

Memoria Cache

Andrei-Leonard PÂRVAN, Marius CARP

Transylvania University of Braşov, România,
Department of Electronics and Computers, Faculty of Electrical Engineering and Computer
Science,
Politehnicii nr. 1, Braşov, România, e-mail: andrei.parvan@student.unitbv.ro

Un calculator are trei sisteme logice: sistemul central de procesare, care procesează datele conform unor instrucţiuni specifice, sistemul de memorare şi stocare care stochează instrucţiunile şi datele (poate fi vorba despre memorie volatilă sau de un sistem de stocare permanentă, precum hard-disk-ul.), iar cel de-al treilea este sistemul de intrare-ieşire, responsabil pentru comunicaţia dintre calculator şi entităţi externe.

I. Introducere

Procesorul, în organizarea sa clasică, obţine performanţa maximă de viteză când în fiecare perioadă/tact de ceas, T_{clk} , este alimentat cu informaţie, care poate fi stocată în memorie. Accesarea acestei informaţii din memoria principală (SDRAM) necesită un timp de ordinul sutelor de tacte (numit latenţa memoriei), ceea ce ar înseamna că procesorul trebuie să se oprească până când primeşte informaţia necesară. Pentru a se elimina această oprire a funcţionării μP , se introduce între memoria principală şi μP un buffer/("rezervor") cu timp de acces redus, un nivel (L1) suplimentar de memorie, sau

două niveluri (L1+ L2), care formează subsistemul **memorie cache**.

Prin introducerea/prezenţa unei memorii cache se "ascunde" latenţa memoriei, dar nu se elimină. Subsistemul memorie cache este situat între μP (CPU) şi magistralele din sistemul principal. Cuvântul adresă generat de procesor se trimite simultan la memoria cache şi la memoria principală, iar dacă cuvântul căutat se află în cache, informaţia necesară se obţine de aici într-un timp scurt, iar adresarea în continuare la memoria principală nu se mai efectuează. Accesarea la memorie principală se va efectua numai când cuvântul căutat nu s-a aflat în cache, evident cu întârziere de sute de tacte.

Memoria cache este o memorie SRAM, care în raport cu memoriei principale (DRAM), este de capacitate redusă şi prezintă un timp de acces 1-2 T_{clk} . Raţiunea şi eficienţa memoriei cache în "ascunderea" latenţei memoriei principale se bazează pe:

- **localizarea (în timp şi spaţiu)** existentă în programele rulate pe calculator

- **incluziunea informaţiei stocate**, $L1 \subset L2 \subset M \subset HD$

II.Operația de citire a informației

Procesorul trimite la memoria cache adresa locației din memoria principală în care se află cuvântul căutat și dacă în cache există conținutul locației (cuvântul căutat) de la adresa solicitată din memorie, atunci acest cuvânt este transferat din cache în procesor în 1-2 tacte de ceas (*transfer pe cuvânt*). Dacă acest cuvânt nu este în cache, se accesează nivelul următor (memoria primară/principală) pentru care se consumă timp care poate fi chiar sute de tacte de ceas; dacă cuvântul căutat se află în memoria principală atunci blocul din memorie, în care se află cuvântul respectiv, este transferat și înscris în memoria cache, de unde, apoi se citește cuvântul de adresă căutat și se aduce în procesor. Dacă blocul căutat nu se află nici în memoria principală atunci se accesează nivelul următor (memoria secundară, hard disk), iar pagina de pe hard disk în care se află blocul căutat este transferată și înscrisă în memoria principală, iar din memoria principală din pagina în care se află blocul care conține cuvântul căutat este transferat și înscris în memoria cache, de unde cuvântul căutat (din blocul adus în cache) se trimite la procesor; tot acest proces de transfer consumă chiar milioane de tacte de ceas. O pagină poate cuprinde, sute sau chiar mii de blocuri.

III.Operația de înscriere a informației

Procesorul trimite la cache adresa din memoria principală unde cuvântul generat

de procesor trebuie înscris. Pentru înscriere există două metode: write-through și write-back.

Prin **write-through**, cuvântul data generat de procesor este înscris întâi la adresa din memoria cache și apoi se continuă cu înscrierea și în memoria principală; deoarece se face acces pentru înscriere și la memoria principală se consumă un timp similar cu cel consumat prin negăsirea cuvântului de la procesul de citire.

Prin metoda **write-back** se înscrie cuvântul generat de procesor numai în memoria cache, urmând ca înscrierea și în memoria principală să se realizeze doar atunci când blocul respectiv, din memoria cache în care s-a înscris cuvântul generat de procesor, este trimis înapoi la memoria principală (aceasta se realizează doar când respectivul bloc din cache este înlocuit de către un alt bloc adus din memoria principală); evident că prin write-back, până la înscrierea și în memoria principală a cuvântului generat de procesor, conținutul blocului din memoria cache și conținutul blocului de aceeași adresă din memoria principală nu este același.

Dacă la înscriere există miss în cache, adică blocul în care se află locația de adresă din memoria principală la care procesorul trebuie să înscrie nu a fost adus încă în memoria cache, atunci se caută blocul respectiv în memoria principală, se înscrie cuvântul generat de procesor în memoria principală, iar blocul respectiv se aduce în memoria cache, ceea ce este un proces consumator de multe sute de tacte de ceas.

IV.Tipuri de mapare

Există o multitudine de posibilități de definire a aplicației/mapării între

mulțimea numerelor blocurilor din memoria principală și mulțimea liniilor/intrărilor/blocurilor din memoria cache. Practic, sunt realizabile doar trei aplicații/mapări care determină următoarele modalități de organizare/structurare pentru memoriile cache: memorie cache cu mapare directă, memorie cache asociativă și memorie cache set-asociativă.

Memoria cache cu mapare directă

Regula de mapare, pentru o organizare de memorie cache cu mapare directă, permite conținutului unui bloc din memorie, **număr bloc din memorie**, să fie plasat doar într-o singură intrare/linie din numărul total de intrări ale memoriei cache, iar acea adresă/intrare **în memoria cache** a blocului respectiv se determină cu relația următoare:

adresa intrării în memoria cache = $(\text{Număr bloc din memorie}) \bmod (\text{număr intrări din cache})$

Memoria cache complet asociativă

Memoria complet asociativă (sau full-asociativă) se situează la extrema opusă în raport cu memoria cu mapare directă. Dacă la memoria cu mapare directă un bloc din memoria principală poate fi mapat/plasat doar într-o singură intrare a memoriei cache, la cea full-asociativă un bloc din memoria principală poate fi plasat în oricare intrare din memoria cache. Pentru că blocul de memorie poate fi plasat oriunde în cache, nu mai trebuie căutată intrarea/liniei din cache în care se plasează blocul.

Memoria cache set asociativă

Memoria set asociativă se situează intermediar în raport cu cea full asociativă și cea cu mapare directă. Memoria set asociativă prezintă un număr de seturi; fiecare dintre aceste seturi conține un număr n de căi/linii de intrare, iar aceste n căi ale unui set formează o memorie full-asociativă cu n intrări; o astfel de organizare este referită ca memorie set asociativă cu n -căi (n -way set-associative). Pentru o memorie cache cu un număr total de m intrări, fiecare dintre seturile sale având n căi, numărul total de seturi (segmente de mici memorii asociative) este egal cu raportul m/n . Fiecare set al memoriei cache este selectat de subcâmpul index din cuvântul de adresă, la fel ca și la memoria cache cu mapare directă, iar adresa setului (dintre cele m/n) din memoria cache set asociativă, în care va fi plasat un bloc din memoria principală, se determină conform relației următoare: adresa setului din cache = $(\text{număr bloc din memorie}) \bmod (\text{numărul seturilor din cache})$.

V.Referințe

- 1.Richard Stacpoole, Tariq Jamil - Cache Memories - <https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=839642>
- 2.Gheorghe Toacșe, Brașov-2020, Curs-Aritectura și organizarea microprocesoarelor