

## ELEMENTE DE GESTIUNE A MEMORIEI

De regulă la microprocesoarele de 8 biți, adresa generată și depusă pe magistrala externă de adrese este folosită direct în adresarea unei locații de memorie sau a unui dispozitiv de I/O. Începând cu microprocesoarele de 16 biți apar noțiunile de adresă logică și adresă fizică. Adresele generate de un program sunt considerate *adrese logice* și totalitatea acestora formează *spațiul adreselor logice*. Totalitatea adreselor ce corespund memoriei (și/sau dispozitivelor de I/O) formează *spațiul adreselor fizice*. Cele 2 spații, al adreselor logice și fizice, pot să fie egale sau inegale. Ca urmare, trebuie să existe un mecanism de conversie a adreselor, de translatare din adrese logice în adrese fizice. La microprocesoarele Intel din seria x86 mecanismul de translatare este inclus pe același substrat de siliciu cu microprocesorul. La unele din microprocesoare (ca de exemplu Motorola MC68000) mecanismul de translatare din adresele logice în cele fizice este preluat de un circuit extern microprocesorului denumit "unitate de administrare a memoriei" (MMU - Memory Management Unit).

În general, tehnicile folosite pentru lucrul cu sistemul de memorie într-un calculator, fie că este vorba de selecție, adresare, alocare de spațiu în memoria principală, folosesc un ansamblu de resurse hardware și software, operațiile executate de acestea fiind incluse în noțiunea generală de *management al memoriei*. Administrarea întregii memorii (de pe toate nivelurile ierarhice) a calculatorului, incluzând și modalitățile de partiționare logică și adresare în spațiul adreselor logice (*L*) și în spațiul de memorare (*M*, adrese fizice) se face cu ajutorul unei unități hardware pe care o notăm în continuare MMU (Figura 1).

Când un program trebuie rulat (lansat în execuție), el este mai întâi adus din memoria auxiliară în memoria principală (în întregime sau parțial).

Pentru MMU se impun următoarele cerințe de bază:

1. să realizeze translatarea adreselor și să susțină alocarea dinamică a memoriei

Înainte ca un program să fie executat, este necesar ca acestuia să i se aloce un spațiu în memoria principală (spațiu de adrese în memoria fizic existentă în calculator). Dacă această atribuire de spațiu (adrese) se face doar la încărcarea programului în memoria principală, procedeul este numit *alocare statică* și prezintă dezavantajul că adresele de încărcare sunt fixe. *Alocarea dinamică* a memoriei se caracterizează prin atribuirea spațiului necesar programelor, sau proceselor, în timpul execuției acestora. Adresa la care se încarcă programul în memoria principală nu mai este fixă, ea putând fi modificată, la momente de timp diferite.

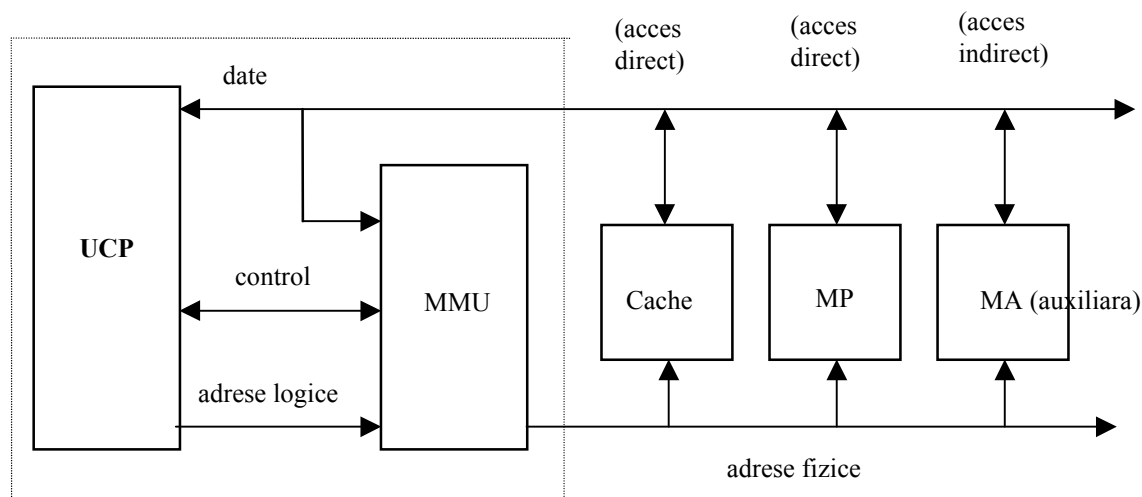


Figura 1. Organizarea ierarhică a memoriei

Se folosesc trei procedee de translatare a adreselor logice în adrese fizice toate bazându-se pe o divizare a spațiului adreselor logice în blocuri de adrese logice continue, după cum urmează:

- *segmente* de lungime arbitrară;
- *pagini* (în spațiul adreselor logice) și *cadre-pagină* / *blocuri* (în cadrul adreselor fizice), de lungime fixă;
- combinații de adresare *segmentat-paginată*.

Mecanismul de traducere, indiferent că se folosesc pagini sau segmente, se produce în faza de execuție a programului și include printre altele, un *tabel de traducere* prin care se face corespondența între adrese logice și adrese fizice (principiul este prezentat în fig. 2).

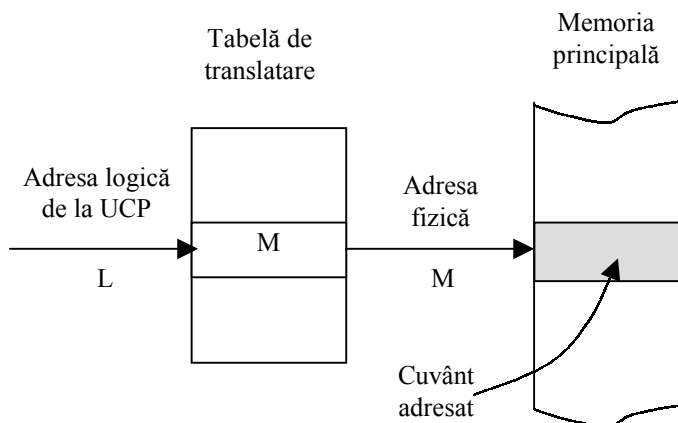


Figura 2. Principiul tehnicilor de traducere a adreselor

De la microprocesor se recepționează adresa logică L. Dacă elementul adresat nu se găsește în memoria principală, atunci registrul (de adresă) L va conține adresa din memoria externă unde se găsește elementul respectiv și el va fi adus de rutine ale sistemului de operare în memoria principală, la o adresă fizică a cărei valoare M se înscrie în registrul de adresă L din tabela de traducere. Dacă elementul adresat este în memoria principală rezultă că adresa fizică este M. Acest mecanism este totuși inefficient, căci tabela de traducere are dimensiunile egale cu ale memoriei principale. Pentru micșorarea dimensiunii tabelor se folosește divizarea spațiului logic în segmente / pagini, iar referirile se fac doar la nivelul de adresă logică de bloc. Trebuie menționat că dacă sunt lansate două procese în execuție, fiecare cu spațiul său de adrese logice, există două tabele de traducere spre memoria principală. Sistemul de operare va reactualiza conținutul tabelii de traducere ori de câte ori va avea loc o relocare (schimbare de poziție în spațiul adreselor fizice) a programelor în memorie.

2. *să susțină mecanismele de implementare ale memoriei virtuale*

Mecanismul memoriei virtuale permite programelor să fie executate chiar dacă numai o parte din instrucțiunile lor se încarcă în memoria principală, iar restul se găsesc în memoria auxiliară de tip disc magnetic. Adresele logice sunt numite adrese virtuale, în cazul implementării memoriei virtuale, pentru că nu există fizic în memoria principală ci ele se găsesc în memoria auxiliară.

3. *să poată furniza protecție și securitate pentru memorie*

Cerința se referă la protecția informației, la alocarea de drepturi de acces la zone din memoria principală și la zone din spațiul adreselor logice, (protecție la citire, scriere, ștergere, copiere, execuție). Se protejează memoria disponibilă alocată programelor de sistem și programelor utilizator, asigurându-se securitatea informației (acces limitat la informații), prin controlul accesului la resursele sistemului.

Așa cum s-a pomenit și mai sus există două strategii folosite pentru implementarea MMU:

- a) MMU se află pe chip-ul procesorului; se aplică la microprocesoarele cu spațiul segmentat al adreselor logice. De exemplu: Intel x86, Pentium, Zilog Z800, Z8000)
- b) MMU este construită separat de UCP, variantă aleasă la microprocesoarele cu spațiu liniar (organizare liniară a) al adreselor logice. De exemplu: MC 680x0, Z8001).

În cazul în care organizarea logică a memoriei este liniară, adresele încep în mod obișnuit de la 0 și avansează liniar. Memoria poate fi apoi structurată prin software, la nivelul translației de adrese. În cazul în care organizarea memoriei este segmentată, programele nu sunt scrise ca secvențe liniare de instrucțiuni și date, ci ca segmente separate de cod, date, sau stivă. Spațiul adreselor logice este spart în mai multe spații cu adresare liniară, fiecare cu lungime diferită, definită de compilator sau de programator. O adresă logică efectivă este calculată ca o combinație între numărul segmentului (care indică adresa de bază a unui bloc de memorie) și un deplasament (offset) în cadrul segmentului.

În general schemele de adresare liniară sunt mai potrivite pentru manipularea structurilor mari de date, în timp ce segmentarea facilitează programarea, permițând programatorului să structureze programele în module segment. În plus adresarea segmentată simplifică protecția și relocarea obiectelor în memorie. Segmentarea facilitează gestiunea memoriei în sistemele de calcul multiuser, multitasking.

### TRANSLATAREA ADRESELOR

Indiferent de schema de organizare a memoriei (liniară sau segmentată) procesorul trebuie să aibă un mecanism de traducere a adreselor, util în implementarea memoriei virtuale. Acest mecanism este de asemenea util pentru protejarea informației din memorie. Traducerea de adrese este un proces atribuire și organizare ("mapare") a adreselor logice în adrese fizice de memorie. Mecanismul de traducere împarte memoria principală în *blocuri (cadre pagină)*. Așa cum am amintit mai sus se folosesc 2 scheme de traducere:

- (1) traducere prin paginare
- (2) traducere prin segmentare

În *sistemele paginate*, memoria principală este împărțită în blocuri de *lungime fixă* în timp ce în *sistemele segmentate* blocurile sunt de *lungime variabilă*.

Paginile au în general lungimi de ordinul 256 - 4096 cuvinte, în timp ce segmentele cu lungimea definită de compilator sau programator au 64 K-cuvinte, sau mai mult. În sistemele cu multiprogramare și time-sharing<sup>1</sup>, mai mulți utilizatori folosesc aceleași programe cum sunt editoare, compilatoare, programe utilitare, biblioteci de programe etc. Atât sistemele cu paginare cât și cele cu segmentare permit mecanisme de partajare, între procesele utilizator, a paginilor, respectiv a segmentelor. Aceste mecanisme se bazează pe intrări în tabele de mapare de pagină (sau segment) în care intrările diferitelor procese indică către același bloc din memoria principală. Combinația între segmente și pagini presupune că un segment conține una sau mai multe pagini virtuale. Mecanismul de segmentare administrează spațiul virtual, împărțind programele în segmente, în timp ce paginarea este destinată administrării memoriei fizice care este împărțită în cadre pagina (blocuri). Vom descrie în continuare mecanismele folosite pentru maparea segmentelor și paginilor.

#### Maparea adreselor folosind pagini

În cazul paginării, implementarea tabelului pentru maparea adreselor este simplă pentru că informația din spațiile de adresare și memorare este divizată în blocuri de dimensiune fixă. Memoria fizică este împărțită logic în blocuri de aceeași dimensiune (64-4096 cuvinte fiecare). Termenul *pagină* se refera la blocuri, *de aceeași dimensiune*, de adrese din spațiul de adresare. De exemplu, presupunem un calculator care are 20 de biți de adresă și folosește doar 32KB în memoria principală (memorie implementată fizic prin circuite de memorie); dacă o pagină, respectiv un bloc, au dimensiuni egale cu 1 kB, atunci spațiul de adrese e divizat în 1024 pagini, iar spațiul de memorie e divizat în 32 de blocuri.

Se consideră că și programele sunt împărțite în pagini. Porțiuni din programe sunt mutate din memoria auxiliară în memoria principală în înregistrări egale cu mărimea paginii. În loc de bloc de memorie e folosit uneori termenul de cadru pagină ("page frame").

În cazul mapării prin paginare adresa virtuală este reprezentată printr-un singur cuvânt de adresă împărțit într-un câmp corespunzător numărului paginii (adresa paginii) și un câmp pentru deplasament. La maparea prin segmentarea un singur cuvânt de adresă nu mai este suficient; dimensiunea variabilă a segmentelor conduce la existența a două cuvinte de adresare, în care primul indică numărul (adresa) segmentului, iar cel de-al doilea deplasamentul în cadrul segmentului.

Într-un calculator cu  $2^p$  cuvinte pe pagină,  $p$  biți sunt folosiți pentru a specifica o adresă de linie iar cei mai semnificativi biți rămași în adresa virtuală specifică numărul de pagină. Fie de exemplu un sistem simplu cu o adresă virtuală cu dimensiunea de 16 biți și pagini cu 4096 cuvinte. Pentru că o pagină are  $2^{12}$  cuvinte, cei patru biți mai semnificativi vor specifica una din cele 16 pagini, iar cei 12 biți mai puțin semnificativi indică adresa liniei în cadrul paginii.

---

<sup>1</sup> time-sharing = partajarea timpului UCP

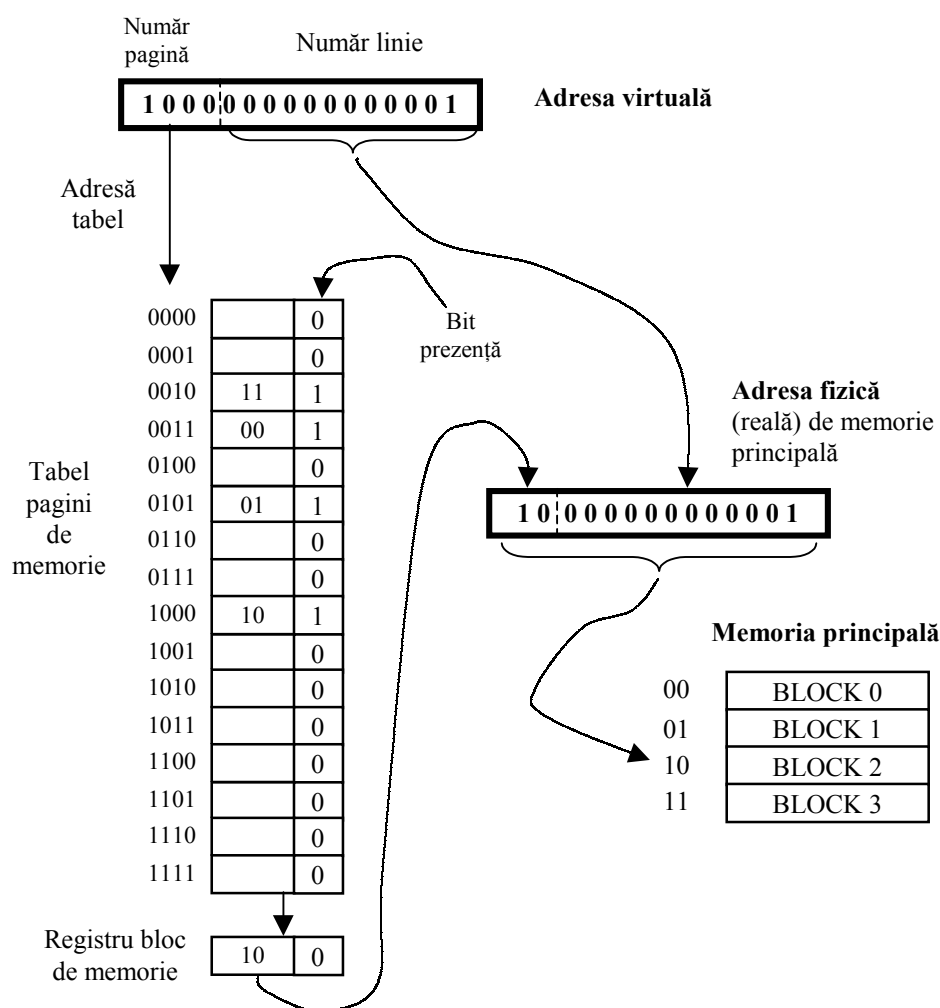


Figura 3. Paginarea memoriei

Ca urmare maparea trebuie făcută doar de la un număr de pagină la un număr de bloc din memoria principală, pentru că adresa liniei e aceeași pentru ambele spații. Organizarea de principiu a tabelului de mapare a memoriei într-un sistem paginat poate arăta ca în figura 3.

În tabelul paginii de memorie, adresa conduce la numărul paginii, iar conținutul indică numărul blocului unde pagina este stocată în memoria principală. Se observă că paginile virtuale 2, 3, 5 și 8 se află în memoria principală, în blocurile 3, 0, 1 și respectiv 2. Un bit de prezență adăugat fiecărei locații a tabelului indică dacă pagina respectivă a fost mapată în memoria principală (i s-a alocat spațiu și a fost transferată din memoria auxiliară în memoria principală). Valoarea 0 a bitului de prezență indică că pagina nu este în memoria principală. În exemplul dat, cei 4 biți ai adresei virtuale specifică atât numărul paginii cât și adresa în tabelul paginilor de memorie. Conținutul tabelului este citit în registrul ce indică blocul de memorie. Dacă bitul de prezență este 1, numărul de bloc e transferat în registrul ce păstrează adresa fizică din memoria principală. Dacă bitul de prezență este 0 rezultă că adresa virtuală se referă la un articol ce nu se găsește în memoria principală. În acest caz, se generează o cerere (în forma unei întreruperi software, de tip “eroare de pagină<sup>2</sup>”) către sistemul de operare (SO) pentru a aduce pagina cerută din memoria auxiliară în memoria principală, înainte de reluarea programului care a accesat adresa din pagina respectivă.

O astfel de organizare în memoria principală a tabelului de traducere este însă inefficientă, căci tabelul are multe locații goale. Dacă presupunem un calculator cu spațiul de adrese de 4 GB ( $2^{32}$ ), cu pagini de memorie de 32 KB ( $2^{15}$ ), iar spațiul memoriei fizice este 32 MB ( $2^{25}$ ) rezultă 128 K-pagini virtuale de memorie și 1024 blocuri (1 K-blocuri). Conform principiului din figura 6.24 tabelul de traducere ar avea

<sup>2</sup> page fault = eroare de pagină

131071 ( $2^{17} = 128 \text{ K}$ ) linii și doar 1024 locații ar fi ocupate (cu bit de prezență 1), restul fiind goale, neutilizate. O soluție a acestei probleme o constituie folosirea unei memorii asociative, ca în figura 4, (*tabel asociativ de pagini de memorie*) cu număr de locații egal cu numărul de blocuri din memoria principală.

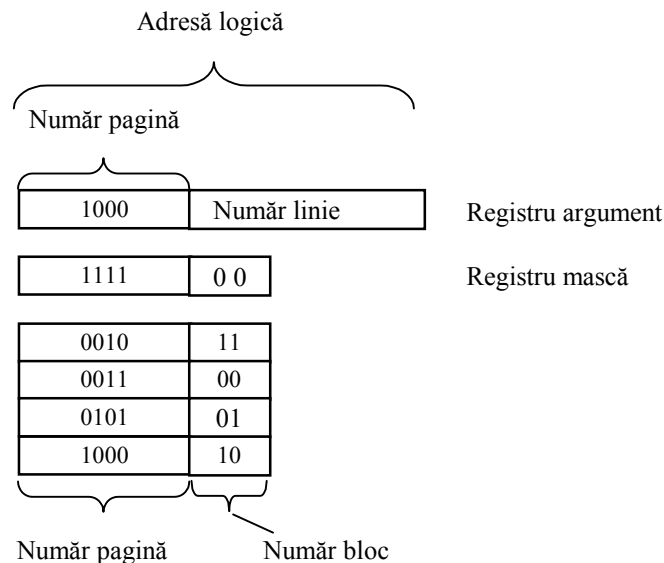


Figura 4. Traducerea adresei de pagină cu memorie adresabilă prin conținut

În fiecare locație de memorie se stochează numărul paginii virtuale și numărul corespunzător al blocului de memorie principală alocat.

Adresa logică se compară combinațional cu conținutul tabelului asociativ, doar pentru pozițiile binare ce corespund la 1 logic în registrul de mascare, deci în figura 4, doar pentru câmpul de adresă al numărului de pagină virtuală. Dacă la apelarea unei pagini virtuale, aceasta nu este mapată în memoria principală, operația de adresare este o rată (miss) și evenimentul este rezolvat printr-o rutină de servire a sistemului de operare care aduce pagina respectivă din memoria auxiliară. Și acest tabel asociativ de traducere poate avea dimensiuni mari uneori. De aceea se folosesc tabele asociative de mare viteză la care numărul liniilor este mai mic decât numărul de blocuri (cadre pagină) și în care se salvează numai o parte din mapările realizate. Acest registru asociativ este numit TLB (Translation lookaside buffer) și este descris mai jos.

### Mapare prin segmentare

La acest tip de mapare spațiul adreselor logice este împărțit în *segmente*, care spre deosebire de pagini, nu au o dimensiune fixă. Segmentele sunt spații de adrese continue (liniare), care includ programe sau porțiuni din programe cu atribute comune. Mecanismul de traducere este prezentat în schema de principiu din figura 5.

Spre deosebire de operația de paginare, la segmentare nu se mai face concatenarea (alăturarea) adresei de deplasament (din adresa logică) la adresa fizică. Aici adresa de deplasament este adunată cu adresa fizică a segmentului (adresa sa de bază) citită din tabelul de traducere al segmentelor, specific unui anumit proces.

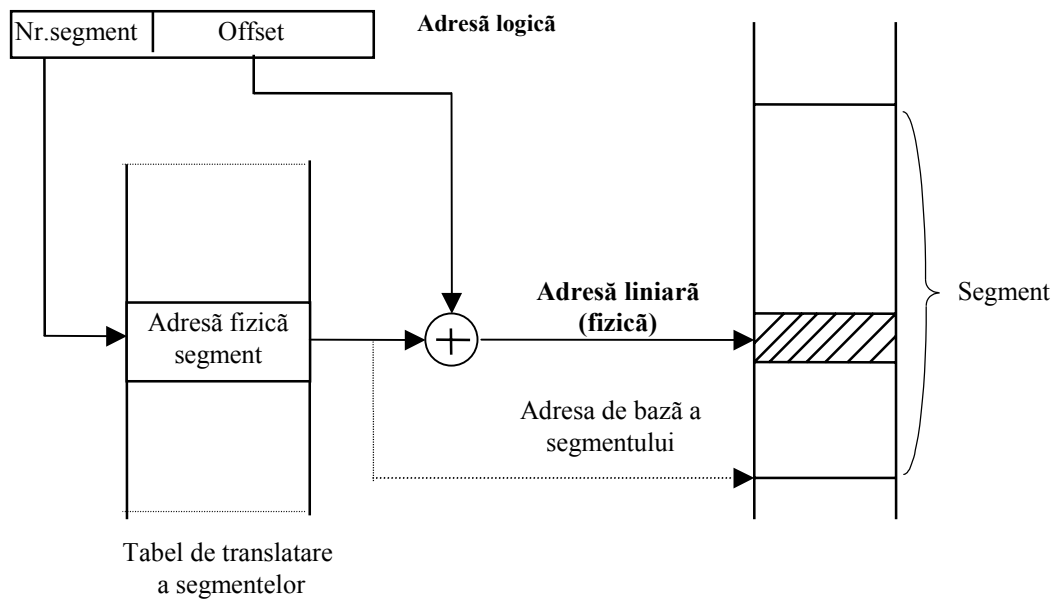


Figura 5: Segmentarea memoriei

## Noțiuni privind protecția memoriei

Inventarea multiprogramării, la care un procesor este partajat între mai multe programe aflate în execuție, a condus la noi cerințe privind protecția între programe. Mecanismele de protecție sunt strâns legate de implementarea memoriei virtuale.

Multiprogramarea conduce la conceptul de proces. Un proces este constituit dintr-un program aflat în execuție plus orice stare necesară pentru continuarea rulării sale. Partajarea timpului ("time-sharing") este o variantă a multi-programării care partajează UCP și memoria între diferiți utilizatori simultani, iar mecanismul de partajare dă iluzia că toți utilizatorii au propriile mașini. Astfel că, în oricare dintre variante, trebuie să fie posibil să se comute de la un proces la altul. Această comutare este numită *comutare de proces* sau *comutare de context*.

Un proces trebuie să se execute corect indiferent dacă se execută continuu de la start până la terminare, sau dacă este întrerupt în mod repetat pentru a se face comutarea către alte procese. Responsabilitatea menținerii comportării corecte a proceselor cade atât în sarcina resurselor hardware, care asigură salvarea și restaurarea corectă a proceselor ce rulează pe UCP cât și în sarcina sistemului de operare care trebuie să garanteze că un proces nu interferează cu celelalte și că un proces nu poate modifica datele altor procese. Aceasta înseamnă asigurarea protecției datelor de către sistemul de operare. În afară de protecție, sistemul de operare permite, în anumite condiții, și partajarea codului și a datelor între diferite procese pentru salvare de spațiu de memorie prin reducerea numărului de copii de informații identice.

Cel mai simplu mecanism de protecție este constituit de o pereche de registre care verifică fiecare adresă, astfel încât să nu se permită accesul în afara unor limite alocate în spațiul de adrese. Aceste registre sunt numite registru - bază și registru limită. O adresă este validă dacă se încadrează între cele două adrese conținute în registre:

$$\text{Bază} \leq \text{Adresă} \leq \text{Limită}$$

La unele sisteme adresa este considerată ca un număr fără semn, care se adună întotdeauna la o adresă de bază, astfel că testul de adresă limită se reduce la:

$$(\text{Bază} + \text{Adresă}) \leq \text{Limită}$$

Dacă proceselor utilizator li s-ar permite să modifice registrele de adrese de bază și limită, atunci acest mecanism de protecție nu ar putea funcționa. Modificarea respectivelor registre este permisă doar sistemului de operare, pentru asigurarea protecției între procese.

Pentru asigurarea protecției sistemul de operare (SO) are trei responsabilități principale:

1. Să furnizeze cel puțin două moduri de execuție indicând că procesul aflat în rulare este un proces utilizator, sau un proces al sistemului de operare. La diferite SO ultimul tip de proces este numit în diferite feluri: *proces kernel* (nucleu), *proces supervisor*, sau *proces executiv*.
2. Să prevadă o porțiune a stării CPU pe care procesul utilizator o poate doar citi, dar nu o poate scrie. Aceasta include registrele bază / limită, indicatori (biți) pentru moduri utilizator / supervisor și indicatori pentru validare / invalidare evenimente de tip excepție.

3. Să prevadă un mecanism prin care UCP poate trece din mod utilizator (user) în mod supervizor și viceversa. Primul sens de trecere este specific apelurilor sistem ("system calls" – apeluri de servicii oferite de SO), implementate ca instrucțiuni speciale care transferă temporar controlul la o locație precisă din spațiul de cod al supervizorului. Registrul contor de program corespunzător locului unde se face un apel sistem este salvat, iar UCP trece în mod supervizor. Revenirea la modul utilizator este similară cu cea produsă la o revenire din procedură.

Adresele de bază și limită constituie minimul unui sistem de protecție. Mecanismul memoriei virtuale oferă alternative mai complexe și mai sigure decât modelul simplu prin bază și limită. Așa cum s-a văzut, adresele virtuale sunt translatate în adrese fizice pe baza unor tabele de traducere. Acest mod de mapare, prin tabele, oferă posibilitatea introducerii de informații pentru controlul erorilor de program (intenționate sau nu) care încearcă să treacă peste mecanismele de protecție. Cea mai simplă cale este introducerea unor indicatori de permisie pentru fiecare pagină sau segment. De exemplu pot exista indicatori care să permită doar citire, doar execuție, sau care să interzică accesul unui proces utilizator la anumite pagini / segmente. Fiecare proces poate face adresare doar către paginile proprii de memorie, procesul utilizator neavând dreptul să modifice tabelele sale de pagină / segment.

Protecția poate fi extinsă chiar pe mai mult decât două niveluri (nu doar utilizator și supervizor), privite ca și inele concentrice de protecție, în centru găsiindu-se nivelul de protecție cel mai înalt. În această ierarhie de niveluri de protecție un program poate accesa doar date de pe nivelul său de protecție și de pe nivelurile inferioare în ierarhie. Poate face însă operații de apelare (call) a serviciilor sistemului de operare, servicii oferite de rutine ce se află pe niveluri superioare de protecție. Adesea se face comparația, în oarecare măsură nefericită, cu clasificări de tip militar ale sistemelor: top – secret, secret, confidențial și neclasificat. Programele utilizator ("civilii" în exemplul militar) au doar dreptul de acces la nivelul de protecție cel mai de jos: "ne-clasificat. La sistemele de protecție de tip inele concentrice, deosebirea față de exemplul anterior, este că se pot face apelări la rutine situate pe niveluri superioare de protecție, dacă există "chei" de acces către acele niveluri. Poate exista dreptul de apelare a unor servicii ale sistemului de operare, prin mecanismul de comutare a proceselor.

Așa cum am pomenit și mai sus, informația de protecție este setată în registre speciale atașate fiecărei intrări în tabelele de traducere, registre setate doar de rutine de control al sistemului de operare. Drepturile de acces pot fi de tipul:

- a. Atribuirea de privilegii complete de citire și scriere. Aceste drepturi se atribuie programului atunci când execută propriile instrucțiuni.
- b. Read-only (protecție la scriere). Protecția la scriere este utilă la operațiile de partajare a unor rutine de sistem (utilitare, biblioteci, etc.).
- c. Execute only (program protection). Protejează programul la copiere. Restricționează referirea la segment doar în timpul fazei de fetch a instrucțiunii și nu și în timpul fazei de execuție. Asta permite utilizatorului să execute instrucțiunile segmentului de program, dar nu permite citirea instrucțiunilor ca date cu scopul de a copia conținutul lor.

Sistemul de protecție a memoriei, în sensul celor spuse mai sus, se construiește pentru:

- (a) memorie (detectează orice eroare de adresare înainte ca aceasta să creeze erori accidentale sau voite);
- (b) programe (previne ca programele utilizator (de aplicații) să facă modificări ilegale în rutinele SO);
- (c) utilizatori (programele utilizatorilor între ele);
- (d) securitate informație (acces limitat la informațiile unde utilizatorul nu are drept de acces).