 - comunicatie
 - tratarea deadlock-urilor

3.1.2 Gestiunea memoriei principale

Memoria are un caracter volatil. Isi pierde continutul in caz de esec. De aceea este necesara o gestiune a memoriei.

SO se ocupa in cadrul gestiunii de memorie cu urmatoarele:

- Sa stie ce parte a memoriei este utilizate si de cine
- Sa decida care procese sa fie incarcate in memorie cand un spatiu devine disponibil.
- Sa aloce si dealoce memoria in functie de necesitati.

3.1.3 Gestiunea fisierelor

Este componenta cea mai vizibila a unui SO. Un fisier este o colectie, cu informatii asemanatoare, definita de utilizator.

SO este responsabil de urmatoarele:

- Crearea si stergerea fisierelor
- Crearea si stergerea directoarelor
- Oferirea sprijinului pt primitive pentru manipularea fisierelor si directoarelor.
- Maparea fisierelor in memoria secundara.
- Realizarea unui back-up a fisierelor pe medii de stocare stabile(non-volatile)

3.1.4 Gestionarea sistemelor I/O

Un sistem I/O este compus din:

- O componenta de gestionare a memoriei care inculde buffering, caching si spooling
- O interfata generala cu driver-ul dispozitivului.
- Drivere ce specifica dispozitivele hardware.

3.1.5 Gestionarea memoriei secundare

Mem principala este volatilă. Dimensiunea ei este redusă și nu poate păstra toate datele și programele permanente. Sistemul trebuie să dispună de un mediu secundar de stocare pentru a putea păstra toate datele. Calculatoarele moderne folosesc discurile ca mediu principal de stocare a datelor și a programelor.

În legătură cu gestiunea memoriei secundare SO este responsabil cu:

- Gestiunea spațiului liber
- Alocarea spațiului de stocare
- Disk scheduling

3.1.6 Reteaua

Un sistem distribuit este o colecție de procesoare care nu împart memorie, dispozitive periferice sau clock. Fiecare CPU are memoria lui. Procesoarele comunică între ele cu ajutorul unei rețele de comunicare. Comunicarea are loc folosind un protocol. Un sistem distribuit oferă posibilitatea ca user-ul să acceseze diferite resurse.

3.1.7 sistem de protecție

Protecția se referă la mecanisme de control a accesului programelor, proceselor sau utilizatorilor la resursele sistemului.

Mecanismul de protecție trebuie să ofere distincție între acces autorizat și neautorizat.

3.1.8 Sistem interpretor de comenzi

Unul din cele mai importante programe dintr-un SO este interpretorul de comenzi. Multe comenzi sunt date unui SO printr-o declarație de control, care se ocupă cu: crearea proceselor, manevrarea proceselor I/O, gestiunea memoriei, accesul la sistemul de fișiere... Programul care citește și interpretează declarațiile numeste command-line interpreter sau shell(UNIX). Funcția lui este simplă, să aducă noua comandă și să o execute.

3.2 Serviciile unui SO

Un SO facilitează programarea prin aceste servicii.

Execuția programului: Capabilitatea sistemului de a încărca programul în memorie și să îl execute.

Operații I/O: Din moment ce programele user nu pot executa operații de I/O în mod direct, SO trebuie să ofere mijloace pentru a putea executa operații I/O

Manipularea sistemului de fișiere: Programul trebuie să poată citi, scrie, crea și șterge fișiere

Comunicatie: Schimbul de informație între procese care se execută fie pe același procesor sau pe sisteme diferite conectate prin rețea. Poate fi implementat fie prin memorie partajată sau prin trimiterea de mesaje.

Detectia erorilor: Sa ofere o executie corecta prin detectia erorilor in CPU, in hardware-ul memoriei, in dispozitivele I/O sau in programele user.

Exista si functii aditionale nu insa pt a ajuta user-ul ci pt a asigura operatii sistem eficiente.

Alocarea resurselor: Alocare resurselor la mai multi useri sau la mai multe job-uri care ruleaza simultan.

Raportul resurselor: Tine evidenta care utilizatori folosesc cate si care resurse. Pentru statistica sau intocmirea unei note de plata.

Protectie: Se asigura ca accesul ala resursele sistemului este controlat.

3.3 Apeluri de sistem

Asigura interfata dintre un proces si SO.

Sunt disponibile ca instructiuni in limbaj de asamblare. Limbajele de programare care inlocuiesc permit ca apelurile sistem sa fie facute in mod direct(ex., C, C++).

Sunt folosite 3 metode sa se transmita parametrii de la un program(care ruleaza) la SO:

- Prin registrii
- Stocarea parametrilor intr-o tabele de mamorie, si adresa tabelei este pasata ca un parametru intr-un registru.
- Parametrii se pun in stiva de catre program si se face pop de catre SO.

pag 65 poza

Tipuri de apeluri sistem: controlul proceselor, gestiunea fisierelor, managementul dispozitivelor, intretinerea informatiei si comunicare.

3.3.1 Controlul proceselor

Daca e face un apel de sistem sa se termine programul care ruleaza in mod anormal se pune continutul memoriei din momentul aparitiei apelului si se afiseaza un mesaj de eroare. Continutul memoriei poate fi exminat de un debugger.

Cand un proces se termina(fie in mod normal sau abnormal) SO trebuie sa transfere controlul interpretorului de comenzi. Acesta citeste urmatoarea comanda. Intr-un sistem interactiv se continua cu urmatoarea comanda. Intr-un sistem Batch interpretorul opreste tot job-ul shi continua cu urmatorul.

MS-DOS sistem mono-tasking-> interpretorul de comenzi este incarcat in mem cand se deschide calculatoru.(pg 68 poza 3.3 a) Interpretorul de comenzi incarca programul de rulat in mem scriind peste o mare parte a sa. Se ruleaza programul. Cand se termina din cauza unei erori sau a unui apel de sistem, codul de eroae est esalvat. Partea care a mai ramas din interpretorul de comenzi incarca restul de pe disc si face public codul de eroare pt user sau prg urmator. Cu toate ca MS-DOS = sist monotasking acesta are o metoda pt executie concurenta limitata. Un program TSR iasa cu apelul sistem terminate ans stay resident cand este activata o intrerupere. MS-DOS rezerva spatiul ocupat de programul TSR, astfel incat nu este suprescris cand interpretorul de comenzi este reluat.

FreeBSD sistem multitasking -> Cand userul se logheaza la sistem (interpretorul de comenzi) shell-ul este rulat. Pt a starta un nou proces shell-ul executa un apel sistem de tip fork iar pt a executa procesul un apel de tip exec. shell-ul poate astepta sa se execute

procesul sau poate sa il ruleze in background. Daca procesul ruleaza in background nu poate fi accesat direct. Operatiile de I/O sunt facute prin fisiere. Cand procesul a fost rulat, executa un apel de tip exit ca sa se termine.

poza pag 66

3.3.2 Gestiunea fisierelor

User-ul trebuie sa creeze, sterga, citeasca, modifice, etc fisiere. Fiecare apel sistem are nevoie de numele fisierului si poate de unele din attributele sale.

3.3.3 Gestiunea dispozitivelor(device)

Un program in timp ce ruleaza poate avea nevoie de resurse suplimentare sa poate continua(memorie, acces la fisiere, etc). Daca resursele nu sunt disponibile programul trebuie sa astepte pana acestea devin disponibile. Fisierelor pot privite ca dispozitive reale sau virtuale. Daca sistemul are mai multi useri, mai intai trebuie sa "cerem" dispozitivul pentru a asigura o utilizare exclusiva. Cand terminam cu dispozitivul trebuie eliberat.

3.3.4 Mentinerea informatiei

Multe apeluri sistem exista pt a transfera informatie intre user si SO. Apel sistem pt returnarea timpului si a datei curente, marimea memoriei, nr de procese, bla bla.

3.3.5 Comunicarea

Doua modele de comunicare: prin mesaje si prin memorie partajata.

Comunicare prin mesaje(message passing): Mai intai trebuie deschisa o conexiune. Trebuie sa cunoastem numele celui cu care comunicam. Fiecare PC are un nume(host name), fiecare proces are un nume(process name). Pentru aflarea numelor avem nevoie de apeluri sistem. Procesul recipient(care primeste mesajele) trebuie sa isi dea acordul pt ca o comunicare sa aiba loc(apel accept). Procesele care accepta conexiuni sunt programe sistem create special pt asta(executa apelul wait pana cand se face o conexiune) Sursa comunicatiei = client, iar cel care primeste = server. Clientul si serverul schimba mesaje prin apelurile read si write message. La inchiderea conexiunii se executa apelul close connection.

Comunicare prin memorie partajata(shared memory): Procesele folosesc apeluri sistem pt a avea acces la regiuni de memorie ale celui alt proces. SO incearca sa previna ca un proces sa acceseze memorie celui alt proces. Comunicarea prin memorie partajata are nevoie ca mai multe procese sa fie de acord sa inlature aceasta restrictie. Procesele sunt responsabile ca sa nu scrie in aceeasi locatie simultan.

pg 72 poza 3.5

3.4 Programe sistem

Ofera un mediu convenabil pt dezvoltarea si executia programelor. Programele pot fi impartite in: Gestiunea fisierelor, Informatia starii sistemului, modificare de fisiere, suport pt limbaje de programare, incarcarea si executia programelor, Comunicatie. Cei mai multi utilizatori vad SO ca fiind definit de programele sistem si nu de apelurile de sistem.

3.5 Structura sistemului

Un sistem asa de mare si complex ca un SO trebuie sa fie gestionat cu grija ca sa functioneze corect si sa fie modificat cu usurinta. Sistemul este partitionat in componente mici.

3.5.1 Structura simpla

Multe sisteme comerciale nu au o structura bine definita. Multe SO au pornit ca un sistem simplu si limitat si au fost apoi dezvoltate.

MS-DOS: a fost scris sa ofere cea mai buna functionare in cel mai mic spatiu. Nu este divizat in module si desi are o anumita structura, interfata sa si nivelele de functionare nu sunt indeajuns separate. pag 75 poza 3.6 structura MS-DOS

UNIX: Initial era limitat de functionarea hardware-ului iar UNIX avea o structura limitata. Este format din 2 parti distincte: Programele sistem si kernel-ul. Kernel-ul este si el divizat intr-o serie de interfete si drivere de dispozitive. Putem privi UNIX ca un SO structurat(layered) pag76 poza.

Totul sub interfata apelurilor sistem si peste hardware este kernel. Kernelul furnizeaza sistemul de fisiere, CPU scheduling, gestiunea memoriei, si alte functii ale SO; un numar mare de functii pt un nivel. Acest lucru face UNIX sa fie dificil de dezvoltat, deoarece schimbarile dintr-o sectiune pot avea efecte adverse asupra alte zone.

3.5.2 Structura stratificata

Modularizarea sistemelor poate fi facuta in mai multe moduri. Una este stratificarea unde SO este impartit intr-un numar de straturi sau nivele. Fiecare nivel este construit desupra celui alt. Nivelul de jos (nivelul 0) este hardware-ul iar nivelul cel mai de sus este interfata cu utilizatorul. Principalul avantaj al structurii pe nivele este modularitatea.

Nivelele sunt selectate in asa fel incat fiecare foloseste doar functiile(operatiile) si serviciile nivelelor de jos. Aceasta abordare simplifica verificarea sistemului. Fiecare nivel este implementat doar cu operatiile oferite de nivelul de jos. Un nivel nu trebuie sa stie cum sunt implementate functiile ci doar ce fac. Cea mai mare dificultate cu abordarea stratificata este definirea nivelelor, deoarece un nivel poate folosi doar nivelele de sub el. Alta problema este ca sunt mai ineficiente decat alte modele. Cand un proces executa o operatie I/O executa un apel de sistem care se regaseste in nivelul I/O, care la randul lui apeleaza alte nivele. La fiecare nivel parametrii pot fi modificati sau datele sunt pasate. Trecerea prin fiecare nivel intarzie apelul => un apel sistem dureaza mai mult decat pe un sistem fara straturi.

Sistemul OS/2 este un descendent al MS-DOS cu multitasking. Sistemul a fost implementat intr-un mod stratificat. pag 78 poza

3.5.3 Microkerneluri

Cu dezvoltarea sistemului UNIX kernelul a devenit tot mai mare și dificil de întreținut. S-a dezvoltat un kernel mai mic -> microkernel. Aceasta abordare scoate toate componentele neesențiale din kernel și le implementează în programele user și sistem. Rezultă un kernel mai mic. Funcția principală a kernelului este să asigure comunicarea între programul client și diferitele servicii care rulează în spațiul utilizator. Comunicarea se realizează prin trimiterile de mesaje. Avantajul abordării microkernel este ușurința cu care poate fi extins SO. SO poate fi adaptat mai ușor unor noi arhitecturi. Microkernelul oferă mai multă securitate și siguranță, deoarece cele mai multe servicii sunt rulate ca procese user. Dacă un serviciu se deteriorează nu influențează restul SO.

3.6 Masini virtuale

O implementare stratificată este dusă la un nivel logic cu conceptul de mașină virtuală. Schedulingul CPU-ului poate crea impresia că fiecare user are propriul procesor. Spooling și sistemul de fișiere poate furniza card reader și imprimante virtuale. Un sistem normal cu divizare de timp (time sharing) poate servi ca consolă operatorului mașinii virtuale.

3.6.1 Implementare

Deși conceptul de mașină virtuală este folosit, se implementează greu. Mașina care stă la bază are 2 moduri de funcționare: modul user și modul monitor. Software-ul mașinii virtuale poate rula în mod monitor, deoarece este SO. Mașina virtuală în sine poate rula doar în mod user. Un transfer de la modul user la modul monitor la o mașină reală trebuie să cauzeze un transfer la mașină virtuală. Transferul poate fi făcut relativ simplu. Dacă un apel sistem este făcut de către un program care rulează pe mașină virtuală are loc un transfer la monitorul mașinii virtuale în mașină reală. Când monitorul mașinii virtuale preia controlul poate modifica registrii și PC-ul pentru ca mașină virtuală să simuleze efectul apelului de sistem. La sfârșit poate fi reluată rularea mașinii virtuale. Starea mașinii virtuale a fost modificată în conformanță cu apelul. poza pag 81

3.6.2 Avantaje

Conceptul de mașină virtuală furnizează protecție totală a resurselor sistemului din moment ce fiecare mașină virtuală este izolată de celelalte mașini virtuale. Aceasta izolare însă nu permite o distribuție directă a resurselor.

Un sistem mașină virtuală este un mod optim de a crea și dezvolta SO. Dezvoltarea sistemelor este făcută pe o mașină virtuală și nu întrerupe operațiile de sistem normale. SO controlează toată mașina. La dezvoltare tot sistemul trebuie oprit pentru a putea face modificări și teste. O mașină virtuală poate elimina o mare parte din problemă. Programatorii au propria mașină virtuală și dezvoltarea se face pe mașină virtuală.

3.6.3 Java

Compilerul Java produce bytecodes(coduri in byte) cu un aspect neutral fata de aritectura. Acestea sunt rulate apoi pe Java Virtual Machine(JVM)
JVM consta din: class loader, class verifier si un interpretor runtime. compilatoare Just-In-Time cresc performanta.
poza pag 85

3.7 Designul si implementarea sistemului

3.7.1 Tinta designului

La nivelul ce mai de sus modelul sistemului va fi infuentat de alegerea hardware-ului si tipul sistemului: batch, time-shared, etc.

Cerintele de design pot fi impartite in:

- scopul userilor: So sa fie usor de utilizat, usor de invatat, de incredere, sigure si rapide
- scopul sietm: SO sa fie usor de modelat, implementat si intretinut precum si flexibile, fara erori si eficeinte.

3.7.2 Mecanisme si politici

Mecanismele determina cum este implementat ceva, politicile decit ce va fi implementat. Separarea macanismelor de politici este un principiu f. important, ofera flexibilitate maxima daca deciziile referitoare la politici sunt schimbate ulterior.

3.7.3 Implementarea sistemului

Dupa ce un SO a fost modelat el trebuie implementat.

Traditional ele erau scrise in limbaj de asamblare. Acuma ele sutn scrise in limbaje de nivel inalt. Codul scris in limbaj de nivel inalt poate fi scris mai repede, este mai compact si este mai usor de inteles si debugenit. Un SO poate fi mutat pe ul alt hardware mult mai usor daca este scris intr-un limbaj de nivel inalt.

3.8 System Generation(SYSGEN)

Un SO este facut sa poate rula pe orice tip de masina, sistemul trebuie configurat pt fiecare calculator. Programul SYSGEN obtine informatia referitoare la cofigurarea sistemului hardware. Booting = pornirea calculatorului prin incarcarea kernelului. Bootstrap program = cod stocat in ROM care este in stare sa localizeze kenelul, sa il incarce in memorie si sa ii starteze executia.

Capitolul 4 PROCESE

Un sistem este compus dintr-o colectie de procese.

4.1 Conceptul de Proces

Un SO executa tipuri diferite de programe: sisteme batch executa job-uri iar sisteme time shared executa programe user sau task-uri.

4.1.1 Procesul

Procesul este un program care se afla in executie. Executia proceselor trebuie procesata in mod secvential. Un proces are inclus un program counter, o stiva(care patreaza datele temporare) si o sectiune de date(unde sunt pastrate varibilele globale).

4.1.2 Starea unui proces

In timp ce un proces se executa el isi schimba starea. Fiecare proces se poate afla in una din starile:

new: Procesul este creat.

running: sunt executate instructiuni.

waiting: procesul asteapta sa intervina un eveniment.

ready: procesul asteapta sa fie atribuit unui procesor.

terminated: procesul a terminat executia

poza 4.1 pg 97

4.1.3 Blocul de control al proceselor(PCB)

fiecare proces din SO este reprezentat de un bloc de control al procesului.

poza 4.2 pg 97

Acesta contine informatii asociate unui anumit proces:

Starea unui proces

Program counter

Registrii CPU

Informatia de scheduling a CPU-ului

Informatie asupra gestiunii memoriei

Informatie resurselor

Informatie referitoare la I/O

4.2 Schedulingul proceselor

Ideea este sa se ruleze mai multe procese simultan pt a marii utilizarea CPU.

4.2.1 Cozi de scheduling

Cand un proces este introdus in sistem el este pus in coada job(job queue). Aceasta coada contine toate procesele din sistem. Procesele rezidente in memorie care sunt ready si asteapta sa fie executate sunt pastrate intr-o lista nimita coada ready(ready queue).(pag

100 poza). Lista proceselor care asteapta dupa un dispozitiv I/O este numita coada dispozitivelor(device queue). Fiecare dispozitiv are proprie coada. Reprezentarea comuna a schedulingului proceselor este o diagrama, queueing diagram(pag 102 poza). (Patratele sunt cozi. Sunt reprezentate 2 tipuri de cozi ready si I/O)

Un nou proces este pus initial in coada ready unde sta pana este selectat sa fie executat.

Odata ce procesul este pasat procesorului si se executa avem cazurile:

- procesul poate avea o cerere I/O, este pasat atunci cozii I/O
- procesul poate crea un subproces si va astepta terminarea acestuia
- procesul poate fi inlaturat fortat de la CPU(interrupere) si este pus inapoi in coada ready

In primele 2 cazuri procesul isi schimba la candva starea din waiting in ready si este pus inapoi in coada ready.

4.2.2 Scheduling

Deoarece un proces migraza intre cozile de scheduling SO trebuie sa poata selecta procese din aceste cozi in scopul schedulingului. Procesul de selectie este realizat de scheduling.

Scheduler pe termen lung(long term): job scheduler, selecteaza care proces va fi adus in coada ready din memorie. Se executa cu o frecventa mica. Poate fi si mai incet. El controleaza gradul de multiprogramare(numarul proceselor din memorie). Este invocat doar cand un proces este scos din sistem.

Scheduler pe termen scurt(short term): CPU scheduler, selecteaza unul din procesele gata de executie si alocă unul din ele CPU-ului. Este executat cu o frecventa mare => trebuie sa fie rapid.

Procesele pot fi descrise ca procese legate de I/O sau procese legate de CPU.

Procesele legate de I/O petrece mai mult timp facand o operatie de I/O decat calculand.

Iar un proces legat de CPU petrece mai putin timp executand I/O. Schedulerul pe termen lung trebuie sa faca o combinatie buna intre cele 2 tipuri de procese. Daca toate procesele sunt I/O coada ready va fi mai tot timpul goala iar schedulerul pe termen scurt nu va avea de lucru. Daca procesele vor i CPU coada I/O va fi goala iar sistemul nu va fi echilibrat.

pag 102 poza

4.2.3 Comutarea de context

Cand procesorul comuta de la un proces la altul starea fostului proces trebuie salvata si incarcata starea noului proces. Comutarea de proces este o "pierdere de timp", deoarece sistemul nu face lucru util in timp ce comuta. Timpul de comutare sunt dependenti de suportul hardware. Cu cat este mai complex SO cu atata trebuie depus mai mult lucru la comutare.

4.3 Operatii asupra proceselor

Procesele din sistem se executa concurent si trebuie create si distruse dinamic.

4.3.1 Crearea proceselor

Procesele parinte creaza procese fiu care la arndul lor pot creea alte procese. Se formeaza astfel un arbore de procese(poza pag 104).

Exista mai multe cazuri de partajare a resurselor: parintele si fiul impart toate resursele, fii impart un subset din resursele parintelui, parintii si fii nu impart nici o resursa. Cand un proces parinte creaza altu avem 2 cazuri: Parintele isi continua executia in paralele cu fii sau parintle asteapta pana unii din fii sai au terminat. Referitor la spatiului de adrese avem 2 cazuri: procesul fiu este un duplicat al procesului parinte sau procesul fiu are incarcat in el en program.

Exemplu UNIX: un nou proces este creat cu apelul de sistem **fork**, **exec** este folosit dupa fork ca sa se inlocuiasca spatiul de memorie a procesului cu un nou program

4.3.2 Terminarea proceselor

Un proces se termina daca si-a terminat executia si ii cere SO sa fie sters(apelul **exit**). Procesul poate returna data procesului parinte(prin apelul **wait**). Toate resursele procesului sunt dealocate de SO. Terminarea unui proces mai apare shi in cazurile cand un proces cauzeaza terminarea unui alt proces. De obicei doar procesul parinte poate termina un proces fiu. Parintele poate termina(apelul **abort**) executia proceselor fiu din motivele:

- fiul a depasit utilizarea resurselor alocate lui
- taskul ascoiat fiului nu mai este necesar
- parintele si-a terminat executia si SO nu permite continuarea executiei unui fiu daca parintele a terminat. Daca unproces termina toti fii trebuie sa se termine -> terminare cascadata

4.4 Procese cooperante

Procese independente nu se pot afecta sau nu pot fi afectate de alte procese. Un proces cooperant poate afecta si poate fi afectat de celelalte procese.

Avantajele cooperari proceselor: Distribuita informatiei, marirea vitezei de computatie, modularitate, convenienta. Executia concurenta trebuie sa ofere mecanisme de comunicare si sincronizare intre procese. Un proces proucator produce informatie care este consumate de un proces consumator. Ca cele 2 procese sa ruleze simultan avem nevoie de buffere.

Unbound buffer nu pune nici o restructie de marime a bufferului. Consumatorii trebuie sa astepte noile obiecte dar producatorul poate produce tot timpul noi obiecte.

Bounded buffer: marimea buffer trebuie sa fie fixa. Consumatorul trabuie sa asepte daca bufferul este iar producatorul trebuie sa astepte daca bufferul este plic.

4.5 Comunicatia intre procese(IPC)

Este o alta metoda de a obtine rezultate identice cu procesele cooperante.

4.5.1 Sistemul de pasare a mesajelor

Trebuie sa furnizeze un mecanism de comunicare si sincronizare intre procese. Procesele comunica fara a fi nevoie sa aiba variabile comune. IPC furnizeaza 2 operatii: send (mesaj) – mar mesajului este fixa si receive(mesaj).

Dac a2 procese doresc sa comunice, ele trebuie sa stabileasca o conexiune intre ele si sa schimbe mesaje intre ele. Linia de comunicatie poate fi implementata fizic sau logic. linia de comu

Acesta poate fi direct sau indirect, simetric sau asimetric.

4.5.2 Numele

Procesele care doresc sa cominica trebuie sa stie sum se refere pe celalalt.

4.5.2.1 Comunicatie directa

Fiecare proces care vrea sa comunice trebuie sa numeasca numele procesului cu care comunica:

send (*P, message*) – send a message to process P

receive(*Q, message*) – receive a message from process Q

O linie de comunicatie are proprietatile:

- Linia este stabilita in mod automat intre fiecare pereche de comunicare. procesele trebuie sa cunoasca identitatea celuiilalu.

- O legatura este asociata cu fix 2 procese

- Exact o legatura exista intre o prereche.

- Legatura poate fi unidirectionala dar de obicei este bidirectionala.

4.5.2.2 Comunicatie indirecta

Mesajele sunt directionate si primite de la un mailbox (port). Fiecare port are un ID unic. Procesele pot comunica doar daca au un mailbox comun. Prmitivele sent si receive sunt de felul urmator:

send(*A, message*) – send a message to mailbox A

receive(*A, message*) – receive a message from mailbox A

Un mailbox poate avea ca proprietar SO sau un proces. Mailboxurile cu proprietar un proces trebuie sa distingem intre proprietar si user. Cand un proces care are un mailbox termina, mailboxu dispare. Toate procesele care comunicau cu mailboxu trebuie instiintate ca acesta nu mai exista.

Proprietatile liniei de comunicatie:

- Legatura este facuta doar daca procesele au un mailbox comun.

- O legatura poate fi asoiata mai multor procese.

- Fiecare pereche de procese poate avea mai multe legaturi de comunicatie.

- Legatura poate fi uni- sau bidirectionala.

Impartirea mesajelor. Cand procesele P1, P2 si P3 impart mailboxu A iar P1 trimite un mesaj si P2 si P3 il primesc. Problema este cine intra in posesia mesajului. Solutiile la aceasta problema ar putea fi:

- Doar doua procese sa fie asociate cu o legatura de comunicatie.

- Doar un proces sa poata executa operatia receive.

Permiterea sa se aleaga destinatatul in mod arbitrar de catre sistem. expeditorul este instiintat cine a fost destinatarul.

4.5.3 Sincronizare

Trimiterea mesajelor poate fi blocanta sau nebloanta adica sincrona sau asincrona.

Trimitere blocanta: procesul expeditor este blocat pana cand mesajul este primit de proces sau mailbox.

Trimitere nebloanta: Procesul trimite mesajul si isi continua executia.

Receptionare blocanta: Destinatarul se blocheaza pana cand un mesaj este disponibil.

Receptionare nebloanta: Destinatarul primeste fie un mesaj valid fie unul nul.

Sunt psibile de combinatii diferite de send si recieve. Cand si send si recieve sunt blocante avem un rendezvous intre expeditor si destinatar.

4.5.4 Buffering

La comunicatie mesajele schimbate de procesele de comunicare sunt pastrate in cozi temporare.

Avem 3 cazuri de implementare unei asemenea cozi:

Capacitate zero: Coadă are lungimea de 0. In acest caz expeditorul trebuie sa blocheze pana cand este primit mesajul. Aceasta implementare este privita ca una fara buffer.

Capacitate limitata: Coadă are lungime finita n. Cel mult n mesaje pot fi stocate in ea. Daca linia de comunicatie este plina expeditorul trebuie sa se blocheze pana cand se face loc in coada.

Capacitate infinita: Coadă are lungime infinita. Expeditorul nu trebuie sa astepte.

4.6 Comunicatie in sisteme client server

4.6.1 Socketuri

Un socket este un capat al comunicatiei. O pereche de procese care comunica are nevoie de doua soketuri. Se foloseste in general o arhitectura client server. O conexiune trebuie sa fie unica. Comunicarea se realizeaza concatenand adresa IP cu portul(socketul

161.25.19.8:162 se refera la port **1625** host **161.25.19.8**).

Trei tipuri diferite de socketuri: - Orientate pe conexiune, socketi faca conexiune(Datagram sockets) si multicast socket(poate transmite la mai multi destinatari).
poza pag 118

4.6.2 Proceduri cu apelul la distanta(RPC)

RPC este un apel de procedura intre procesele dintr-o retea. Semantica RPC permite invocarea unei proceduri la distanta in acelasi mod ca una locala. Sistemul RPC ascunde detaliile necesare permitand ca comunicarea sa aiba loc. Acest lucru este posibil prin furnizarea unui stub pe partea clientului. Un stub separat exista pt fiecare procedura la distanta. Cand o procedura la distanta este apelata, RPC apeleaza stub-ul potrivit pasand parametrii trimisi catre procedura. Acest stub localizeaza portul pe server si **marshalls**

parametrii. Acest proces include inpachetarea parametrilor intr-o forma in care pot fi transmisi prin retea. Stub-ul transmite apoi un mesaj serverului. Un stub similar al serverului primeste mesajul si apeleaza procedura. Daca este necesar valorile rezultate sunt trimise inapoi clientului intr-un mod asemanator. Sistemele RPC definesc o reprezentare a datelor independente de masina: reprezentare externa a datelor(XDR). In partea clientului marshalling presupune convertirea datelor dependente de masina in XDR inainte sa fie trasmise serverului. Pe partea serverului, datele in reprezentare XDR sunt transformate in reprezentare locala a serverului. poza pg 123

4.6.3 Invocarea metodelor la distanta(RMI)

poza pg 124

RMI este un mecanism Java similar cu RPC.

RMI permite programului Java sa invoce o metoda al unui obiect la distanta. Un stub este un proxy pentru obiectul la distanta si apartine de client. Cand un client invoca o metoda la distanta, stub-ul este apelat. Stub-ul client este responsabil cu crearea unui pachet(parcel) care consta din numele metodei si parametrii metodei. Acest pachet este trimis la server unde un schelet al obiectelor la distanta il primeste. Acesta cheama metoda dorita. Rezultatul este apoi trimis intr-un pachet la client. poza pg 125

THREAD-uri

5.1 Introducere

Thread = numit si **lightweight process(LWP)**, este unitatea de baza a utilizarii CPU ; contine un ID al thread-ului, un numarator de program (PC), un set de registri si o stiva. Imparte cu thread-urile apartinatoare aceluiasi proces: segmental de cod, de date si alte resurse de operare a sist (de ex. : fisiere deschise si semnale).

Un proces normal are un singur thread de control. Daca procesul are mai multe thread-uri de control atunci el poate realiza mai multe task-uri deodata.

5.1.1. Motivatii

Multe programe care ruleaza pe un calculator modern sunt multithread. O astfel de aplicatie este implementata cu mai multe thread-uri de control. In anumite situatii unei singure aplicatii i se cere sa efectueze mai multe task-uri similare. O solutie ar fi sa ruleze ca un proces care accepta cereri. Aceasta solutie era de actualitate inainte ca thread-urile sa devina populare. Insa crearea de procese este destul de dificila, de aceea este mai simplu ca un proces care contine mai multe thread-uri sa serveasca aceluiasi scop, de ex. daca a fost

facuta o solicitare in loc sa creeze un alt proces creaza un alt thread care va raspunde cererii.

Thread-urile joaca un rol semnificativ in sistemele RPC (remote procedure call). In mod normal serverele RPC sunt multithread. Daca serverul primeste un mesaj, ii raspunde folosind un alt thread. Acest lucru permite serverului sa raspunda in acelasi timp la mai multe cereri.

5.1.2 Beneficii

Reactie : cu ajutorul multithreading o aplicatie permite unui program sa functioneze chiar daca o parte din ea este blocata sau executa operatie lunga, crescand astfel reactia.

Impartirea resurselor : in mod prestabilit thread-urile isi impart memoria si resursele procesului de care apartin. Beneficiul: o aplicatie poate avea mai multe thread-uri ale activitatii la aceiasi adresa de memorie

Economie : alocarea memoriei si a resurselor pentru crearea unui proces este costisitoare. De aceea e mai economic sa creezi si sa comuti contextual thread-uri deoarece ele impart resursele procesului.

Utilizarea arhitecturii multiprocesor : beneficiile multithreading cresc semnificativ atunci cand avem o arhitectura multiprocessor, unde fiecare thread poate rula in parallel pe procesoare diferite.

5.1.3 Thread-uri user si Kernel

Thread-uri user: sunt suportate deasupra thread-urile Kernel si sunt implemetate de library thread la nivel user. Libraria ofera posibilitatea crearii, aranjarii si administrarii thread-urilor fara ajutorul Kernel, de aceea sunt usor de folosit si administrat

Thread-uri Kernel: sunt suportate direct de catre sistemul de operare: realizeaza creari, administrari si aranjari ale thread-urilor directin spatial Kernel; din aceasta cauza thread-urile Kernel necesita mai mult timp pentru creare si administrare decat thread-urile user.

5.2 Modele Multithreading

Multe modele suporta atat thread-uri user cat si Kernel.

5.2.1 Modelul Many-to-One

Acest model mapeaza (insumeaza) thread-urile user la un singur thread Kernel. Administrarea thread-ului se face in spatiul user deci este efficient, dar intreg procesul se va bloca daca un thread face un apel de bllocare a sistemului. Deoarece doar un user thread poate accesa thread-ul Kernel la un moment, thread-urile multiple nu pot rula pe multiprocesoare.

5.2.2 modelul One-to-One

Mapeaza pentru fiecare thread user un thread Kernel. Asigura o mai multa concurenta decat modelul de dinainte deoarece daca un thread user solicita un apel de blocare a sistemului alte thread-uri pot functiona in continuare. Singura deficianta este ca creand un thread user trebuie creat si thread-ul Kernel corespunzator, acest lucru ingreutand performantele aplicatiei.

5.2.3 Modelul Many-to-Many

Multiplexeaza mai multe thread-uri de nivel user cu o cantitate mai mica sau egala de thread-uri de nivel Kernel. Acest model permite programatorului sa creeze cate thread-uri vrea, iar concurenta nu este castigata deoarece Kernel poate aranja doar cate un thread pe rand. Acest model permite cresterea concurentei dar programatorul trebuie sa aiba grija sa nu creeze multe thread-uri intr-o aplicatie.

5.3 Restrictii de threading

In acest cap. vom discuta despre restrictiile ce apar in programele multithreading.

5.3.1 fork and exec System call

Spre deosebire de cap 4 semantica fork si exec System call se modifica. Daca un thread dintr-un program face o cerere fork, se dubleaza thread-urile sau este noul proces single-threaded? Unele sisteme UNIX au 2 versiuni de fork unul care dubleaza toate thread-urile si unul care dubleaza doar thread-ul care a cerut fork. Apelul exec functioneaza in acelasi mod descries in cap 4. inseamna ca daca un thread cere exec, programul specificat ca si parametru inlocuieste intreg procesul inclusive toate thread-urile si LWP-urile.

Utilizare celor doua versiuni de fork depinde de aplicatie. Daca exec este cerut imediat dupa fork atunci nu mai este necesar dublarea fiecarui thread, deoarece programul dat ca parametru va inlocui procesul, cel mai bine este dublarea thread-ului care a facut cererea.

5.3.2 Anularea

Anularea threadurilor reprezinta terminarea unui thread inainte ca acesta sa se fi executat complet. Thread-ul anulat este deseori referit ca thread-ul tinta.

Anulare asincrona: un thread anuleaza imediat threadul tinta

Anulare cu intarziere: Threadul tinta poate sa verifice daca ar trebui sa termina, oferind modalitatea threadului tinta se se anuleze in mod ordonat.

Dificultate la anulare: Cand resurse sunt alocate unui thread anulat sau daca threadul a fost anulat in timp ce actualiza date ce le imparte cu alte threaduri. Aceasta problema apare la anulare asincrona. SO pune stapanire pe resurse dar de multe ori nu pe toate. Astfel anularea asincrona nu va elibera resursele necesare. La anularea cu intarziere anularea va avea loc doar cand threadul tinta verifica acest lucru. Anularea are loc la un moment cand ea este sigura. Pthreadurile numeste astfel de puncte puncte de anulare.

5.3.3 Manipularea semnalelor

La UNIX un semnal este folosit sa instiinteze un proces ca a avut loc un anumit event. toate semnalele urmaresc acelasi model:

- Un semnal este generat la aparitia unui anumit event.
- Semnalul generat este transmis la proces.
- Dupa ce a fost transmis, semnalul trebuie tratat.

Semnale sincrone: cand avem divison by 0, acces ilegal la mem...Semnalele sincrone sunt trimise aceleiasi proces care a executat operatia care a generat semnalul.

Semnale asincrone: sunt generate de evenuri externe unui proces care ruleaza. De obicei semnalul asincron este trimis unui alt proces.

Fiecare semnal trebuie sa fie tratat de unul din handlers: hanlerul default sau handlerul definit de user. Fiecare semnal are un handler de semnal default care este rulat de kernel cand tratam semnalul. Actiunea default poate fi suprascrisa de un handler de semnal definit de utilizator.

Transmiterea semnalelor este mai dificila la programele multithreading. In general exista urmatoarele optiuni unde sa se trimita semnalul:

- trimiterea semnalului la threadul la care se aplica semnalul
- trimiterea semnalului la toate threadurile din proces
- trimiterea semnalului la anumite threaduri din proces
- desemnarea unui thread special care va primi toate semnalele din proces

Versiuni multithreading ale lui UNIX permit threadului sa specifice ce semnale va accepta si care le va bloca.

Desi win 200 nu are suport pt semnale, ele pot fi emulate folosind apeluri asincrone de procedura(APC). APC permite unui thread user sa specifice functia care va fi apelata cand threadul este anuntat de un event pariticular.

5.3.4 Grupuri de threaduri

Cand un server primeste o cerere el creeaza un thread care sa rezolve cererea.

Dezavantaje: Timpul pt crearea unui thread este destul de mare, threadul va vi sters indata ce iti va termina lucrul, prin existanta simultana a unui numar nelimitat de threaduri resursele sistemului se epuizeaza.

Solutia: Folosirea unui grup de threaduri(thread pool)

La inceperea procesului sunt create un nr de threaduri care sunt puse in acest grup. Cand serverul primeste o cerere el "trezeste" un thread. Cand thread-ul termina el este pus din nou in grup in starea de asteptare.

Avantajele thread pool: o cerere poate fi tratata mult mai repede, numarul de threaduri este limitat-> nu apare o supraincarcare a sistemului.

5.3.5 Date specifice threadului

Threadurile unui proces impart datele procesului. Este posibil ca threadul sa aiba nevoie de o copie proprie de date. Acest tip de date se numesc date specifice threadului.

5.4 Pthreaduri

Pthread se refera la standardul POSIX (IEEE 1003.1c) si defineste un API care se refera la crearea sincronizarea threadurilor. API specifica comportamentul threadului din librarie, implementarea depinde de dezvoltarea librariei. Pthreads sunt deseori regasite in SO UNIX.

5.5 Thread-uri SOLARIS 2

Capitolul 6 (CPU SCHEDULING)

6.1 Concepte de bază

Ideea de bază este că fiecare proces este planificat înainte de a fi executat pentru ca procesorul să nu stea degeaba când procesul în execuție trece în starea de *wait*.

6.1.1 CPU – I/O Burst Cycle

Procesele au următoarea prioritate: execuția lor constă într-un ciclu de execuție CPU și apoi un I/O *wait*. Procesele încep cu o izbucnire (burst) de CPU urmată de o izbucnire (burst) de I/O *wait* după care cer iară CPU și tot așa. (a se urmări și diagramele de la pag. 152-153)

6.1.2 CPU Scheduler

Atunci când procesorul devine nefolosit (idle), SO trebuie să aleagă un proces din coada de *ready*. Acest lucru este făcut de **short-term scheduler** (planificator de termen scurt, scrie de asta și prin cap 3 parca). Coada de ready nu este neapărat de tip FIFO ci depinde de algoritmul fiecărui planificator.

6.1.3 Preemptive Scheduling

Deciziile de planificare pentru execuția CPU sunt luate în următoarele circumstanțe:

1. Când un proces trece din *running* în *waiting* (cerere de I/O)
2. Când un proces trece din *running* în *ready* (apare o întrerupere)
3. Când un proces trece din *waiting* în *ready* (terminarea unei cereri de I/O)
4. Când se termină un proces

În cazurile 1 și 4 nu există termen de planificare (trebuie ales un proces din coada de ready). Când planificarea are loc în aceste cazuri se spune ca ea este nepreemptivă (**nonpreemptive**), altfel planificarea este preemptivă (**preemptive**). Sub planificare nepreemptivă, odată ce CPU a fost alocat unui proces, acesta deține CPU până când se termină sau până când trece în *waiting*. Planificarea preemptivă are un dezavantaj (cost) și anume: în momentul în care un proces este în mijlocul unui update de date și este scos de un alt proces (prin preempțiune) care vrea să citească acele date care nu sunt completate trebuiesc folosite mecanisme care fac coordonarea acceselor la datele comune (shared). Preempțiunea are efect și în design-ul kernel-ului (în timpul unei procesări de apel de sistem, kernel-ul poate fi ocupat cu un proces (ex schimbarea datelor importante de kernel, cum ar fi cozile I/O). Dacă un proces este preemptat în mijlocul acestor schimbări și kernel-ul trebuie să modifice aceeași structură ar fi haos).

Cocluzie: Un sistem este preemptiv atunci când SO poate decide să oprească un proces (face planificare împotriva voinței procesului).

6.1.4 Dispatcher (expeditorul??)

El este modulul care predă controlul CPU unui proces selectat de short-term scheduler (traducerea e mai sus) și trebuie să: schimbe contextul, schimbe utilizatorul (user mode), sară la locația adecvată a programului utilizator pentru a-l restarta. Timpul care îi ia dispatcher-ului să oprească un proces și să pornească altul se numește **dispatch latency**.

6.2 Scheduling Criteria

Caracteristicile folosite la comparația diferiților algoritmi de planificare pot face diferențe substanțiale în determinarea celui mai bun algoritm. Criteriile includ următoarele:

Utilizarea CPU (acesta trebuie să fie cât mai ocupat posibil; într-un sistem real utilizarea lui este de la 40% până la 90%)

Throughput (nr-ul proceselor completate per unitate de timp)

Turnaround time (timpul total de execuție al unui proces; suma timpilor de așteptare pentru a intra în memorie, *waiting* în coada de ready, execuție CPU și I/O)

Waiting time (timpul de așteptare în coada de ready)

Response time (timpul de cerere a unui răspuns până când primul răspuns este primit).

Este mai important a se cunoaște timpii de răspuns a cererilor (sau micșorarea variațiilor timpilor de răspuns) decât micșorarea timpului mediu de răspuns deoarece astfel se pot face predicționările mai corect.

6.3 Scheduling Algorithms

Trebuie ales ce proces să intre în execuție din coada de *ready*.

6.3.1 Planificarea Primul venit, Primul servit (FCFS – first come, first served)

Este de departe cel mai simplu algoritm; se folosește o coada de ready de tip FIFO.

Este dezavantajoasă ca timp mediu de așteptare deoarece unele procese pot dura foarte mult iar altele foarte puțin. Explicația, mai greu de reprodus, este la pagina 157.

Algoritmul FCFS este nepreemptiv, odata ce a fost CPU alocat unui proces, acesta îl ține ocupat până când se termină sau are cerere I/O.

6.3.2 Cea mai scurtă treabă, primul la execuție (SJF – Shortest job-first)

Algoritmul asociază fiecărui proces lungimea următoarei izbucniri CPU (CPU burst) și îl alege totdeauna pe cel cu cea mai mică izbucnire. Dacă 2 procese au aceeași izbucnire (sau mai bine zis durată de execuție CPU) se folosește algoritmul FCFS. Deoarece este dificil de știut durata cererii următoare de CPU este frecvent folosit la planificarea pe termen lung (long-term scheduling). Acest algoritm poate fi preemptiv sau nepreemptiv. Dacă este preemptiv, se mai numește și **shortest-remaining-time-first scheduling**. (explicația cu calcule este la pag 159-160).

6.3.3 Planificarea prioritară (Priority Scheduling)

Algoritmul precedent este un caz special al organizării prioritare. O prioritate este asociată fiecărui proces și i se alocă CPU procesului cu prioritatea cea mai mare.

Prioritățile pot fi definite intern (timp limită, cereri memorie, rata medie I/O și CPU burst, etc) și extern sistemului de operare (importanța procesului,...).

- poate fi preemptiv sau nepreemptiv

Problema mare a acestui algoritm este blocarea infinită (**indefinite blocking**). Procesele cu prioritate mică pot ajunge să nu fie servite niciodată. O soluție la această problemă ar fi folosirea vârstei (**aging**) și anume din timp în timp se crește prioritatea procesului din coada de ready cu un număr de unități.

6.3.4 Round-Robin Scheduling (RR)

Este folosită în special la sistemele cu time sharing;

Este similară cu FCFS dar este adăugată preempțiune la schimbarea între procese

Este definită o cantitate mică de timp numită **time quantum** (sau **time slice**), fiecare cuantă având între 10 și 100 de milisecunde. Coada *ready* este tratată ca circulară de tip FIFO iar planificatorul de CPU trece pe la fiecare proces din coadă și îi alocă procesor un interval de o cuantă. (calcule și explicații sunt la pag 163-164).

6.3.5 Multilevel Queue Scheduling

Algoritmul separă coada de *ready* în mai multe cozi în funcție de clasele proceselor (sau proprietățile lor) (ex: foreground și background processes care au cerințe de timpi de răspuns diferite). În plus, trebuie făcută planificare între cozi cel mai des implementată fiind fixed-priority preemptive scheduling. Mai există posibilitatea ca fiecare coadă să primească o anumită cantitate de timp pentru execuție iar avansarea în cozi se face după FIFO.

6.3.6 Multilevel Feedback Queue Scheduling

Diferența între acest algoritm și cel precedent este că procesele pot migra dintr-o coadă în alta. Ideea este de a separa procesele cu durată de execuție CPU diferită. Dacă un proces utilizează prea mult timp CPU, va fi mutat într-o coadă cu prioritate mai mică și invers. Acest algoritm este definit de următorii parametri: numărul cozilor, algoritmul de planificare al fiecărei cozi, metoda folosită pentru promovarea unui proces într-o coadă mai prioritară și invers, metoda folosită pentru intrarea unui nou proces într-o anumită coadă.

6.4 Planificarea multi-procesor (Multiple-Processor Scheduling)

Se presupune că procesoarele sunt identice (omogene) în funcționalitate și se presupune **uniform memory acces(UMA)**. Toate procesele se introduc într-o coadă de unde sunt distribuite la procesoarele disponibile (libere), acest lucru este pentru a evita **load sharing**. Trebuie asigurat să nu aleagă două procesoare același proces și ca procesele să nu se piardă din coadă. Acest lucru se poate face prin folosirea unui procesor ca planificator => ia naștere o structură master-slave (devine multiprocesare asimetrică).

6.5 Planificare în timp real (Real-Time Scheduling)

Există sisteme **hard real-time** și sisteme **soft real-time**. Primele sunt solicitate la terminarea unei sarcini critice într-o cantitate de timp garantată. În general, un proces vine cu o declarație a timpului necesar pentru a îndeplini un I/O. Planificatorul poate admite procesul, garantând că procesul se va termina la timp, sau poate să-l respingă. Acest lucru se cunoaște ca rezervarea resurselor. Aceste garanții presupun ca planificatorul să cunoască exact cât durează fiecare funcție a SO ceea ce este imposibil în sistemele cu memorie virtuală.

Evaluarea **soft real-time** este mai puțin restrictivă; necesită ca procesele critice să primească prioritate peste cele mai puțin norocoase. Implementarea unei funcționalități necesită atenție la design-ul planificatorului și la aspecte legate de SO. SO trebuie să aibă planificare prioritară și procesele în timp real trebuie să aibă prioritatea cea mai mare iar aceasta nu trebuie să scadă în timp.

Pentru a menține latența de dispatch la nivel jos, chemările sistem trebuie să fie preemptibile. Acest lucru poate fi făcut prin crearea de puncte de preempție în cadrul kernel-ului sau prin facerea a tot kernel-ului preemptibil. La ultima variantă poate apărea fenomenul de **priority inversion** atunci când un proces cu prioritate mai mare vrea să citească sau să modifice date din kernel care sunt accesate de către un proces cu prioritate mai mică. Problema poate fi rezolvată cu **priority-inheritance protocol**, în care procesele cu prioritate mai mică care accesează datele unui proces cu prioritate mai mare moștenesc prioritatea mare până când sunt gata cu resursele.

6.6 Evaluarea algoritmilor

Criteriile pentru selectarea unui algoritm pot fi următoarele:

- maximizarea utilizării CPU sub constrângerea că timpul maxim de răspuns este de 1 secundă.
- maximizarea throughput a.î. timpul de execuție al unui proces este în medie dir. prop. cu timpul total de execuție

6.6.1 Modelul determinist (Deterministic Modeling)

Evaluarea analitică folosește un algoritm dat și încărcarea unui sistem (workload) pentru a produce o formulă sau un număr care evaluează performanța algoritmului pentru acea încărcare. Un tip de astfel de evaluare analitică este modelul determinist. Această metodă ia o încărcătură (workload) particulară predeterminată și definește performanțele fiecărui algoritm. (calculare la pag 173-174)

Este simplă și rapidă. Produce numere exacte permițând ca algoritmi să fie comparați. Dezavantajul este că metoda are nevoie cifre exacte ca intrare și de aceea nu prea este folosită.

6.6.2 Modele tip coadă (Queueing Models)

Deoarece procesele care rulează pe multe sisteme variază de la o zi la alta nu se poate folosi modelul determinist. Ce poate fi determinat este distribuția izbucnirilor (burst) CPU și I/O. Sistemul calculatorului este descris ca o rețea de servere. Fiecare server are o coadă de procese în așteptare. CPU este un server cu coada lui de ready precum și sistemul I/O cu cozile dispozitivelor. Cunoscând ratele de sosire și ratele de servire, se pot calcula utilizarea de calcul, media lungimilor cozilor, media de așteptare, etc. Acest lucru se numește **queueing-network analysis**.

Ex: n = media aritmetică a lungimilor cozilor

W = media aritm. a timpului de așteptare în coadă

λ = media ratei de sosire a noilor procese în coadă

Dacă sistemul este într-o stare sigură (steady state), atunci numărul proceselor care părăsesc coada trebuie să fie egal cu numărul proceselor care ajung. $\Rightarrow n = \lambda \times W$ (această ecuație este cunoscută sub numele de **Little's formula**).

Dezavantaje: matematica algoritmilor sau distribuțiilor poate fi dificilă, sosirea și servirea poate fi definită des nerealistă, trebuie făcute anumite presupuneri care nu sunt precise.

6.6.3 Simulări

Pentru a obține evaluări mai precise a algoritmilor de planificare se pot folosi simulări. Se folosesc simulatoare care folosesc variabile reprezentând ceasul care sunt incrementate pe măsură ce simulatorul își schimbă starea. În timp ce se execută simularea, statistici care indică performanța algoritmului sunt adunate și printate.

Totuși, datorită relațiilor între evenimentele succesive dintr-un sistem în timp real, o simulare distribution-driven poate fi imprecisă. Pentru a corecta această problemă se pot folosi **trace tapes** (fig. de la pag 176). Se poate crea un trace tape prin monitorizarea sistemului real, înregistrând secvența evenimentelor actuale.

6.6.4 Implementarea

Până și simularea are acuratețe limitată. Cea mai precisă metodă de evaluare ar fi implementarea practică a algoritmului într-un SO. Dezavantajul acestei abordări este costul. Acesta nu constă doar din scrisul codului algoritmului și modificarea SO pentru a-l suporta ci și reacția utilizatorilor la un SO în continuă schimbare. Mai există dezavantaj și în faptul că odată fiecare algoritm de evaluare se schimbă și mediul în care este folosit algoritmul. Astfel este nevoie de scrierea a noi programe și apariția a probleme diferite.

6.7 Process Scheduling Models

În această secțiune sunt tratate planificarea proceselor în Solaris 2, Windows 2000 și Linux. Librăria de thread-uri organizează user-level threads ca să ruleze pe un LWP (lightweight process) disponibil, schemă cunoscută ca **process local scheduling**, în care planificarea thread-urilor se face local în aplicație. Invers, kernel-ul folosește **system global scheduling** pentru a decide care thread de kernel să fie planificat.

6.7.1 Exemplu: Solaris 2

Planificarea proceselor se face după prioritate. Are patru clase de planificare care sunt, în ordinea priorității: real time, system, time sharing și interactive. Fiecare clasă include priorități diferite și algoritmi de planificare diferiți (a se vedea fig. de la pag 179).

Un proces pornește cu un LWP și poate crea oricâte LWP are nevoie. Fiecare LWP moștenește clasa de planificare și prioritatea de la procesul părinte. Default, clasa de planificare a unui proces este time sharing.

6.7.2 Exemplu: Windows 2000

Windows 2000 planifică thread-uri folosind o planificare preemtivă bazată pe priorități. Planificatorul lui win2000 asigură că thread-ul cu cea mai mare prioritate va fi întotdeauna executat. Porțiunea kernel-ului care se ocupă cu planificarea se numește *dispatcher*. Un thread selectat pentru rulare se va executa până când va fi preemtat de un thread cu prioritate mai mare, până când se termină, până când își termină cuanta sau până când cere I/O.

Dispatcher-ul folosește un plan de priorități pe 32 de nivele pentru a determina ordinea execuției thread-urilor. Prioritățile sunt împărțite în 2 clase: **variable class** (thread-uri cu priorități de la 1 la 15) și **real-time class** (thread-uri cu priorități de la 16 la 31). Când nu este găsit nici un thread în *ready* dispatcher-ul va executa un thread special numit **idle thread**. Există o relație între prioritățile numerice a kernel-ului lui win2000 și Win32 API (carte pag 180-181).

6.7.3 Exemplu: Linux

Linux are doi algoritmi de planificare a proceselor separați. Unul este algoritmul time-sharing (pentru planificare preemtivă corectă între procese multiple) iar celălalt este proiectat pentru cereri în timp real unde prioritățile absolute sunt mai importante decât corectitudinea. Linux permite ca doar procesele care rulează în modul utilizator să fie preemate. Un proces care rulează în modul kernel nu poate fi preemtat chiar dacă este disponibil un proces în timp real cu prioritate mai mare.

Prima clasă de planificare este pentru procesele time-sharing. Linux folosește un algoritm prioritizat, bazat pe credite (**credit-based**). Fiecare proces procesează un anumit număr de credite planificate; când trebuie aleasă o nouă clasă pentru execuție, procesul cu cele mai multe credite este selectat. La fiecare unitate de timp un proces în execuție pierde un credit până ajunge la 0 și un alt proces este selectat. Dacă nici un proces rulabil nu are credite, SO face o operație de recreditare fiecărui proces din sistem după următoarea regulă: $\text{credits} = \text{credits}/2 + \text{priority}$.

Planificarea în timp real a lui Linux implementează 2 clase: FCFS și RR. În ambele cazuri, fiecare proces are adițional la clasa de planificare o prioritate.

Capitolul 8 DeadLocks (Puncte Moarte)

DeadLock- semnificatie:

Intr-un mediu multiprogram, mai multe procese concureaza pt un numar finit de resurse. Daca resursele nu sunt disponibile in acel moment procesul respectiv intra in stare de asteptare. Dar s-ar putea ca acest proces sa stea mereu in stare de asteptare pt ca resursele de care are el nevoie sunt folosite de alte procese. Aceasta situatie se numeste **DeadLock**

8.1 DeadLock-uri in cadrul unui sistem

Un sistem este alcatuit dintr-un nr finit de resurse, care sunt impartite in mai multe tipuri, fiecare tip continand mai multe instante identice. Exemple de tipuri de resurse: spatiu de memorie, cicluri procesor, fisiere, dispozitive I/O(imprimante, tape drivers). Daca sistemul are 2 CPU atunci resursa de tip CPU are 2 instante, la fel resursa de tip imprimanta poate avea 5 instante.

Daca un proces solicita o instanta a unei resurse, alocarea *oricarei* dintre instante ar trebui sa satisfaca cererea. Daca nu se intampla asa atunci instantele nu sunt identice, si clasele pentru tipuri de resurse nu au fost definite corect.

Un proces trebuie sa ceara o resursa inainte de a o folosi si trebuie sa elibereze resursa respectiva dupa ce nu mai are nevoie de ea. Un proces poate face cereri pentru atatea resurse de cate are nevoie, dar in mod evident cererea sa nu poate depasi numarul total de resurse din sistem.

Intr-un mod normal de operare un proces utilizeaza o resursa in urmat secventa:

1. **Request:** daca cererea nu poate fi satisfacuta imediat, procesul intra in starea de asteptare pana cand resursa este eliberata.
2. **Use:** procesul poate folosi resursa.
3. **Release:** procesul elibereaza resursa.

Cererea si Eliberare unei resurse(Request & Release) sunt apeluri sistem(vezi capitolul 3). Exemple: *requeste* and *release device*, *open* and *close file*, *allocate* and *free* memory-system calls. Cererea si eliberarea altor resurse poate fi facuta prin intermediul operatiilor *wait* and *signal* cu ajutorul semafoarelor(cap anterioare). Sistemul are o tabela unde unde este memorata starea fiecarei resurse(free or allocated). Daca resursa solicitata este alocata deja atunci procesul in cauza este adaugat intr-o coada de asteptare(*queue*).

Un set de procese se afla in stare de *deadlock* cand fiecare proces din acel set asteapta dupa un eveniment ce poate fi cauzat doar de catre un proces din cadrul aceluiasi set.

Exemplu: avem un sistem cu 3 tape drivers. Fiecare unitate este manipulata de catre un proces. Daca fiecare proces face o cerere pentru alta unitate atunci cele 3 procese se afla in deadlock.

8.2 Caracterizarea DeadLock-urilor

8.2.1 Conditii necesare:

Un deadlock poate aparea daca sunt indeplinite simultan cele 4 conditii de mai jos:

1. **Excluziune Mutuala:** cel puțin o resursă trebuie ținută într-o stare “non-sharabila”, adică nu poate fi partajată. Un singur proces o poate folosi la un moment dat.
2. **Ocupare și Așteptare(Hold and Wait):** un proces trebuie să ocupe cel puțin o resursă și să aștepte după alte resurse care sunt momentan folosite de alte procese.
3. **Fără Preemptiuni:** resursele nu pot fi eliberate numai în mod voluntar de către procesul care le ținea ocupate, numai după ce acesta și-a terminat “treaba”.
4. **Așteptare circulară:** trebuie să existe un set de n procese $\{P_0, P_1, \dots, P_n\}$ care să aștepte în următorul mod: P_0 așteaptă după o resursă care este ocupată de P_1 , P_1 așteaptă după o resursă care este ocupată de P_2 , ..., P_{n-1} așteaptă după P_n . P_n la rândul lui așteaptă după P_0 .

8.2.2 Graful Alocării Resurselor:

Este un graf prin care deadlock-urile pot fi descrise mai precis. El conține noduri, notate cu V și muchii, notate cu E . Nodurile V sunt împartite în două categorii: $P=\{P_1, P_2, \dots, P_n\}$ care reprezintă toate procesele active din sistem, și $R=\{R_1, R_2, \dots, R_n\}$ care reprezintă toate tipurile de resurse din sistem.

Dacă avem o muchie $P_i \rightarrow R_j$ asta înseamnă că procesul P_i solicită o resursă de tipul R_j , și acum așteaptă după acea resursă. $R_j \rightarrow P_i$ înseamnă că o resursă de tip R_j a fost alocată unui proces P_i . $P_i \rightarrow R_j$: **request edge(muchie de cerere)**, $R_j \rightarrow P_i$: **assignment edge(muchie de atribuire)**.

Grafic fiecare proces este reprezentat printr-un cerc, iar o resursă printr-un patrat. Fiecare instanță a unei resurse este desenată ca un punct în interiorul patratului. Muchiile de cerere pointează doar spre punctul resursei corespunzătoare. Când un proces P_i cere o resursă de tipul R_j este desenată în graf o muchie de cerere. Când această cerere poate fi satisfăcută muchia de cerere desenată este instantaneu transformată într-o muchie de atribuire.

O situație concretă este prezentată la pg 247 fig 8.1.

Dacă graful de alocare a resurselor nu conține nici un ciclu atunci nici un proces din sistem nu se află în deadlock.

Dacă fiecare resursă are exact o instanță atunci un ciclu implică apariția unui deadlock. Dacă ciclul implică doar un set de resurse de un anumit tip, fiecare din ele având doar o instanță atunci a apărut un deadlock. Fiecare proces ce face parte din acest ciclu este în deadlock. Dacă fiecare tip de resursă are mai multe instanțe atunci nu este neapărată sigură apariția unui deadlock. Aceste concepte sunt ilustrate mai bine în fig 8.2 pg 248.

8.3 Metode de rezolvare a Deadlock-urilor

Pentru a rezolva problema putem adopta 3 strategii:

- Putem folosi un protocol pt a preveni sau evita deadlock-urile, astfel ne asigurăm că sistemul nu va ajunge niciodată în această stare.
- Putem permite ca sistemul să intre în deadlock, să îl detectăm și apoi să rezolvăm situația.
- Putem să ignorăm problema, și să presupunem că nici un deadlock nu a avut loc în sistem. Această soluție este folosită de majoritatea sistemelor de operare.

În continuare ne vom ocupa pe rând de fiecare dintre aceste strategii.

8.4 Prevenirea DeadLock-urilor

Pentru a putea prevenii aparitia acestor situatii trebuie sa ne asiguram ca cel putin una dintre conditiile prezentate in cap 8.2.1 nu este indeplinita:

8.4.1 Excluziune Mutuala

Aceasta conditie trebuie indeplinita numai pentru resurselor nepartajabile. Resursele partajabile nu necesita excluziune mutuala(exemplu fisierele read-only: ele pot fi folosite simultan de mai multe procese). In general nu putem preveni aparitia unui deadlock asigurand excluziune mutuala pt ca unele resurse sunt prin definitie nepartajabile.

8.4.2 Hold and Wait

Pentru a ne asigura ca nu este indeplinita aceasta conditie trebuie ca un proces sa nu ceara o resursa in timp ce tine ocupata o alta. In acest caz putem folosi un protocol care presupune ca un proces sa ceara si sa primesca resurse inainte sa isi inceapa executia.

Un alt protocol viabil ar fi ca sa permitem unui proces sa ceara resurse numai daca acesta nu mai foloseste si altele. Aceste protocoale au 2 mari dezavantaje: utilizare resurselor nu se face eficient, pt ca procesele tin ocupate resurse dar nu le folosesc; este posibil ca un proces sa astepte la infinit dupa o resursa ce este ocupata de alt proces.

8.4.3 Fara Preemptiuni

Pentru a fi siguri ca nu se indeplin aceasta cond folosim protocolul urmator: daca un proces retine anumite resurse si cere alte resurse care nu ii pot fi alocate in acel moment, atunci toate resursele detinute sunt preemptionate(sunt eliberate). Resursele preemptionate sunt adaugate la lista de resurse pe care le asteapta procesul. Acesta va fi restartat numai dupa ce va putea folosi vechile resurse precum si cele noi pt care a facut cerere.

Deasemenea daca un proces cere niste resurse care nu sunt disponibile verificam daca acele resurse sunt alocate unui proces care asteapta dupa alte resurse. Daca este asa eliberam resursele de procesul vechi(le preemptiam) si le alocam procesului nou(cei care a facut cererea).

8.4.4 Asteptare Circulara

Fie $R = \{R_1, R_2, \dots, R_n\}$ setul de resurse al sistemului. Fiecarei resurse ii asignam un numar si obligam procesele sa ceara resurse doar in ordinea crescatoare a numerelor. De exemplu $F(\text{disk drive})=1$, $F(\text{imprimanta})=5$. Astfel procesul care vrea sa foloseasca diskul si imprimanta trebuie sa ceara mai intai accesul la disk pt ca numarul asignat acestuia este mai mic. Cand un proces cere o resursa ce are asignat, sa zicem numarul i , atunci el trebuie sa elibereze toate resursele ce au numere $\geq i$.

Daca sunt indeplinite aceste doua conditii atunci **asteptarea circulara** nu va avea loc niciodata.

8.5 Evitarea DeadBlock-urilor

Daca aplicam algoritmi descrisi la 8.4 vom putea evita deadlock-urile dar pot aparea si efecte secundare cum ar fi utilizarea scazuta a resurselor sistemului.

O alternativa la aceasta metoda ar fi sa cunoastem despre cum vor fi alocate resursele. Daca cunoastem exact secventa de cereri si eliberari ale resurselor(de exemplu avem doua procese: primul are nevoie mai intai de floppy si apoi de imprimanta si al

doilea exact invers), atunci putem decide la fiecare cerere daca procesul va intra in stare de asteptare sau nu.

Cel mai simplu algoritm de evitare presupune ca fiecare proces sa decare un **numar maxim** de resurse de care are nevoie.

8.5.1 Stari Sigure

O stare a unui sistem este sigura daca sistemul poate aloca resurse fiecarui proces intr-o anumita ordine astfel incat sa evite un deadlock. Cu alte cuvinte trebuie sa existe o **secventa sigura** de procese. Orice proces P_i din secventa poate primi resursele cerute, iar un proces P_j , $j < i$ ii va ceda resursele detinute de el. Daca nu exista o asemenea secventa atunci starea nu este sigura(unsafe). Nu toate stările nesigure sunt deadlock.

8.5.2 Algoritmul grafului de alocare a resurselor

Daca avem un sistem cu cate o singura instanta din fiecare tip de resursa, atunci putem folosi o varianta a grafului de alocare a resurselor pentru a evita deadlock-urile. In plus fata de graful descris anterior mai adaugam un nou tip de muchie, **muchia de cerere posibila**($P_i \rightarrow R_j$), care indica ca procesul P_i s-ar putea sa aiba nevoie pe viitor de resursa R_j . Se reprezinta printr-o linie punctata.

Cand cand aceasta cerere chiar are loc atunci muchia este convertita in muchie de cerere, iar cand resursa este eliberata este convertita inapoi in muchie de cerere posibila. Daca nu exista nici un ciclu atunci sistemul se afla intr-o stare sigura.

8.5.3 Algoritmul bancherului(sau algoritmul bancar)

Acest algoritm se poate aplica si pt mai multe instante aleunei resurse dar nu este la fel de eficient ca cel prezentat mai sus. Orice proces trebuie sa declare mai intai numarul maxim de instante de resurse de care s-ar pute sa aiba nevoie.

Consideram n numarul de procese din sistem, si m numarul de tipuri de resurse. Pentru a implementa acest algoritm trebuie sa declarăm mai multe structuri de date:

- **Available**: un vector de dimensiune= m care indica numarul de instante disponibile dintr-o anumita resursa.
- **Max**: o matrice $n \times m$ care indica cererea maxima de instante ale unei resurse facuta de un proces. Daca $Max[i,j]=k$, procesul P_i va cere cel mult k instante ale resursei R_j .
- **Allocation**: o matrice $n \times m$ care indica numarul de instante de fiecare tip alocate in prezent unui proces. Daca $Allocation[i,j]=k$, procesul P_i are in prezen alocate k instanet ale resursei R_j .
- **Need**: o matrice $n \times m$ care indica numarul de resurse de care mai are nevoie un proces. Daca $Need[i,j]=k$, procesul P_i mai are nevoie de inca k instante din resursa R_j pt a-si termina misiunea.

Notatii: X is Y vectori de lungime n . Spunem ca $X \leq Y$ daca $X[i] \leq Y[i]$, unde $i=1,2,\dots,n$. Facem aceasta notatie pt a putea trata fiecare linie din matricile Allocation si Need ca pe niste vectori corespunzatori unui proces. Allocation(i) indica resursele alocate procesului P_i .

8.5.3.1 Algoritmul de siguranta- algoritm care ne spune daca sistemul este sau nu in stare sigura

1. Fie Work[m] si Finish[n] doi vectori. Se initializeaza Work:=Available, Finish[i]=false $i=1,n$

2. Se cauta un i astfel incat:
 - a. $Finish[i]=false$
 - b. $Need(i) \leq Work$
3. $Work := Work + Allocation(i)$
 $Finish[i] := true$
 Go to step 2
4. If $Finish[i]=true$ pentru toti i atunci sistemul e in stare sigura.

Acest algoritm are nevoi de $m \times n^2$ operatii.

8.5.3.2 Algoritmul cererii de resurse

Fie $Request[i]$ vectorul de cereri pt procesul P_i . Daca $Request[j]=k$, atunci procesul P_i cere k instante ale resursei R_j . Cand P_i face o cerere de resurse se executa urmatoarii pasi:

1. If $Request(i) \leq Need(i)$, go to step 2 Else se activeaza un event de eroare pt ca procesul a depasi numarul maxim de resurse de care a spus ca are nevoie.
2. If $Request(i) \leq Available$, go to step 3 Else P_i trebuie sa astepte pt ca resursele nu sunt disponibile.
3. $Available := Available - Request$;
 $Allocation := Allocation + request$;
 $Need(i) := Need(i) - Request(i)$;

Daca se intra intr-o stare sigura atunci P_i primeste resursele, daca nu P_i trebuie sa astepte dupa $Request(i)$ si este restaurata vechea stare.

8.5.3.3 Un exemplu ilustrativ pg 259

8.6 Detectia DeadLock-urilor

Avem nevoie de:

- Un algoritm care sa examineze starea unui sistem pt a determina daca a aparut un deadlock
- Un algoritm care sa revina din deadlock(recover)

8.6.1 Cate o singura instanta din fiecare resursa

Daca fiecare resursa are exact o instanta atunci putem defini un algoritm de detectie care foloseste o varianta a grafului de alocare a resurselor numita **graf de asteptare(wait-for graph)**. Obtinem acest graf din grful de alocare din care nodurile pentru resurse si colapsam muchiile corespunzatoare. Daca in graf gasim un ciclu atunci a aparut un deadlock. Exemplu fig 8.7 pg 261.

8.6.2 Mai multe instante pentru fiecare resursa

Se folosesc aceleasi structuri de date prezentate la algoritmul bancherului(8.5.3).

1. Fie $Work[m]$ si $Finish[n]$ doi vectori. Se initializeaza $Work := Available$
 For $i=1,2,\dots,n$
 - if $Allocation(i) \neq 0$
 - $Finish[i] = false$
 - else
 - $Finish[i] = true$
2. Se cauta un i astfel incat:
 - a. $Finish[i]=false$
 - b. $Request(i) \leq Work$

- Daca nu exista un astfel de i go to step 4
3. $Work := Work + Allocation(i)$
 $Finish[i] := true$
 Go to step 2
 4. If $Finish[i] = false$ pentru anumiti i atunci sistemul e in stare deadlock. Mai exact daca $Finish[i] = false$ atunci P_i este deadlock.

Acest algoritim are nevoie de $m \times n^2$ operatii.

8.6.3 Eficienta algoritmului de detectie

Apelul acestui algoritim se face in functie de 2 factori:

1. Cat de des a aparut un deadlock
2. Cate procese vor fi afectate de deadlock in momentul cand acesta va aparea

Intr-un caz ideal am putea apela acest algoritim la fiecare cerere de resursa, dar ar rezulta o mare pierdere de putere de calcul. O solutie mai ieftina ar fi sa il apelam la anumite intervale de timp.

8.7 Rrevenirea din starea de DeadLock(Recovery)

Cand s-a detectat un DeadLock avem mai multe posibilitati:

- Informam operatorul, deci va intra in sarcina acestuia sa rezolve problema manual
- Sistemul isi revine automat(automatically recover)

Sunt 2 solutii de a revenii dintr-un Deadlock:

1. Renuntam la unul sau mai multe procese pt a intrerupe asteptarea circulara.
2. Preemptionam resurse de la procesele care cauzeaza deadlockuri(le luam resursele)

8.7.1 Terminarea proceselor

Metode de abandonare a proceselor:

- **Abandonarea tuturor proceselor ce cauzeaza deadlock:** are dezaavantajul ca intrerupand multe procese pierdem rezultatele pe care acestea trebuiau sa le returneze, deci irosim timp.
- **Abandonarea tot a cate unui proces pe rand pana cand dealock-ul este eliminat:** de fiecare data trebuie apelat si un algoritim de detectie pt a vedea daca mai avem deadlock

Factori pe care ii luam in considerare cand alegem procesul pe care il intrerupem:

1. Ce prioritate are procesul
2. De cat timp este pornit acel proces si cat timp mai are pana se termina
3. Cat de multe resurse a folosit acel proces si ce fel de resurse(simple sau preemptionate)
4. De cate resurse mai are nevoie procesul pentru a termina
5. Procesul mai interactioneaza sau nu cu altele

8.7.2 Preemptionarea Resurselor

Prin preemptionarea resurselor unui proces, ii luam acestuia resursele si le atribuim unui alt proces pana cand ciclul deadlock este intrerupt. Trebuie sa rezolvam urmatoarele probleme:

- **Selectionare unei victime:** trebuie sa avem in vedere un cost minim(cat mai putine resurse ocupate de acel proces si un timp de cand se afla in executie cat mai mic)

- **Rollbak(refacerea contextului):** dupa ce am preemptionat resursele de la un process, la revenirea acestuia trebuie sa refacem stare in care se afla. Cea mai simpla solutie este sa il abandonam de tot si pe urma sa il restartam, dar nu este asa de eficient.
- **Infometare(Starvation):** orice proces nu trebuie ales ca victima de prea multe ori pentru a se evita “infometarea” lui.

CAP 9

Managementul memoriei

Metodele de management al memoriei nu tin numai de sistemul de operare si de partea hardware. Multi din algoritmi de management al memoriei necesita support hardware.

9.1.1 Legatura dintre adrese

In mod normal un program se afla pe disc sub forma unui fisier executabil. Programul trebuie adus in memorie pus intr-un process si apoi executat. In functie de managementul memoriei procesul se poate fi mutat de pe disc in memorie in timpul executiei lui.

Colectia de procese de pe disc care asteapta sa fie aduse in memorie pentru a fi executate formeaza coada de intrare.

Datorita faptului ca spatiul un program se poate afla (de obicei) in orice zona a memoriei fizice este nevoie de o mapare a adreselor facandu-se o legatura intre doua spatii de adrese (adresele logice sunt transformate in adrese fizice).

Legatura dintre instructiuni (si date) si adresa de memorie se poate face astfel:

1. *In timpul compilarii:* se genereaza codul absolut (adresele din program coincide cu cele fizice). Daca programul este mutat de la locatia initiala el trebuie sa fie recompilat.
2. *In timpul incarcarii (load time):* Compiler-ul genereaza un cod care isi poate muta pozitia (relocatable code). Daca adresa de start se schimba vom reincarca codul adugand noua valoarea.
3. *In timpul executiei:* Legatura dintre adrese se face in timpul executiei.

9.1.2 Spatiul de adrese logice vs spatial de adrese fizice

CPU genereaza de obicei adrese logice. Schimarea adreselor logice in adrese fizice se face cu ajutorul MMU (memory management unit).

O schema simpla de mapare a adreselor se face cu ajutorul unui registru de relocare. In acest registru se memoreaza o valoare (ex 14000) care se adauga la fiecare la fiecare adresa logica (ex 14000 +456 unde 456 adresa logica).

9.1.3Incarcarea dinamica

In schema de mai sus ca un process sa fie executat el trebuie sa se afle in memoria fizica. Deci procesele vor fi limitate de dimensiunea fizica. Pentru a obtine o mai buna utilizare a spatiului de memorie se poate folosi incarcarea dinamica (dynamic loading). In acest caz un subprogram (rutina) nu este incarcata in memorie pana cand nu este chemata. Aceasta metoda se foloseste atunci cand avem programe cu cod mult.

9.1.4Legaturi dinamice si Impartirea librariilor

Pentru economisirea spatiului se folosesc legaturi catre librarii. In cod se face doar legatura catre o librerie. Aceasta metoda permite modificarea doar a librariilor (repararea de bug-uri) fara ca legaturile sa se modifice (se pot face update-uri la program folosind aceasta metoda).

Impartirea librariilor(shared libraries) se refera la modificarile mari aduse la o librerie care ar putea afecta executia programelor. In acest caz se retine pe langa legatura si versiunea libreriei. Programele care au fost compilate cu o versiune mai veche vor continua sa ruleze pe versiunea veche fara ca programul sa sufere modificari.

9.1.5Extindere

Pentru a permite unui proces sa fie mai mare decat spatiul de memorie alocat lui se foloseste extinderea (overlay). Ideea este de a tine in memorie doar acele instructiuni si date care sunt necesare in orice moment. Cand o alta instructiune trebuie incarcata ea este pusa peste spatiul ocupat de o instructiune de care nu mai este nevoie.

9.2 Transferul de memorie(swapping)

Un process poate fi scos afara din memorie intr-un spatiu de stocare (backing store) si adus inapoi in memorie pentru continuarea executiei. Acest transfer se foloseste la algoritmul round-robin si algoritmi bazati pe prioritate (eu cred ca la toti algoritmi preemptivi-precauti). In mod normal un process care a fost scos din memorie va fi pus in acelasi loc de unde a fost scos datorita legaturii dintre adrese. Aceasta restrictie este inlaturata daca se foloseste maparea a adreselor in timpul executiei. Spatiul de stocare (backing store) este de obicei un disc foarte rapid care are un spatiu suficient pentru a retine mai multe procese.

Solutia transferului de memorie (swap) este folosita in putine sisteme datorita timpului mare dintre transferuri si timpului mic de executie. (dureaza mult transferal din memorie in backing store si invers).

9.3 Alocarea contigua a memoriei

Memoria este de obicei impartita in doua zone: una alocata sistemului de operare si una alocata aplicatiilor user. Prin alocarea contigua fiecare process contine o singura sectiune contigua de memorie.

9.3.1 Protectia memoriei

Diferentia dintre o zona de memorie a sistemului de operare si o zona de memorie dedicata user-ului se poate face cu ajutorul unui registru de relocare si un registru de limita. Registrul de relocare contine adresa fizica cea mai mica iar registrul limita contine dimensiunea maxima a adresei logice.

9.3.2 Alocarea memoriei

Cea mai simpla metoda de alocare a memoriei este de a divide memoria in partitii de dimensiune fixa. Fiecare partitie contine exact un proces. Atunci cand o partitie este goala un process este selectat din coada de intrare si incarcat in acea partitie. Sistemul de operare retine intr-o tabela care parte a memoriei e ocupata si care nu.

Initial toata memoria este libera pentru procesele user. Cand un process este incarcat si are nevoie de memorie se cauta un spatiu suficient de mare pentru process. I se alocă doar atat cat este nevoie restul memoriei fiind disponibila pentru celelalte procese. La terminarea procesului spatiul ocupat de process poate fi reutilizat.

Sistemul de operare alocă memorie proceselor pană ce nu se mai găsește spațiu pentru process. În acest caz sistemul de operare poate să aștepte până se găsește suficient spațiu pentru process sau să aloce spațiu pentru un process care are nevoie de mai puțină memorie.

În cazul în care sunt mai multe locuri în care se poate pune procesul în memorie se pot folosi trei strategii:

- a) First fit: se pune în primul spațiu găsit
- b) Best fit: se pune în spațiul cel mai mic (suficient pentru process)
- c) Worst fit: se pune în spațiul cel mai mare

Acești algoritmi suferă însă de o fragmentare externă datorată încărcării și ștergerii din memorie a proceselor. O fragmentare externă există atunci când spațiu liber din memorie este suficient pentru încărcarea unui process dar zona nu este contigă.

9.3.3 Fragmentarea

Fragmentarea poate fi internă sau externă. Fragmentarea internă este diferența dintre spațiul ocupat de o partitie și spațiul ocupat de procesul care trebuie adus în memorie. Fragmentarea externă este spațiul liber total dar necontigă. Reducerea fragmentării externe se face prin compactare. Compactarea se poate face doar dacă alocarea este dinamică. Scopul compactării este de a pune toate locurile libere într-un singur bloc mare.

9.4 Paginarea

Permite un spatiu de adrese logice necontinue. Prin paginare memoria fizica e impartita in blocuri de dimensiune fixa numite frame (dimensiune fiind o putere a lui 2) iar memoria logica e impartita in blocuri de aceeasi dimensiune numite pagini. Partea hardware care ajuta la paginare e formata dintr-o tabela de paginare (care retine adresa de baza a fiecarui frame din memoria fizica). CPU-ul emite o adresa formata din doua parti : numarul paginii si offsetul paginii. Cu ajutorul numarului paginii se ia din tabela adresa fizica de baza care impreuna cu offsetul formeaza adresa fizica. Lungimea frame-lor este definita prin hard. Pentru a rula un program cu dimensiunea de n pagini vom avea nevoie de n frame libere. Informatia despre numarul total de frame libere , framele libere, framele ocupate se gasesc in tabela de frame.

9.4.2 Implementarea hardware

Implementarea hard a tabelii de paginare poate fi facuta cu ajutorul unui set de registre dedicate. Aceste registre lucreaza la o viteza mare pentru a putea face translatia cat mai completa. Folosirea registrelor este rezonabila daca tabela de paginare are o dimensiune mica (pana la 256 de intrari). Pentru masinile cu tabela de paginare foarte mare , tabela se tine in memoria principala si se foloseste un registru numit PTBR (page table base register) care pointeaza la aceasta zona de memorie. Aceasta metoda este si ea dezavantajoasa deoarece se face doua accese la memorie.

Solutia standard pentru implementarea tabelii de paginare este folosirea TLB-ului (translation look aside buffer). TLB-ul este o memorie asociativa foarte rapida. Fiecare intrare in TLB contine o cheie si o valoare. Cand se acceseaza TLB-ul se compara cheia cautata cu cheile din TLB (principiul memoriilor asociative). Daca se gaseste atunci se returneaza valoarea de la acea cheie (numarul de intrari este de obicei 64-1024).

Daca numarul paginii nu se afla in TLB (TLB miss) atunci trebuie sa se faca o referinta la memorie la tabela de paginare. La miss in TLB se adauga cheia si valoarea pentru ca data viitoare sa nu mai fie miss. Daca TLB-ul este plin atunci se foloseste algoritmul LRU pentru eliminare. Unele informatii din TLB nu se pot elimina fiind folosite de kernel.

Aceste intrari se numesc wired down.

Unele TLB-uri stocheaza un identificator al spatialei de adrese (ASID) in fiecare intrare a tabelii. Cand TLB-ul este folosit atunci pe langa verificarea cheii se verifica si ASID-ul procesului current sa fie egal cu cel din TLB altfel miss in TLB.

Timpul de acces la memorie este : $\text{hit rate TLB} \cdot \text{timp TLB} + \text{miss rate TLB} \cdot \text{timp acces memorie}$.

9.4.3 Protectia

Protectia se face prin asocierea unui bit fiecarui frame (poate fi read-write sau read only). Inca un bit se adauga la intrarea in TLB care indica daca informatia din TLB este valida sau nu. O informatie din TLB este valida daca adresa logica pentru care se cere maparea face parte din spatiul de memorie a procesului.

9.4.4 Structura tabelii de paginare

9.4.4.1 Paginarea ierarhizata

Spatiul de adrese logice este impartit in mai multe tabele de paginare. Un exemplu simplu este un algoritm de paginare pe doua nivele. La aceasta metoda adresa logica (formata din numarul paginii si offset) numarul paginii este impartit in doua parti (ex daca este o adresa pe 32 biti 20 biti pt nr paginii si 12 pt offset. Cei 20 biti se impart in 10b si 10b). Prima parte pointeaza spre o tabela de paginare externa (outer page table) iar a doua parte este deplasamentul paginei din interiorul tablei de paginare.(se foloseste la Pentium 2). Metoda de translatie a adresei:

9.4.4.2 Hashed Page Table(Tabele de dispersie)

Fiecare intrare in tabela de dispersie contine o legatura catre o lista de elemente care contine adrese din acelasi spatiu de adresare. Algoritmul functioneaza astfel: din adresa logica se ia doar numarul paginii q (nu si offsetul) si se trece printr-o functie de dispersie. Valoarea obtinuta este indexul din tabela de dispersie. Fiecare intrare din tabela de dispersie avand atasata o lista se face o cautare in acea lista dupa valoarea paginii q . Fiecare element din lista contine: valoarea paginii(adresa logica), valoarea frame-ului(adresa fizica), adresa next element.

9.4.4.3 Inverted page table

O tabla de paginare inversata are proprietatea ca fiecare pagina sau frame are o singura intrare. In tabela de paginare inversata se face indexarea cu ajutorul PID (process ID). CPU-ul in acest caz returneaza o adresa logica de forma <PID, nr pagina,offset>.

Pentru a micșora timpul de cautare se poate folosi și o tabelă de dispersie.

9.4.4.5 Pagini partajate (shared page)

Se pot folosi în cazul în care procesele nu au un cod care se schimbă în timpul execuției. În memorie fizică sunt reținute o copie a codului program și pentru fiecare user datele introduse. Avantajul constă în faptul că nu se ocupă spațiu în memorie (nu se copiază codul programului pt fiecare user în parte).

9.5 Segmentarea

Un spațiu de adrese logice este împărțit în mai multe segmente. Fiecare segment are un nume și o lungime (segment de date, segment de cod..).

9.5.2 Hardware

Segmentarea se face cu ajutorul unei tabele de segmentare. Adresa logică (emisă de CPU) conține un index și un offset. Indexul pointează la tabelă de segmente în ea aflându-se un tuplu de forma <limită, bază>.

9.5.3 Partajarea segmentelor

Segmentarea ajută foarte mult în partajare. Cu ajutorul segmentelor se pot partaja doar datele unui proces și nu și codul acestuia.

9.5.4 Fragmentarea

Are loc o fragmentare externă ca la paginare.

9.6 Paginarea și segmentarea

Se folosesc ambele idei. Adresa logică este împărțită în: page number și page offset. Cu page number se pointează la o tabelă de segmentare care ne dă bază page –ul pt tabelă de paginare. Din tabelă de paginare se ia adresa fizică.

Capitol 10 MEMORIA VIRTUALĂ

Mem virtuală este o tehnică care permite execuția proceselor care nu pot fi completate în memorie.

10.1 Background

Mem virtuală este separația dintre mem logică de mem fizică. Separarea permite ca mem virtuală multă să fie disponibilă la numai o mem fizică mică.

Avantaje:

- doar o parte a programului trebuie sa se afle in mem cand se executa.
- spatiul logic de adrese poate fi mai mare decat spatiul fizic de adrese
- spatiul de adrese poate fi impartit de mai multe procese(prin page sharing)
- Permite o creare de procese mult mai eficienta

Mem virtuala poate fi implementata prin: demand paging sau demand segmentation.

10.2 Demand paging

Este similar cu un sistem de pagini cu schimbare. Procesele se afla in mem secundara si sunt aduse in mem doar cand vor fi executate.

10.2.1 Concepte de baza

poza pg 320

Cand un proces trebuie schimbat, paginatorul ghiceste care pagini vor fi folosite pana cand procesul va fi schimbat iara. Paginatorul aduce doar paginile necesare in mem => Mai putine I/O, mai putina mem folosita, raspuns mai rapid, mai multi useri.

Avem nevoie de o forma de ajutor hardware sa stim daca pagina este in mem sau pe disc -> bit valid/invalid. Daca bitul este setat pe valid -> pagina este valida si in mem. bit = invalid -> pagina este fie invalida sau se afla pe disc. (poza pag 321). Initial bitul valid/invalid va fi pus pe 0 pentru toate paginile. Daca procesul nu va accesa pagina nu conteaza daca pagina este valida sau nu. Daca ghicim paginile corecte procesul va rula ca si cum am fi adus toate paginile.

Daca procesul acceseaza o pagina care nu se afla in memorie(accesul la o pagina marcata ca invalid) => page fault.

Tratarea unui page fault:

- cautam in tabela intern a procesului sa vedem dac este o referinta valida la memorie sau nu
- daca referinta este invalida terminam procesul, daca este valida si pagina nu a fost adusa inca se aduce pagina
- gasim un frame gol
- executam o operatie de disc de citire a paginii in frame-ul nou alocat
- dupa citire modificam tabela interna a procesului indicand ca pagina este in mem
- restartam instr care a fost intrerupta din cauza page fault

poza pag 322

Paginare la cerere = se aduce o pagina in mem doar daca a fost accesata.

O instr trebuie restartata dupa fiecare page fault. Putem rezolva acest lucru prin aducerea inca odata a instructiunii. Dificultatea este cand o instr modifica mai multe locatii diferite.

Solutii:

- se incearca detectarea unui page fault inainte ca sa se scrie in locatii, inainte ca ele sa fie modificate.
- folosirea unor registrii temporari care pastreaza valorile locatiilor scrise. La un page fault toate valorile vechi sunt scrise in mem.

10.2.2 performata paginarii la cerere

Rata de page fault este: $0 \leq p \leq 1.0$

- daca $p = 0$ nu avem page fault
- daca $p = 1$, fiecare accesare este un page fault

timpul efectiv de acces(EAT):

$EAT = (1 - p) \times \text{memory access} + p (\text{page fault overhead} + [\text{swap page out}] + \text{swap page in} + \text{restart overhead})$

Rata de page fault trebuie pastrata cat mai mica in paginarea la cerere.

10.3 Creearea proceselor

10.3.1 Copy on write

Aceasta tehnica permite ca Procesul parinte si procesul fiu sa imparta initial aceleasi pagini. Aceste pagini sunt denumite pagini copy on write. Daca unul din procese scrie intr-o pagina partajata, doar atunci se face o copie a paginii partajate. Tehnica permite o creare eficienta a proceselor deoarece doar pagina modificata este copiată. Paginile copiate trebuie puse in unele libere. Trebuie determinat de unde luam acestese pagini. multe SO au un pool(grupa) de pagini libere pentru asemenea cereri.

10.3.2 Fisiere mapate in Mem

I/O de fisiere mapate in mem permite ca I/O de fisiere sa fie tratat ca un acces de mem prin maparea unui bloc de disc intr-o pagina din mem. Un fisier este citit initial folosind paginarea la cerere. O portiune de marimea unui pagini din fisier este apoi citita din sistemul de fisiere intr-o pagina fizica. Citiri/scrieri secventiale la/din fisier sunt tratate ca accese obisnuite la mem. Simplifica accesul la fisiere prin tratarea fisierelor ca I/O prin memorie decat prin apelurile de sist read() write(). Deasemenea permite ca mai multe procese sa mapeze acelasi fisier prin partajarea paginilor in mem. poza pag 331

10.4 Inlocuirea paginilor

Previne supraincercarea memoriei prin modificarea rutinei de tratare page fault sa includa si inlocuirea paginilor. Supraincercarea memoriei apare cand toata mem este folosita (nu e nici un frame liber pentru a pune noua pagina).

10.4.1 Schema de baza

Daca nici un frame nu este liber gasim unul care nu este folosit si il eliberam. Continutul frame-ului il scriem in spatiul de schimb (swap) si schimbam tabela de pagini indicand ca pagina nu mai este in mem. Frame-ul care va fi scos din memoria are numele de frame victima. Cand nu este nici un frame liber -> doua transferuri de pagina. Acest lucru poate fi prevenit prin folosirea unui bit dirty. Bitul este setat cand se scrie in pagina. Doar pagini care au fost modificate sunt scrise pe disc.

Algoritmii de inlocuire trebuie sa aiba cea mai mica rata de page fault. vezi ex pag 335

10.4.2 Algoritmul de inlocuire FIFO

Cel mai simplu algoritm de inlocuire. Un algoritm fifo asociaza fiecare pagina cu timpul cand a fost adusa in mem. Cand o pagina trebuie inlocuita se alege cea care este de cel mai mult timp in mem. Se poate implementa cu ajutorul unei cozi fifo. ex pag 336

Anomalia lui Belady: pt anumiti algoritmi se pot obtine rez mai slabe cu cresterea frezeurilor.

10.4.3 Inlocuire optima a paginilor

Un algoritm optimal are cea mai mica rata de page fault-uri si nu sufera de anomalia lui Belady. El este denumit OPT sau MIN:

- inlocuieste pagina care nu va fi folosit pt cea mai mare perioada de timp.

Dificil de implementat, deoarece trebuie sa stim ce pagini vor fi referite in viitor.

Ex pag 338

10.4.4 Algoritmul LRU

Pagina care va fi inlocuita este cea care nu a fost folosita de cel mai mult timp. LRU asociaza fiecarei pagina timpul cand a fost utilizata.

Implementare:

- cu numaratoare: fiecare intrare din tabele de pagini are asociat un camp timpul cand a fost folosita pagina. cand se acceseaza pagina acest contor primeste valoarea ceasului logic al CPU. Aceasta schema implica insa cautarea paginii LRU in tabele.
- cu stiva: cand o pagina este referita ea este pusa in varful stivei. Astfel in vf stivei este pagina accesata de curand iar in coada este pagina LRU. Aceasta implementare nu necesita cautari dar mutarea paginii din mijlocul stivei in varf este consumatoare de timp.

10.4.5 Algoritm LRU aproximativ

Introducerea unui bit de accesare(reference bit). Bitul paginii este initial pus pe 0. Cand pagina este accesata el devine 1. Inlocuirea se face cu pagina care are bitul 0. Nu stim insa ordinea in care au fost accesate paginile.

10.4.5.1 Biti aditionali de accesare

In locul unui bit avem 8 biti. La intervale regulate acesti cate un bit este inregistrat si ceilalti shiftati. ex 11000100 a fost accesat mai recent ca 011110111.

10.4.5.2 Second_chance algoritm

Daca pagina care trebuie inlocuita are bitul de acces 1 atunci el este pus pe 0, se lasa pagina in mem si alta este inlocuita.

10.4.5.3 Algoritm imbunatatit second-chance

se folosesc bit dirty si modificat

- (0, 0) – nu a fost accesata nici modificata recent – pagina cea mai buna de inlocuit
- (0, 1) – nu a fost folosita de curand dar este modificata – nu e f bine la inlocuit deoarece trebuie scrisa inapoi
- (1, 0) – folosita dar nemodificata – poate fi folosita incurand
- (1, 1) – recenta si modificata

10.4.6 Algoritmi de inlocuire bazati pe numarare

Tin o evidenta de cate ori au fost accesate paginile:

- LFU: inlocuieste pagina care a fost accesata de cele mai putine ori
- MFU: se bazeaza pe argumentul ca paagina LFU abia a fost adusa si trebuie sa fie folosita

10.4.7 Algoritmi cu Page-buffering

Cand pagina este adusa din mem ea este introdusa intr-un buffer pana cand pagina va fi evacuata din mem. In asa fel nu mai trebuie sa asteptam cu executia pana aceea sa va fi evacuata.

10.5 Frame-uri de alocare

Fiecare proces are nevoie de unmainar min de pagini.

10.5.2 Algoritmi de alocare

Alocare fixa:

- Equal allocation – e.g., if 100 frames and 5 processes, give each 20 pages.
- alocare proportionala – alocarea are loc in functie de marimea procesului

Alocare prioritara:

- se foloseste o schema de alocare care foloseste prioritati si nu marime
- Daca un price P_i genereaza un page fault: se selecteaza una din framele lui sa se inlocuiasca sau una de la un proce cu o prioritate mai mica

10.5.3 Alocare globala VS alocare locala

- Inlocuire locala: fiecare proces selecteaza doar din setul de frame-uri alocate lui
- Inlocuire globala: Procesul selecteaza un frame de inlocuit din setul total de frame-uri; un proces poate lua de la celalalt

10.6 Thrashing

Thrashing este fenomen cand un proces este ocupat cu schimbarea paginilor.

10.6.1 Cauza

Daca un proces nu are destule pagini, rata de page fault este mare. Aceasta conduce la Utilizare slaba a CPU, SO crede ca trebuie sa creasac nivelul de mutiprogramare, alt proces este adaugat sistemului. Fig pag 349

Ca sa marim utilizarea CPU si sa inlaturam thrashing trebuie sa micsoaram gradul de multiprogramare. Putem limita efectele Thrashingului prin folosirea unui algoritm de inlocuire local. Cand un proces incepe thrashing nu poate lua frame-uri de la un alt proces si sa il faca si pe el sa thrashing.

10.6.2 modelul Set de lucru(working set)

Modelul este bazat pe presupunerea de localitate. Folosete parametrul Δ pentru a defini fereastra setului de lucru. Setul de pagini din pagina Δ cea mai recenta se numeste set de lucru. Daca pagina se foloseste activ se afla in setul de lucru. Acuraterea setului de lucru depinde de Δ . Daca e prea mic nu acopera toata localitatea si daca e prea mare poate acoperi mai multe localitati.

SO monitorizeaza setul de lucru al fiecarui proces si alocat setului un numar suficient de frame-uri. Daca mai sunt frame-uri libere poate porni alt proces. Daca nr de seturi de lucru se mareste SO poate decide sa suspende un proces. Aceasta strategie previne Thrashingul si mentine multiprogramarea la cel mai ridicat nivel. Dificultate: Pastrarea evidentei asupra setului de lucru. Setul de lucru este o fereastră care se misca

10.8 Alte considerente

10.8.1 Prepaging

Aducerea in mem la un moment toate paginile care vor fi folosite

10.8.2 Marimea paginii

O marime mica creste numarul paginilor. Deoarece in tabela de pagini fiecare pagina are o copie se doreste un nr mic de pagini -> marime mare a paginii.

Memoria este mai bine utilizata daca folosim pagini mici. Avem rezolutie mai buna, permitandu-ne sa izolam doar portiunea de memorie de care avem nevoie. Cu o pagina mica numarul total de I/O este redus din cauza ca localitatea este imbunatatita.

10.8.3 TLB reach

TLB Reach = se refera la cantitatea de mem ce poate fi accesata din TLB.

$TLB\ Reach = (TLB\ Size) \times (Page\ Size)$

In cazul ideal setul de lucru al fiecarui proces este tinut in TLB. Altfel numarul de page fault va fi mare.

Cresterea marimii TLB-ului:

- Un numar de pagini mai mari: Poate duce insa la cresterea fragmentarii din cauza ca nu toate aplicatiile au nevoie de pagini mari
- Asigurarea marimilor multiple ale paginii: Permite aplicatiilor care au nevoie de pagini mai mare sa le foloseasca fara sa mareasca fragmentarea.

10.8.5 Structura programului

```
F int A[1024][1024];
```

F Each row is stored in one page

```

F Program 1 for (j = 0; j < A.length; j++)
              for (i = 0; i < A.length; i++)
                  A[i,j] = 0;
1024 x 1024 page faults
F Program 2 for (i = 0; i < A.length; i++)
              for (j = 0; j < A.length; j++)
                  A[i,j] = 0;

```

1024 page faults

10.8.6 I/O interlock

- n **I/O Interlock** – Pages must sometimes be locked into memory.
- n Consider I/O. Pages that are used for copying a file from a device must be locked from being selected for eviction by a page replacement algorithm.

Cap 12

12.1 Avantaje discuri: - pot fi rescrise

- se poate accesa direct orice bloc.

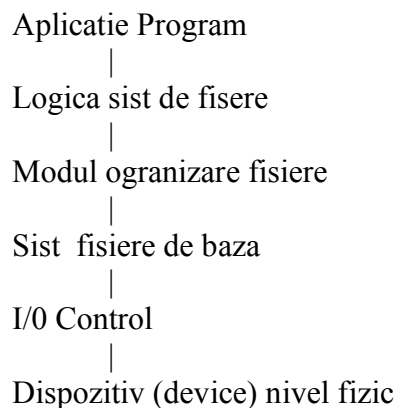
-transferul DISC-Mem : se transfera blocuri (unul sau mai multe sectoare).

Marimea blocurilor variaza intre 32B-4096B, usual 512B.

Sistemul de fisere ofera un acces convenabil la disc. Implementarea unui system de fisere are la baza 2 concepte:

- definirea fisierului in raport cu utilizatorul adica definirea atributelor fisierului si a operatiilor premise precum si directorul din care va face parte.
- crearea unui algoritm care face legatura intre logica sist. de fisere si nivelul fizic al discului.

Sist de fisere este compul dinmai multe nivele. Fiecare nivel se foloseste de serviciile nivelele inferioare si adauga noi servicii care vor fi folosite de nivelele superioare.



I/O Control: Este format din drivere pt dispozitive si handlerului de intrerupere. Un dispozitiv poate fi privit ca un translator el primește mesaje de tipul “citește bloc 113” si transforma aceste mesaje in altele specifice hardware. Un dispozitiv scrie diferite biti specifici in controlerul de I/O pt a alege o anumita locatie si o anumita actiuni care este ceruta.

Sistemul de fisiere de baza: trimite comenzi de scriere/citire a blocurilor fizice de pe disc catre dispozitivele corespunzatoare . (mode de adresare a unui bloc: disc, cilindru, pista sector).

Modul de organizare fisiere: face traslatarea din adresa logica a blocului de fisier in adresa fizica . Fiecare bloc de fisier logic primește un numar de la 1 la N acesta nu este obligatoriu identic la nivel fizic. Tot aici se afla si managerul pentru spatial liber care tine evidenta blocurilor nealocate.

Logica sist de fisiere: se ocupa cu informatiile “metadata”. Acestea contin structura fisierului mai putin data stocata in el . Acest nivel se ocupa de structura directoarelor si ofera nivelului inferior toate informatiile de care are nevoie pentru un anumit fisier. El mentine structura fisierelor cu ajutorul file control blocks (FCB). FCB contine informatii despre fisier (proprietar, permisiile, locatie). Tot aici este implementat si sistemul de protectie si securitate.

12.2 Implementare sistem de fisere

On disk:

Bloc de boot control: are informatii necesare pentru bootare unui sistem operare (doar pe partițiile care contin un SO).

Blocul de control al partiției: detalii cu privire la partiție : nr de blocuri, dimensiune, dimensiune bloc, blocuri libere, FCB, FCB libere, etc.

O structura de directoare care organizeaza fisierele

FCB descrie mai sus. (figura 12.2 pag 415)

In memory: este folosita pentru managementul unui sist de fisiere si pentru imbunatatirea performantelor folosind cacheing.

Tabela de partiție: informatii cu privire la fiecare partiție incarcata

Structura de directoare: retine informatii despre directoarele accesate recent.

Tabela de fisiere deschise la nivelul sist global (TFG): contine FCB pt fiecare fisier deschis si alte informatii.

Tabela de fisiere deschise la nivel de process (TFP): contine un pointer catre randul corespunzator din tabela de fisere globala si alte informatii.

Creaare fisier: Se creaza un nou FCB, se citește directorul corespunzator se fac modificarile si apoi se scrie inapoi pe disc.

Deschidere fisier: Mai intai se verifica daca fisierul nu este deja deschis cautandu-se in tabela TFG, daca este deja deschis atunci se creaza o noua intrare in TFP care va pointa la randul corespunzator din TFG. Daca fisierul nu este deschis atunci se cauta fisierul in structura de directoare (parti din structura de directoare sunt stocate in memorie), dupa ce se gaseste fisierul FCB-ul lui este copiat in TFG si se alocă un rand in TFP cu un pointer catre TFG. Toate operatiile urmatoare vor fi realizate prin intermediul acestui pointer.. (figura 12.3 pag 417)

12.2.2 Partitii si Mounting

Un disc poate fi impartit in mai multe partitii sau o partitie poate contine mai multe discuri.

Partitie “raw”- nu contine un system de fisiere.

“cooked”- contine un system de fisiere.

Informatiile pentru bootare pot fi stocate pe partitii diferite. De obicei aceste informatii sunt niste secvente de blocuri care se incarca in memorie. Executia lor incepe de la o anumita locatie predefinita. Aceste imagini de boot pot contine mai multe instructiuni decat cele necesare pentru a starta un SO. Sunt cazuri in care pe un PC avem 2 sau mai multe sisteme de operare si trebuie sa alegem pe care il vor porni. Un boot loader care intelege conceptual de multiple sisteme de fisiere si de operare poate fi pus in spatial special rezervat pt boot, dupa ce va fi incarcat el va alege SO.

Partitia root contine kernelul SO si alte informatii utile. Este incarcata la momentul sectiunii de boot. Alte partitii sunt incarcate fie in acelasi moment sau manual mai tarziu acest lucru depinde de SO. SO verifica daca partitia incarcata are un sistem valid de fisiere.

12.2.3 Sistem de fisere virtual

Sistemele de operare moderne suporta mai multe sisteme de fisiere in acelasi timp.

Acest deziderat este realizat prin implementarea sistemului de fisiere in 3 nivele majore. (figura 12.4 pag 419)

Nivelul Virtual File System (VFS) are 2 functii:

- separa sistemului de fisiere generic de implementarea prin interfata VFS.

Mai multe implementari ale VFS pot sa existe in paralel astfel permitandu-se accesul transparent la diferite tipuri de sisteme de fisiere locale.

- VFS are la baza o structura de reprezentare a fisierelor numita **vnode**, aceasta contine un cod numeric unic pentru fiecare fisier din retea in acest fel se poate implementa un sistem de fisiere la nivelul retelei. Kernelul mentine un vnode pentru fiecare nod active (fisier sau director).

VFS distinge fisierele locale de cele aflate pe alte hosturi in retea. Iar fisierele locale sunt la randul lor separate dupa modul de implementare a sistemului de fisiere. Astfel ca VFS stie a aleaga handl-ul corespunzator

12.3 Implementare Director

12.3.1 Lista Liniara

Este cea mai simpla impletare. Se foloseste o lista liniara pentru definirea intrarile in director. Accesul se face tot liniar, timpul fiind destul de mare. Pentru a crea un fisier trebuie cautat tot directorul ca sa vedem ca nu exista si apoi sa il adaugam la final. Ca sa stergem un fisier el trebuie cautat si apoi sters. Spatiul ramas liber in urma stergerii poate fi marcat fie printr-un nume special “all-blank”, prin implementarea unui bit de folosit/nefolosit sau se poate crea o lista care sa retina locatiile libere. O alta posibilitate este sa copiem in locatia libera ultimul fisier din lista si sa scadem marimea directorului cu 1.

O reducere a timpul de cautare se poate face prin mentinrea listei sortate.

12.3.2 Hash Table

Tot o lista liniara mentine intrarile in director. Dar pe langa aceasta se implementeaza si un hash table. Hash table adresat cu numele fisierului (aplica o functie de hash) returneaza un pointer catre locatia acestuia din lista liniara. Timpul de cautare este redus. Inserarea si stergerea se produc implementeaza usor. Singurul lucru de care trebuie tinut cont sunt coliziunile adica 2 nume de fisiere sa poarte catre aceasi locatie. Pentru a preveni fenomenul de hash overflow (ex tabela are doar 64 de locatii si sunt 65 de fisiere in director) se poate ca pointerul returnat sa nu adreseze o locatie ci o lista inlantuita.

12.4 Metoda de alocare a spatiului

12.4.1 Alocatie contigua

Fiecare fisier ocupa o zona contigua pe disc. Acest lucru inseamna ca accesul la un fisier va implica o singura miscare a capului de citire (blocul b+1 este accesibil fara miscarea capului dupa ce acesta s-a pozitionat pentru citirea blocului b).

Problema in acest caz este gasirea unui loc liber suficient de mare pt stocarea fisierului pentru aceasta se folosesc strategiile first fit sau best fit. Totusi acesti algoritmi nu pot maximiza folosirea eficienta a spatiului total disponibil pe disc.=> fragmentare externa mare.

La momentul creeri unui fisier nu se stie de cat spatiu va avea nevoie acel fisier, aceasta este o alta problema a alocatiei contigue. Aceasta problema este rezolvata intr-unele SO prin alocarea unui alt spatiu contiguu in cazul depasirii spatiului alocat.

Avantaje: fragmentarea interna potrivita, acces direct, acces secvential

Dezavantaje: fragmentarea externa.

Fragmentarea interna=spatial pierdut in interiorul spatiului alocat.

12.4.2 Alocare inlantuita

In aceasta implementare fiecare fisier este o lista de blocuri de pe disc. Fiecare bloc poate fi oriunde pe disc(figura 12.6). Directorul contine un pointer catre primul fisier primul si ultimul bloc al fisierului. Fiecare bloc contine un pointer catre urmatorul bloc.

La crearea unui fisier se alocă intrarii in director un pointer la primul bloc al fisierului, initial acest pointer este initializat cu nil (end of list pointer value) adica un fisier gol .La scrierea in fisier un bloc liber va fi gasit si este pus un legatura cu sfarsitul fisierului.

Dezavantajul major este ca poate fi folosit efectiv doar in cazul fisierelor secventiale. Nu putem accesa al n-lea bloc din structura unui fisier decat parcurgand cele n-1 blocuri anterioare. Un alt dezavantaj sunt cei 4 B necesari pentru stocarea pointerului in fiecare bloc. Acest inconvenient se rezolva prin alocarea de clustere(grupari de blocuri) in loc de blocuri, pointerii folositi vor ocupa atunci un loc mai mic (acest lucru duce la cresterea fragmentarii interne: adica in cazul in care un cluster nu este plin atunci se pierde loc mai mult decat in cazul in care doar un bloc nu ar fi plin).

Un al dezavantaj al acestei implementari o constituie fiabilitatea. Daca un pointer este pierdut atunci se pierde calea spre urmatorul bloc/cluster al fisierului.

Tabela FAT: Este pastrata o sectiune pe disc de la inceputul fiecare partitii pt tabela.

Tabela este indexata cu numarul fiecarui bloc. Fiecare intrare contine urmatorul bloc al fisierului. Blocurile libere contin valoarea 0 iar blocul final al fisierului contine o valoare speciala end of file.

Tabela fat imbunatreste accesul direct.

Avantaje alocare inlantuita: fragmentare externe inexistentă.

Dezavantaje: accesul direct foarte slab, fragmentarea internă mai ales în cazul clusterelor.

12.4.3 Alocare indexata

În index block vor fi stocați toți pointerii rezolvându-se astfel problema accesării directe avute de alocarea înlantuită. Fiecare fișier are propriul index block. Intrarea n în system block pointează către blockul n al fișierului. Directorul conține adresa index blockului . figura pag 427.

Această implementare rezolvă accesul direct și nu suferă de fragmentare externă.

Se pune problema cât de mare să fie index blockul.

Linked scheme: mărimea index block = mărimea disc bloc. În cazul fișierelor care necesită un index mai mare se vor face legături între index blocuri.

Multilevel index: există mai multe nivele de index blocuri. Primul nivel conține o listă de index blocuri de nivel secund. Care conține lista de index de nivel inferior și așa mai departe. Cel mai inferior nivel va conține pointerii către blocurile fișierului.

Schema combinată: această variantă presupune reținerea unui număr de pointer către index bloc în nodul fișierului. O parte din acești pointeri vor pointa direct către index blocuri care mapează blocurile fișierului, iar o altă parte vor implementa principiul de multilevel index.

12.5 Managementul spațiului liber

Bit vector: fiecărui bloc îi corespunde o intrare în acest bit vector. Dacă blocul este liber atunci va avea valoarea 1 dacă blocul este ocupat va avea valoarea 0.

Avantajul este simplitate și eficiența metodei. Dar rezultatele nu sunt bune decât dacă vectorul este stocat în memorie și copiat la intervale de timp pe disc pentru consistență. Ori vectorul poate avea o dimensiune semnificativă în cazul discurilor de capacitate mare.

Soluția este folosită doar în cazul microcomputerelor.

Liste înlantuite: se reține doar adrese primului bloc liber acesta va adresa următorul bloc liber și așa mai departe.

Schema nu este eficientă accesul la blocurile libere fiind unul secvențial. De remarca e că metoda FAT are implementat acest mecanism nefiind necesară o redefinire a tabelului.

Grouping: Lista înlantuită va menține adresele a n blocuri libere în fiecare nod. Avantajul este că adresele a mai multor blocuri libere pot fi găsite mai repede.

Counting: În această metodă lista va avea în fiecare nod adresa primului bloc liber și numărul de blocuri libere care îi urmează. Metoda se bazează pe principiul de alocare contiguă sau pe alocarea care folosește cluster sustinând că o ștergere din memorie va elibera de cele mai multe ori mai multe blocuri consecutive. Prin păstrarea acestor date lista înlantuită va scădea în dimensiune.

+12.6 Eficiența și performanța facultative

12.7. Recovery

Fișierele și directoarele sunt ținute atât în memorie cât și pe disc din acest motiv trebuie asigurată consistența datelor. Este de presupus că datele din memorie vor fi mai recente decât cele de pe disc. Cu toate acestea în cazul unei căderi a sistemului datele din memorie vor fi pierdute iar modificările la nivelul discului vor fi incomplete ducând la

eventuale stări de inconsistență a informației. În aceste cazuri la repornirea sistemului va fi lansat un program special, acesta va compara data din structura director cu data din blocurile de pe disc și va încerca să elimine orice inconsistență.

12.7.2 Backup and Restore

Se utilizează un sistem de salvare ciclic folosind n copii ale datelor. În prima zi se va salva copia sistemului întreg în a doua zi doar datele care s-au modificat între timp și așa mai departe. Când ciclul ajunge la final el se va relua printr-o nouă salvare completă a sistemului.

12.8. Log structured File System

Toate schimbările la nivel de metadata sunt scrise secvențial în log. Fiecare set de operații este o tranzacție. După ace schimbările sunt scrise în log ele sunt considerate comise. Intrările din log vor schimba sistemul de fișiere actual și pe măsura ce se fac schimbările un pointer indică care acțiune s-a făcut și care urmează să se facă. Când o tranzacție întreagă este comisă ea este ștearsă din log. Logul este de fapt un buffer circular. În cazul în care sistemul cade starea de consistență se va păstra iar modificările vor putea fi implementate la nivelul sist. de fișiere datorită logului. Singura problemă este în cazul în care sistemul cade în mijlocul unei tranzacții, în acest caz toate modificările făcute trebuie să fie înlăturate.

12.9 NFS

NFS este un exemplu bun de un sistem de fișiere dintr-un sistem de rețea client server.

12.9.1 Overview

NFS vede n set de stații de lucru conectate ca un set de mașini independente cu sistem de fișiere independente. Scopul este să se permită un o anumită partajare între sistemele de fișiere în mod transparent. Să se asigure independența dintre mașini, partajarea unui sistem de fișiere afectează doar mașina client și nu alta mașina. Un director la distanță este instalat peste sistemul local de fișiere. Directorul montat arată ca un subarbore integral al sistemului local de fișiere și înlocuiește subarborile care coboară din directorul local.

Specificările directorului la distanță la operația de montare(suprapunere) este netransparentă, numele gazdei directorului la distanță trebuie furnizat. Fișierele din directorul la distanță pot fi accesate într-un mod transparent. Orice sistem de fișiere sau director poate fi montat la distanță peste un director local. Montările cascade sunt și ele permise în unele implementări NFS.

NFS este realizat să opereze într-un mediu heterogen de mai multe mașini, SO și arhitecturi de rețele. Specificațiile NFS sunt independente de acestea. Independența este atinsă prin folosirea primitivelor RPC(proceduri la distanță) construite pe baza protocolului XDR(reprezentării externe a datelor). Dacă sistemul constă din mașini heterogene și sisteme de fișiere interfatate corect, tipuri diferite de sisteme de fișiere pot fi montate(mounted) și local și la distanță. Specificările NFS afce diferența dintre serviciile furnizate de un mecanism de montare și de serviciile de acces la fișierelor la distanță. poze pg 442 și 443

12.9.2 Protocolul de montare(Mount Protocol)

Protocolul de instalare stabileste conectiunea logica initiala intre server si client. Operatiile de montare includ numele directorului la distanta si numele serverului pe care se afla. Cererea de montare este mapata intr-un RPC corespunzator si inaintat serverului care se ocupa de montat care ruleaza pe serverul corespunzator. Serverul gestioneaza o lista de export (export list) care specifica sistemele locale de fisiere care sunt exportate pt a fi montate impreuna cu numele serverelor ca pot sa le monteze. Cand un server primeste o cerere de montare care se coformeaza cu lista de export, acesta returneaza clientului un handle al fisierului care are rolul de cheie pentru accese viitoare la fisiere din sistemul montat. Handle-ul fisierului contine toate informatiile de care are nevoie serverul ca sa face distinctie intre fisierele individuale pe care le are. Handle-ul fisierelor consta dintr-un identificator de sistem fisier si un numar care identifica directorul montat in sistemul de fisiere exportat. Serverul mentine o lista cu clientii si directoarele montate corespunzatoare. Operatia de montare schimba doar vederea din partea user-ului si nu afecteaza partea serverului.

12.9.3 Protocolul NFS

Asigura un set de apeluri de proceduri pentru operatii cu fisiere la distanta. Procedurile suporta operatiile urmatoare:

- cautarea unui fisier in cadrul directorului
- citirea unui set de intrari in director
- manipularea linkurilor si a directoarelor
- accesarea atributelor fisierelor
- citirea si scrierea fisierelor

Serverele NFS nu au stare, fiecare cerere trebuie sa furnizeze un set complet de argumente. Serverele nu gestioneaza informatia despre clientii lui de la un acces la altul. Datele modificate trebuie salvate pe discul serverului inainte ca rezultatele sa fie returnate clientului. Protocolul NFS nu furnizeaza un mecanism de control al concurentei.

Interfata de fisiere UNIX este bazata pe apelurile open, read, write, close si pe descriptorii de fisier. Layer-ul sistemului de fisiere virtual (VFS) distinge intre fisiere locale si la distanta; fisierele locale se disting dupa tipul sistemului de fisiere. VFS activeaza operatii specifice sistemului de fisiere care se ocupa de cererile locale conform tipului sistemului de fisiere. VFS apeleaza proceduri ale protocolului NFS pentru cereri la distanta. O cerere RPC este facuta nivelului serviciu NFS al serverului la distanta. Poza pg 446

12.9.4 Path-Name Translation

Translatia numelui caili se face prin impartirea caili in componente si executarii unui apel separat lookup call pentru fiecare pereche de componente si directoare vnode. Pt ca aceasta cautare sa fie rapida, un cache implementat la client retine nodurile pentru numele directoarelor la distanta. Acest cache mareste viteza de referire la fisiere cu numele caili initiale identice. Cache-ul este indepartat cand atributurile returnate de la sever nu se potrivesc cu cele ale nodurilor din cache.

12.9.5 Operatii la distanta

Exista o corespondenta aproape de 1 la 1 intre apelurile de sistem uzuale UNIX si protocoalele NFS de RPC (cu exceptia deschiderii si inchiderii fisierelor). Nu exista corespondenta directa intre operatii la distanta si un RPC. Blocurile sunt adunate si introduse intr-un cache local. Doau tipuri de cache:

File-blocks cache: cand un fisier este deschis, kernelul controleaza daca atributele din cache trebuie aduse sau revalidate. Blocurile cu fisiere cahced sunt folosite doar odata daca atributele corespunzatoare din cache sunt actuale.

File-attribute cache: Cache-ul cu attribute este actualizat de cate ori apar noi attribute de la server.

Clientii nu elibereaza blocurile intarziate de scriere decat cand serverul confirma ca data a fost scrisa pe disc.

Cap 17

Coordonarea in sistemele distribuite

17.1 Ordonarea evenimentelor

In sistemele distribuite nu avem un ceas comun sau memorie comuna.

17.1.1 Happened Before Relation

Din moment ce consideram doar procese secventiale rezulta ca toate evenimentele dintr-un process sunt ordonate.

A->B inseamna ca A se executa inainte de B. De exemplu un mesaj de raspuns este conditionat de un mesaj cerere. + figura pag 597.

17.1. 2 Implementare

Fiecare proces are asociat un ceas logic P_i are asociat LC_i . Acest ceas este un simplu contor care se poate doar incrementa. El va desemna un numar unic fiecarui eveniment din interiorul procesului.

Ordonarea la nivel global se face in felul urmatoar. Daca un P_j primeste un mesaj A de la un P_i si $LC_i(A) > LC_j$ atunci $LC_j = LC_j + 1$ in caz contrar nu se va intampla nimic.

17.2. Excluderea mutuala

17.2.1 Principiul centralizat

Exista un singur process ales care va coordona intrarea celorlalte procese in sectiunea critica. El va decide care dintre procese va intra in aceasta sectiune dupa urmatorul algoritm:

Un process P_i care vrea sa intre in sectiunea critica va trimite o cerere catre procesul central P_c , acest va verifica daca alt process este in sectiunea critica exista 2 cazuri: 1. nici un proces nu este in SC => va trimite un reply acordand permisiunea procesului P_i . La terminarea sectiunii critice P_i va trimite un mesaj de elibera catre P_c .

2. in caz contrar P_i este pus intr-o coada de asteptare

La receptia unui mesaj de elibarare P_c va scoate din coada primul process si ii va trimite un mesaj de permisiune.

Daca P_c cade atunci un alt proces este ales P_c .

17.2.2 Principiul descentralizat

Nu exista un process central. Cand procesul Pi vrea sa intre in sectiunea critica el va trimite o cerere (Pi, TS) catre toate procesele existente. Daca va primi un reply de la toate procesele atunci inseamna ca nici unul nu este in sectiunea critica si deci are permisiunea sa intre in ea. In cazul in care un process este in sectiunea critica el nu va trimite reply si va pune intr-o coada cererea procesului Pi.

Mai exista cazul in care 2 procese vor sa intre in acelasi timp in SC. Atunci ele compara marca de timp TS, cea cu marca de timp cea mai mica va avea prioritate.

La intrarea unui nou process numele acestuia trebuie trimis tuturor proceselor existente iar el trebuie sa primeasca numele proceselor existente.

17.2.3 Token Passing

Are la baza o topologie logica de tip inel . Un token va trece pe rand pe la fiecare proces. Un proces care are token-ul va putea sa intre in SC, daca nu vrea sa intre in SC va ceda tokenul urmatorului proces. Exista 2 cazuri defavorabile daca tokenul este pierdut atunci trebuie creat altul nou , sau daca un process cade atunci trebuie redesenata topologia .

17.3 Atomicitate

Atomicitatea= o unitate de program care trebuie executata in totalitate sau neexecutata deloc.

Un coordonator de tranzactie intr-un sistem distribuit trebuie sa aiba grija sa pastreze atomicitatea tranzactiei..

17.3.1 Two Phase Commit Protocol (2PC)

T- tranzactie initiate pe procesorul Si cu coordonatorul de process Ci.

Prima faza: Ci adauga in logul sau <pregateste T>. Ci transmite la toate procesoarele care au executat pe T , mesajul pregateste (T) . La receptia acestui mesaj managerul de tranzactie al statiei va decide daca este dispus sa comita sau nu T. Daca nu vrea sa comita T atunci va adauga in log <no T> si va trimite mesajul abort(T), daca este pregatit sa comita T atunci va adauga in log <ready T> si va trimite mesajul ready(T).

Faza 2: Ci primeste mesajele raspuns. Daca cel putin un mesaj este abort(T) atunci va concludia ca T este <abort T> si va pune in log randul acesta. In caz contrar va concludia comiterea lui T si va adauga in log <commit T>. Dup ace a luat decizia Ci trimite mesajul abort(T) sau commit(T). Statile vor stoca acest mesaj in logurile lor.

17.3.2 Erori in 2PC

17.3.2.1 Eroare la nivelul unei statii participante

Logul statie contine <commit T> => statia executa redo(T). (nu a terminat sa efectueze tranzactia complet si o va lua de la capat)

Logul contine <abort T> => undo T

Logul contine <ready T> Statia nu stie daca T este comisa sau nu. Asa ca va cere aceasta informatie la Ci printr-un query. Daca Ci a picat si el el va cere aceasta informatie tuturor statiilor din retea, daca nici un site nu are aceasta informatie el va astepta nestiind ce decizie sa ia nu va elimina aceasta linie din log. Ci va avea intotdeauna in logul lui decizia luata cu privire la T deci la revenirea lui Si statia va stii care a fost decizia.

17.3.2.2 Eroare la nivelul lui Si

Daca o statie active contine <commit T> atunci T trebuie sa fie comisa tuturor statiilor

Daca o statie active contine <abort T> atunci T trebuie abort la nivelul tuturor statiilor.

Daca niste statii active nu contin <ready T> atunci T trebuie abort. tuturor statiilor

Daca toate statiile contin <ready T> atunci nici o decizie nu poate fi luata. Statiile vor trebui sa astepte revenirea lui Si. Intre timp este posibil ca anumite resurse de la nivelul statiilor sa fie blocate in asteptarea unei decizii.

17.4. Controlul Concurentei

17.4.1 Locking Protocols

Se refera la date si algoritmi partial modificati din cap 7.9

17.4.2 Marci de timp

Fiecare tranzactie are atribuit un numar unic astfel se poate decide ordinea de executie.

17.4.2.1 Generalizarea marci de timp unice

Fiecare statie genereaza o marca de timp locala unica. Marca globala este obtinuta prin concatenarea la aceasta marca a indentificatorului statiei care este de asemenea unic (!!!acesta se pune la final pe cei mai putini semnificativi biti). Ajustarea marcilor de timp locale se poate face dupa principiul ceasurilor logice discutate mai sus.

17.4.2.2 Schema de ordonare cu marci de timp

Combina schma centraliza a marcilor de timp cu 2PC pentru a realize un protocol care asigura serializarea si inlatura roll back cascading.

se refera la date si algoritmi partial modificari din cap 7.9

17.5 Deadlock handling

17.5.1 Deadlock prevention

Resource ordering-prevention deadlock: fiecare resursa primeste un numar unic. UN process poate cere o reursa cu numarul i doar daca nu retine o resursa cu un numar mai mare ca i.

Banker algorithm: alege un singur process din sistem care va implementa acest algoritm.

1. Fiecare process are un numar de prioritate.
2. Numerele de prioritate sunt utilizate ca sa se decida daca un process va astepta dupa altul
3. Elimina deadlockurile.

Pentru controlul Preemption (?) fiecare process va avea o prioritate. Se stabileste ca Pi va astepta dup Pj doar daca prioritatea sa este mai mare decat a lui Pj in caz contrar Pi va fi roll back. Desi aceasta schema previne deadlockurile ea poate duce la fenomenul de infometare. Pt eliminarea acestui se combina aceasta tehnica cu marcile de timp.

1. wait –die scheme. Daca Pi cere o resursa detinuta momentan de Pj el va fi pus sa astepte doar daca marca lui de timp este mai mica decat a lui Pj in caz contrar va muri(roll back).
2. wound-wait scheme. Este invers ca la 1 daca Pi are marca de timp mai mare va fi lasat sa astepte.

Fenomenul de informatare se poate elimina prin conditia ca atunci cand un process moare sau ii se alocă alta marca de timp

17.5.2 Deadlock detection

Fiecare statie mentine un graf local wait-for. Nodurile din graf corespund tuturor proceselor locale si nonlocale care mentin sau cer o anumita resursa a statiei. Fig pag 612. Daca graful are un ciclu atunci avem un deadlock. Se demonstreaza ca nu este neaparat necesar ca un graf local sa aiba un ciclu pt a exista un deadlock. Detectia completa a unui deadlock este facuta doar in momentul unirii tuturor grafurile locale intrun graf global.

17.5.2.1 Metoda centralizata

Un singur process va retine graful global acest process se numeste process coordinator de deadlock detection. Din moment ce exista intarzieri de comunicatie acest graf nu este unul in timp real. Daca se detecteaza un deadlock atunci statiile vor fi anuntate si se va intra in starea de deadlock.

Moduri de obtinerea a grafului: update cand o noua muchie este inserata.

Periodic cand se intampla un numar de schimbari in graf

Oricand processor coordinator are nevoie sa invoce un ciclu de detectie.

Primele 2 modalitati pot induce deadlockuri false datorita intarzierilor din retea.

Ultima modalitate este implementata in asa fel incat aceste deadlockuri false sunt eliminate. (explicatie pag 615)

17.5.2.2 Metoda descentralizata

Se mentin mai multe grafuri parțiale la nivelul fiecare statii. Se pleaca de la idea ca daca exista un deadlock atunci el va determina cel puțin un ciclu într-unul din grafurile existente. (pag 616)

17.6 Algoritmul de alegere

In multe implementari din acest capitol un process trebuie ales pentru a coordona o activitate. Se pleaca de la principiul ca fiecare process un numar de prioritate unic (P_i are prioritate i) Procesul coordinator va fi cel cu numarul cel mai mare.

17.6.1 Bully Algorithm

P_i trimite un mesaj de alegere tuturor proceselor cu prioritate mai mare ca a sa si asteapta un raspuns intr-un interval de timp T . Daca nu primeste nici un raspuns se alege pe el process coordinator si trimite un mesaj de instiintare celorlalte procese. Daca primeste un raspuns in acest interval e va astepta un interval de timp T_1 pentru a primi un mesaj de la procesul coordinator. Daca nu primeste acest mesaj atunci el va relua algoritmul de alegere presupand ca procesul cu prioritate mai mare a cazut.

17.6.2 Ring algorithm

Se foloseste de listele active. Aceste liste contin numerele de prioritate ale tuturor proceselor active in momentul in care algoritmul se termina. Fiecare process are propria lista.

Algoritmul :

Daca P_i determina caderea procesului coordinator creaza o noua lista initial goala. Apoi trimite mesajul elect(i) la vecinul din dreapta.

Daca P_i primeste un mesaj de la vecinul din stanga $elect(j)$ avem 3 posibilitati:

1. Este primul mesaj primit. Atunci creazao noua lista adauga numerele I si j si trimtie la vecinul din dreapta $elect(j)$ si $elect(i)$.
2. Daca $i \neq j$ atunci inseamna ca vecinul din stanga nu are numarul lui P_i el va adauga I in lista si va trimtie mai departe lista.
3. daca $i = j$ atunci s-a ajuns la capatul cercului. Si P_i poate sa determine cel mai mare numar din lista si sa identifice procesul coordinator.

17.7 Ajungerea la intelegere

Comunicatie ureliable

Dup ace P_i trimite un mesaj lui P_j el asteapta un timp t sa primeasca un mesaj de confirmare de primire. Daca nu va primi acest mesaj el va retrimite mesajul initial pana cand va primi confirmarea sau va fi notificat de sistem ca statia S_j este cazuta. Daca P_j asteapta si el o confirmare de la P_i ca a primit replay se poate concluziona ca cele 2 procese ajung la o intelegere

17.7.2 Faulty

Presupunem ca mediul de comunicatie este viabil dar procesele pot sa cada. Avem n procese din care doar m sunt supuse esecului. Problema se pune cum determinam pe cele m care pot da rateuri .

Fiecarui process i se asociaza o valoare V_i private. Algoritmu va crea pentru fiecare process nonfualty un vector $X_i = \{A_{i1}, A_{i2}, .. A_{in}\}$ cu prorprietatile

1. daca P_j nonfaulty atunci $A_{i,j} = V_j$
2. daca P_i nonfaulty atunci $X_i = X_j$.

+ explicatie pentru $m=1$ si $n=4$ pag 623.