

MEMORIA CACHE

Memoria cache a apărut din necesitatea de a elimina diferența de viteză dintre procesorul rapid și memoria sistem, mult mai lentă. Este o memorie statică (SRAM) de dimensiuni mici dar foarte rapidă, folosită pentru a aduce datele mai aproape de unitățile de execuție din procesor. Ca urmare a performanțelor de funcționare a memoriei cache, actualmente doar un procent de aproximativ 10% din accesările la memorie ale procesorului se fac din memoria sistem, iar restul de 90% se fac din memoria cache.

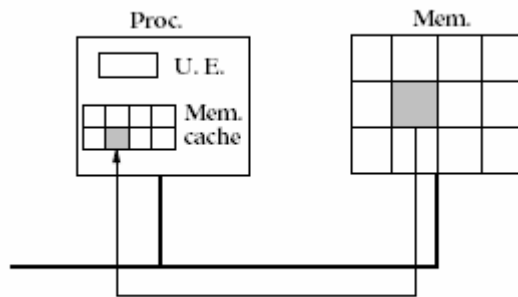


Figura 1. Memoria cache

Memoria cache se bazează pe principiul localității. Avem două tipuri de localitate:

- localitate temporală – locațiile de memorie accesate la un moment dat tind să fie accesate din nou în viitor;
- localitate spațială – datele aflate în vecinătatea unor date accesate tind să fie și ele accesate.

Exemple de memorii cache:

- memoria TLB (*Translation Look-aside Buffer*) este un cache special pentru traducerea între adresele fizice și adresele virtuale din calculator. Datorită folosirii memoriei TLB, calculul adreselor pentru memoria virtuală nu se face la fiecare acces al memoriei, ci se ia din tabela TLB.

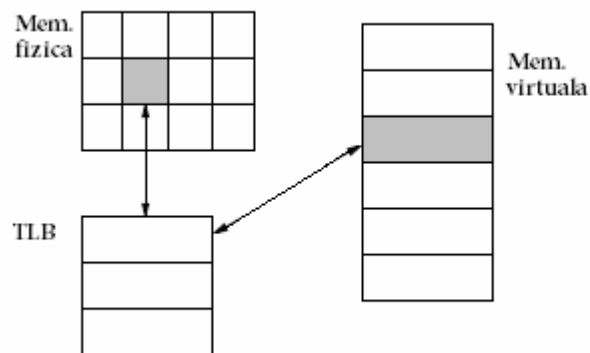


Figura 2. Memoria TLB

- memoria cache pentru instrucțiuni - este folosită pentru încărcarea în avans a instrucțiunilor în procesor ;
- memoria cache de date - folosită pentru încărcarea în avans în datelor necesare din memoria sistem.

Există 3 moduri de mapare a datelor din memoria fizică în memoria cache:

- 1) mapare directă
- 2) mapare cu asociere completă
- 3) mapare cu asociere pe blocuri (set asociativ)

1) Memoria cache cu mapare directă

În acest caz există o mapare directă, secvențială între adresele din memoria principală și locațiile memoriei cache. Aici memoria fizică este împărțită în pagini, fiecare pagină fiind de dimensiunea memoriei cache. Atâtea paginile din memoria sistem cât și memoria cache sunt împărțite în blocuri de aceeași dimensiune.

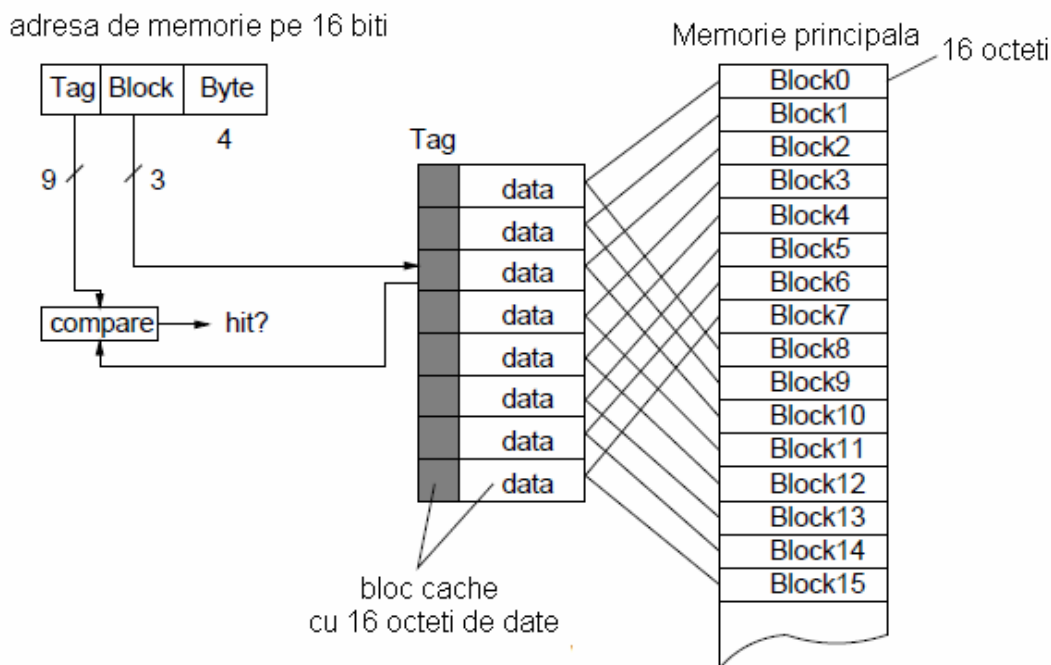


Figura 3. Memoria cache cu mapare directă

În exemplul de mai sus, o pagină conține 8 blocuri; astfel, blocul 0 din memoria cache corespunde cu blocurile 0, 8, 16, etc. din memoria sistem.

Pentru a ști exact cu ce bloc din memorie corespunde un anumit bloc din memoria cache, se folosește un tag (indicator) la fiecare intrare din memoria cache. Tag-ul conține cei mai semnificativi biți ai adresei blocului din memoria principală, reprezentând de fapt adresa de pagină a blocului.

Dacă procesorul vrea să acceseze un bloc de date din memorie situat la o anumită adresă, mai întâi trebuie să verifice dacă blocul cerut se află sau nu în memoria cache. Pentru aceasta, se va proceda astfel:

- a) se selectează o intrare din memoria cache pe baza informației de bloc din adresa cerută;

Tag	Block	Byte
-----	-------	------

- b) se compară câmpul tag din adresa cerută de procesor cu cel din intrarea cache;
- c) dacă cele două sunt egale, avem așa-numitul cache-hit, adică blocul cerut din memoria principală se regăsește în memoria cache;
- d) dacă cele două tag-uri sunt diferite, atunci blocul cerut nu există în memoria cache și el va fi încărcat din memoria principală.

Avantajul modului de mapare direct :

Memoriile cu mapare directă sunt memorii simple și de mare viteză. Ele se pot utiliza ca memorii cache de nivel 2, putând avea dimensiuni mai mari.

2) Memorii cache cu asociere completă

Aici, un bloc din memoria sistem poate fi plasat oriunde în memoria cache.

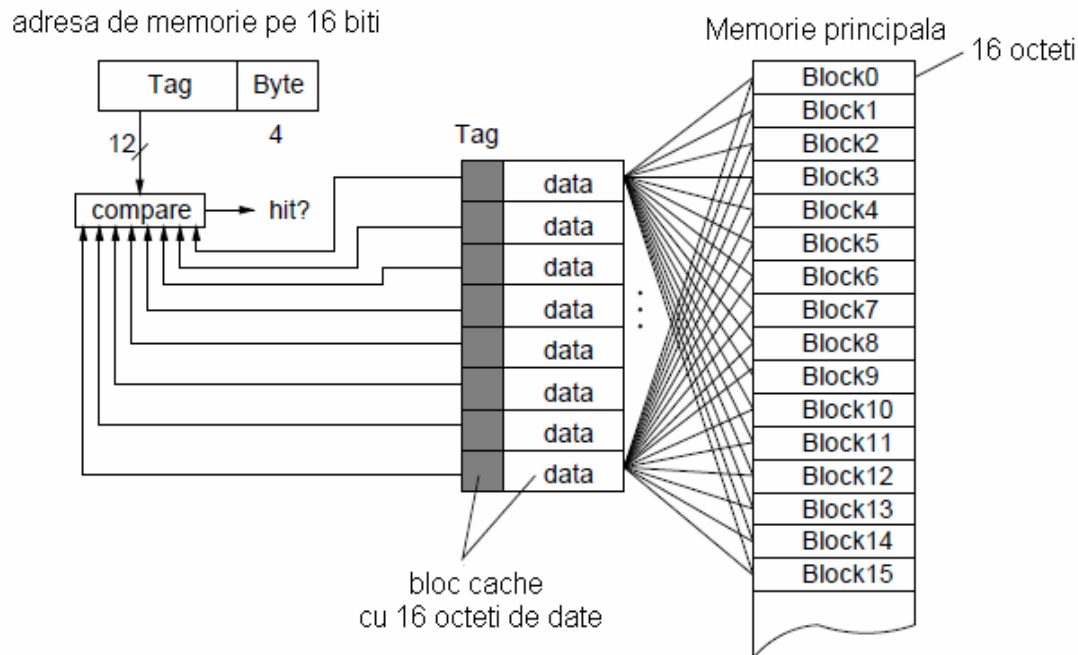


Figura 4. Memoria cache cu asociere directă

Pentru a detecta dacă un bloc din memoria sistem există în memoria cache, se face o comparație în paralel a tag-urilor din toate blocurile memoriei cache. Dacă se identifică un tag, înseamnă că s-a găsit blocul căutat.

Avantajul acestui tip de memorie: folosește un mod de mapare mai flexibil.

Dezavantaj: este mai costisitor de implementat datorită comparațiilor care se fac în paralel.

3) Memorii cache set asociative

Acestea folosesc o combinație între maparea directă și maparea cu asociere completă. Aici, memoria cache este împărțită în mai multe seturi de înțiriri, existând o mapare directă între blocurile din memoria principală și seturile din memoria cache. În interiorul seturilor se face o mapare cu asociere completă.

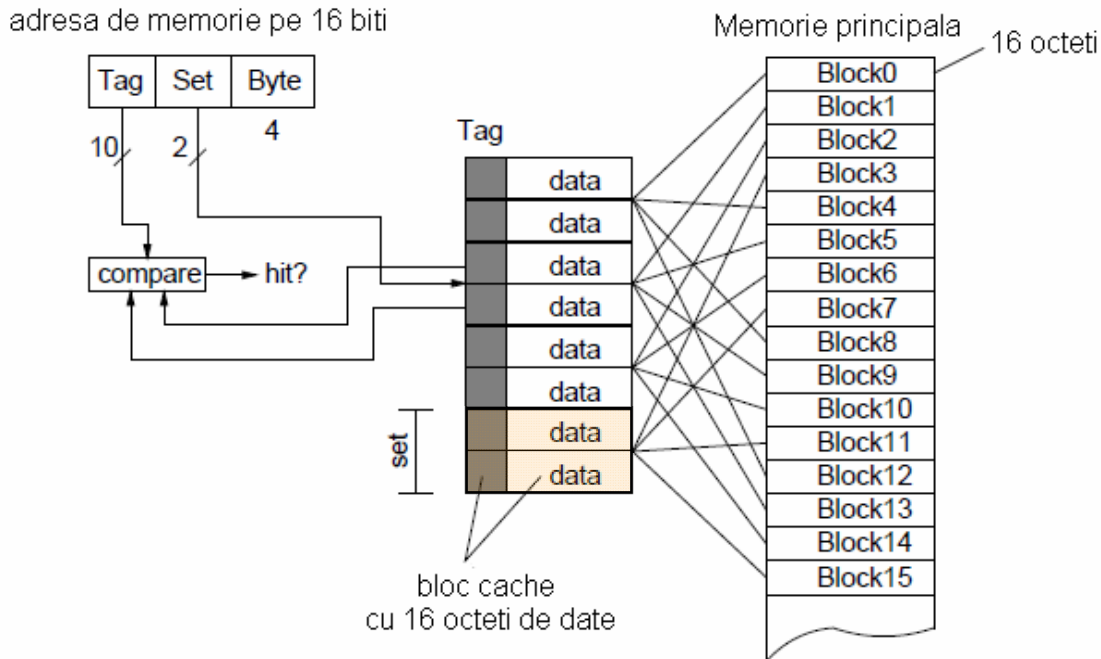


Figura 5. Memoria cache set-asociativ

Pentru a vedea dacă blocul dorit se găsește în memoria cache, se face o comparație în paralel a tag-urilor doar pentru înțirile dintr-un singur set.

Avantajul memoriilor set-asociative: se reduce numărul de comparații efectuate în paralel; sunt mai puțin costisitoare decât memoriile complet asociative; se folosesc ca memorii cache de nivelul 1, având performanțe ridicate.

Observa ie: În sistemele actuale exist o ierarhie de memorii cache organizat pe mai multe nivele; în func ie de pozi ia memoriei cache fa de procesor, putem avea: memorie cache de nivel 1 (care se g se te în interiorul capsulei procesorului), memorie cache de nivel 2, de nivel 3, etc.

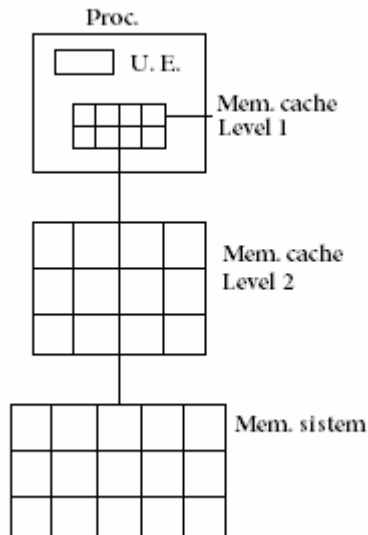


Figura 6. Ierarhia de memorii cache într-un sistem de calcul

Strategii utilizate pentru memoriile cache:

În leg tur cu anumite evenimente ce pot sa apar la accesarea memoriei cache, se stabilesc urm toarele strategii:

S1) Strategia *cache-miss*. Aceast strategie se refer la memorii cache complet asociative. Astfel dac e nevoie de un bloc de memorie care nu se reg se te în cache (*cache-miss*), trebuie s se stabileasc care anume dintre blocurile existente în memoria cache va fi înlocuit de noul bloc care va fi adus din memoria sistem.

Strategia de înlocuire poate fi:

- aleatorie
- *first in, first out* - se înlocuie te primul bloc introdus în cache
- se înlocuie te cel mai pu in recent utilizat. Aici se înlocuie te blocul care nu a fost accesat de cel mai mult timp (cel mai vechi).

S2) Strategii de scriere a memoriilor cache. S presupunem c procesorul efectueaz o opera ie STORE (scriere în memorie). Exist 2 posibilit i:

a) valoarea se scrie în memoria cache și în memoria sistem. Aceste memorii se numesc *write-through*; valoarea respectiv se scrie prin cache în memoria sistemului.

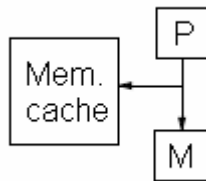


Figura 7. Memoria *write-through*

b) scrierea se face doar în memoria cache. Scrierea în memoria sistem se face ulterior doar când este necesar acest lucru. Aceste memorii se numesc *write-back*.

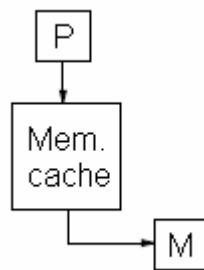


Figura 8. Memoria *write-back*

Avantajele memoriei *write-through*:

- memoria principală este întotdeauna consistentă cu memoria cache.

Dezavantajele memoriei *write-through*:

- această memorie este mai încetă pentru că se așteaptă scrierea în memoria sistem. La o operație STORE se așteaptă scrierea în memoria sistem, care este mai lentă.
- se mărește traficul între magistrală și memoria principală datorită faptului că fiecare acțiune STORE a procesorului se propagă în memoria principală.

Acest lucru nu se petrece la memoria *write-back*. În schimb, la acest tip de memorie este necesar un bit suplimentar pentru fiecare bloc din memoria cache, care să specifice dacă blocul este consistent cu memoria principală. Acest bit se numește bit *dirty*.

Dacă bitul *dirty* este 0, înseamnă că blocul din memoria cache este consistent (identic) cu blocul din memoria sistem. Dacă bitul *dirty* este 1, înseamnă că blocul nu este consistent (nu este actualizat în memoria sistem). Când un astfel de bloc *dirty* trebuie să fie eliminat din memoria cache, el va fi scris mai întâi în memoria principală.

Dezavantaj al memoriei *write-back*:

La schimbarea contextului program de către un sistem de operare multi-tasking se folosește întotdeauna un context diferit de memorie, de aici rezultă că un mare număr de blocuri din memoria cache vor trebui să fie înlocuite și astfel timpul de comutare a contextului program crește foarte mult.

S3) Strategii *write-miss*. Evenimentul *write-miss* apare când exist un *cache-miss* la o instruc iune STORE, adic blocul dorit nu se afl în memoria cache atunci când procesorul vrea s scrie în acesta.

Memoriile *write-back* folosesc 2 tipuri de strategii:

- a) Strategia “aloc la scriere” (*allocate on write*). Aici memoria cache aloc un bloc cache, dup care realizeaz ac iunea de a scrie în acel bloc.
- b) Strategia “încarc la scriere” (*fetch on write*). Aici memoria cache mai întâi cite te blocul din memoria principal , apoi scrie data în blocul citit.

Memoriile *write-through* folosesc strategia “nu se aloc la scriere” (*no allocate on write*). Aici, dac exist un *cache-miss* la scriere, atunci nu se mai aloc un bloc în memoria cache, ci data trimis de procesor va fi scris direct în memoria principal a sistemului.