

4

Memorii și ierarhii de memorii

1. Prezentare teoretică

Memoria reprezintă una din componentele esențiale ale unui calculator numeric. Rolul său este de a memora programe și date. La modul general, memoria se împarte în memorie principală și memorie secundară. Memoria principală stochează informațiile și datele curente necesare procesorului cât timp memoria secundară memorează date care nu sunt folosite în mod curent.

Subsistemul de memorie al unui calculator numeric este văzut ca o ierarhie de module de memorie. La baza ierarhiei se află memoria auxiliară/secundară (lentă dar de capacitate foarte mare, de ordinul zecilor de GB), iar în vârful ierarhiei se află memoria **cache** (foarte rapidă dar de capacitate foarte mică, de ordinul KB).

Memoria cache este o memorie specială, utilizată pentru mărirea vitezei de prelucrare a procesorului. Aceasta este dispusă între procesor și memoria principală în vederea compensării diferenței de viteză dintre cele două componente. Viteza de funcționare a memoriei cache este foarte apropiată de viteza de funcționare a procesorului, acesta fiind un motiv pentru care, de obicei, această memorie este încapsulată împreună cu procesorul.

Memoria principală

Este o memorie cu un timp de acces de aproximativ 700ns. Capacitatea ei variază de la 32MB până la 512MB. Tehnologia folosită în construcția acestui tip de memorie se bazează pe circuite integrate semiconductoare. Circuitele integrate RAM pot opera în două moduri:

1. **static** - memoria statică RAM este formată din bistabile care memorează informația în formă binară. Informația rămâne disponibilă atât timp cât circuitul este alimentat cu curent electric.
2. **dinamic** - memoria dinamică RAM, memorează informația binară în formă de sarcini electrice ce sunt aplicate unor condensatori. Condensatorii, realizați cu ajutorul tranzistoarelor MOS, tind să se descarce în timp și astfel apare riscul ca informația să se distrugă. Pentru a elimina acest inconvenient, în aceste memorii se realizează un ciclu de reîmprospătare a informației la intervale de câteva milisecunde.

Memoria **ROM** (Read Only Memory) este o memorie cu acces aleator și este utilizată pentru memorarea programelor, care sunt rezidente permanent în calculator și pentru stocarea tabelor de constante definite de către producătorul calculatorului. Tot în această memorie se păstrează și programul inițial de pornire a calculatorului (bootstrap loader) responsabil cu lansarea sistemului de operare. Conținutul acestei memorii se păstrează chiar dacă circuitul nu este alimentat. Singura metodă de a șterge informația dintr-un astfel de circuit constă în iradierea sa cu raze ultraviolete.

Simbolul și tabela de funcționare a unui circuit de memorie RAM sunt prezentate în figura 4.1.

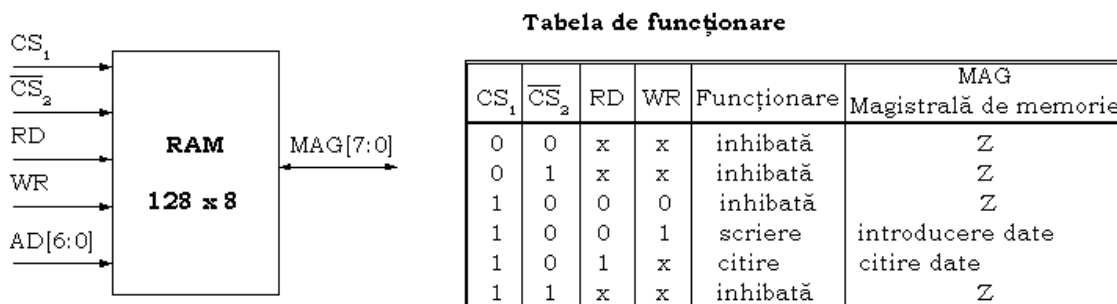


Figura 4.1: Memorie RAM 128 x 8.

Odată ce proiectantul unui calculator dispune de modulele de memorie fizice, el trebuie să asigneze memoriilor RAM sau ROM întreaga cantitate de memorie, determinată în prealabil, necesară pentru rularea aplicațiilor. Adresarea memoriei de către procesor este stabilită pe baza unei table în care se specifică adresele de memorie asigurate fiecărui dispozitiv/tip de memorie. Această tabelă poartă denumirea de "hartă adreselor de memorie". Modalitatea de construire a acestei table este prezentată prin intermediul unui exemplu.

Exemplu Se presupune că un calculator necesită 512 octeți de RAM și 128 octeți de ROM. Pentru implementare, la dispoziția proiectantului există circuitele RAM și ROM cu capacitate de 128 x 8 fiecare. Numărul circuitelor de memorie este limitat la 5. Harta adreselor de memorie este prezentată în tabelul 4.1.

Componenta	Adresa hexazecimală	Magistrala de adrese									
		10	9	8	7	6	5	4	3	2	1
RAM_1	0000 - 007F	0	0	0	x	x	x	x	x	x	x
RAM_2	0080 - 00FF	0	0	1	x	x	x	x	x	x	x
RAM_3	0100 - 017F	0	1	0	x	x	x	x	x	x	x
RAM_4	0180 - 01FF	0	1	1	x	x	x	x	x	x	x
ROM	0200 - 03FF	1	x	x	x	x	x	x	x	x	x

Tabelul 4.1: Harta adreselor de memorie.

Memoria totală este: $512 + 128 = 640$ octeți. Pentru a accesa o memorie de 640 octeți este necesar un cuvânt-adresă de 10 biți. Astfel, dacă se va considera o magistrală de adrese de 16 biți, primii 6 biți vor fi întotdeauna egali cu 0. Având în vedere că memoria RAM ocupă 512 octeți rezultă necesitatea utilizării a 9 linii de adresă. Deoarece sistemul conține și memorie RAM și memorie ROM, o linie de adresă este necesară pentru a se realiza o distincție între ele. Deci în total vor fi 10 linii de adresă ocupate. Restul de 6 linii de adresă vor fi tot timpul 0. Pentru ușurința proiectării, valoarea de pe magistrala de adrese este exprimată în hexazecimal.

Memoria asociativă

Memoria asociativă se bazează pe conceptul că timpul de căutare a unei informații date într-o memorie este redus considerabil dacă datele memorate pot fi identificate după conținut și nu după adresă.

Când se scrie un cuvânt într-o memorie asociativă, nici o adresă nu este furnizată. Memoria este capabilă să determine locația goală și neutilizată în vederea memorării noului cuvânt. La citirea unui cuvânt dintr-o memorie asociativă, conținutul cuvântului sau o parte a cuvântului este specificată. Memoria localizează toate cuvintele care se potrivesc cu conținutul specificat și ele vor fi marcate pentru citire.

Organizarea hardware a unei memorii asociative este prezentată în figura 4.2.

În figura 4.2 registrele A și K sunt registre de n biți. Registrul M , registrul de potrivire, este de m biți, unul pentru fiecare cuvânt de memorie.

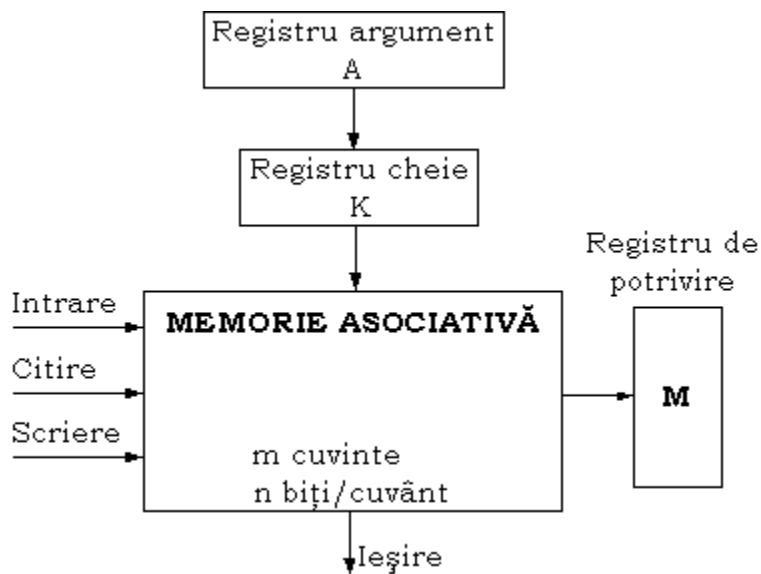


Figura 4.2: Organizarea hardware a unei memorii asociative

Fiecare cuvânt din memorie este comparat cu conținutul registrului A . Cuvintele care se potrivesc setează un bit corespunzător în registrul M a cărui valoare este 1. În final, registrul M va conține doar biții care indică ce cuvinte s-au potrivit. Citirea se realizează printr-un acces secvențial la memorie, doar pentru cuvintele ai căror biți corespunzători în registrul M sunt 1.

Registrul K este un registru de mascare în vederea stabilirii unui câmp particular sau a unei chei din cuvântul memorat în registrul A . Compararea întregului cuvânt de memorie cu conținutul registrului A se va realiza doar dacă registrul K conține toți biții egali cu 1. În mod normal se compară doar cuvintele care au biții 1 în pozițiile în care există 1 în registrul K .

Organizarea hardware a unei celule de memorie asociative precum și relația dintre vectorul memorie și registrele externe sunt prezentate în figura 4.3.

Așa cum se observă în figura 4.3, fiecare bit A_j din registrul A este comparat cu toți biții coloanei j din matricea vector, dacă $k_j = 1, \forall j = \overline{1 \dots n}$. Dacă toți biții registrului K sunt egali cu biții din cuvântul i , atunci se setează $M_i = 1$, altfel $M_i = 0$.

Circuitul de potrivire logică din figura 4.3 se poate deduce matematic. Cuvântul i este egal cu conținutul registrului A dacă:

$$A_j = F_{ij} \quad \forall j = \overline{1 \dots n} \quad (4.1)$$

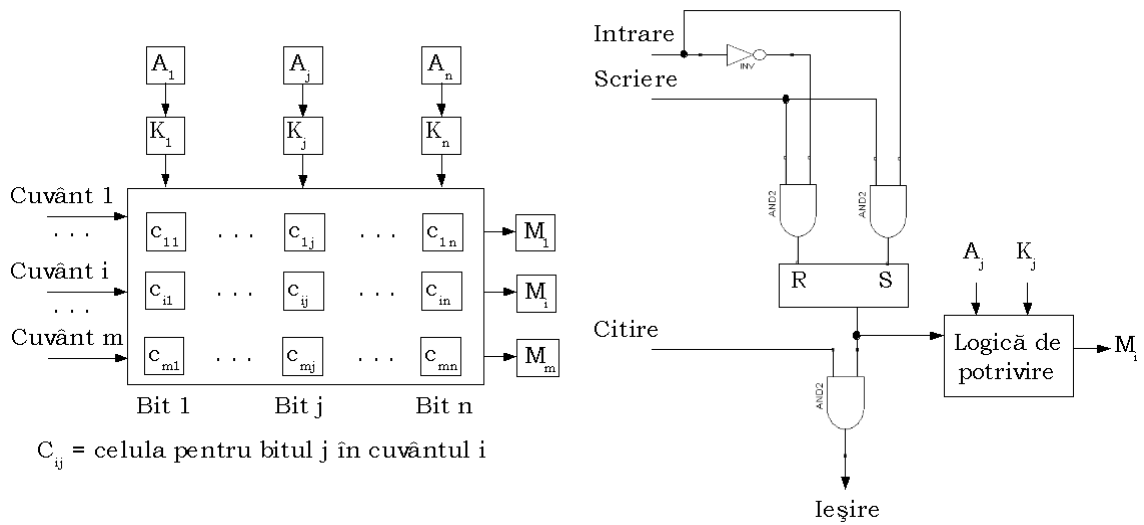


Figura 4.3: Organizarea hardware a unei celule de memorie

Egalitatea a doi biți poate fi exprimată ca $x_j = A_j \cdot F_{ij} + \overline{A_j} \cdot \overline{F_{ij}}$, unde $x_j = 1$ dacă toți biții de pe poziția j sunt egali. Pentru ca bitul M_i să fie 1 trebuie respectată următoarea egalitate:

$$M_i = x_1 \cdot x_2 \cdot x_3 \cdot \dots \cdot x_n \quad (4.2)$$

În funcție de valoarea lui K_j , există două situații posibile:

1. $k_j = 0$, atunci biții corespunzători A_j și F_{ij} nu trebuie comparați;
2. $k_j = 1$, atunci biții corespunzători A_j și F_{ij} trebuie comparați.

Se poate concluziona că:

$$x_j + \overline{k_j} = \begin{cases} x_j & \text{daca } k_j = 1 \\ 1 & \text{daca } k_j = 0 \end{cases} \quad (4.3)$$

Ecuția (4.2) se rescrie ca:

$$M_i = (x_1 + \overline{k_1}) \cdot (x_2 + \overline{k_2}) \cdot (x_3 + \overline{k_3}) \cdot \dots \cdot (x_n + \overline{k_n}) \quad (4.4)$$

sau într-o formă mult mai compactă:

$$M_i = \prod_{j=1}^n (A_j \cdot F_{ij} + \overline{A_j} \cdot \overline{F_{ij}} + \overline{k_j}) \quad \forall i = \overline{1, m} \quad (4.5)$$

În vederea realizării operației de citire din memoria asociativă, se scanează conținutul registrului M și se citește doar un bit al său la un moment dat, obținând astfel o secvență de cuvinte aflate în memoria asociativă și care se potrivesc cu cuvântul de memorie.

În cazul operației de scriere apar următoarele situații posibile:

1. *memoria este vidă* - în acest caz scrierea poate fi realizată prin adresarea fiecărei locații într-o anumită secvență. Astfel memoria devine o memorie cu acces aleator în cazul scrierii și o memorie de tip adresabilă după conținut în cazul operației de citire. Avantajul este că adresele cuvintelor de intrare pot fi decodificate ca în cazul memoriilor cu acces aleator, rezultând astfel numai d linii de adrese în loc de m ($m = 2^d$).
2. *suprascrierea unui cuvânt* - este situația în care memoria este complet ocupată și se dorește introducerea unui nou cuvânt în memorie. În această situație se mai adaugă un registru special (registru de etichete) care memorează cuvintele active și cele inactive asignând valoarea 1 pentru cuvintele active și 0 pentru cuvintele inactive. Numărul biților acestui registru trebuie să fie același cu numărul de cuvinte din memorie. Se înlocuiesc doar cuvintele cu eticheta 0 utilizând în acest scop un algoritm FIFO, LRU etc. După realizarea operației de suprascriere, eticheta corespunzătoare va avea valoarea 1.

Memoria cache

Memoria cache este o memorie de mică capacitate dar cu o viteză de accesare foarte mare. Când CPU dorește să acceseze memoria principală, întâi se examinează memoria cache. Dacă, cuvântul solicitat se găsește stocat în memoria cache, atunci el este preluat de către procesor, din această memorie. În caz contrar este accesată memoria principală, se citește cuvântul dorit și este stocat apoi în memoria cache.

Pentru a respecta principiul "localității spațiale", se transferă din memoria principală în memoria cache nu doar cuvântul solicitat de procesor, ci un bloc de cuvinte care conține cuvântul solicitat. Dimensiunea blocului transferat variază de la o arhitectură la alta, dar uzual, blocul transferat are o dimensiune cuprinsă între 1 și 16 cuvinte.

Performanța unei memorii cache este măsurată în termeni de "hit ratio". Dacă, cuvântul solicitat de CPU se află memorat în memoria cache, atunci este vorba de un succes **HIT**, iar altfel, de un insucces **MISS**. **HIT RATIO** se definește ca fiind raportul dintre numărul de succese și numărul total de referiri CPU la memorie.

Procesul de amplasare a datelor citite din memoria principală, în memoria cache se numește *mapare*. Uzual în proiectarea și organizarea memoriilor cache se folosesc

următoarele tipuri de proceduri de mapare:

1. mapare complet asociativă;
2. mapare directă;
3. mapare asociativă pe mai multe căi.

Pentru o descriere cât mai completă a fiecărui mod de mapare se va utiliza următorul exemplu.

Exemplu Se presupunem existența unui calculator care posedă o memorie principală cu capacitatea de $32k \times 12$ și o memorie cache cu capacitatea $512 \text{ cuvinte} \times 12$. Procesorul calculatorului poate comunica cu ambele memorii. El trimite o adresă de 15 biți către memoria cache. Dacă se obține HIT, CPU acceptă 12 biți de date de la memoria cache. În caz de MISS, CPU citește cuvântul din memoria principală și cuvântul va fi apoi memorat în memoria cache.

Maparea complet asociativă

Acest tip de mapare este cea mai flexibilă și mai rapidă organizare a memoriei cache. Organizarea acestui tip de memorie poate fi observată în figura 4.4. În memoria complet asociativă sunt memorate adresele și conținutul cuvintelor de memorie.

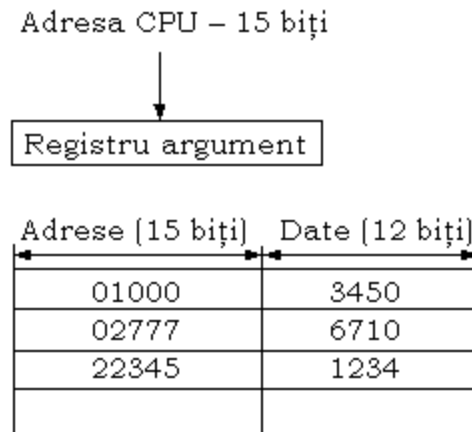


Figura 4.4: Maparea asociativă.

Adresa CPU de 15 biți este memorată în registrul argument/descriptor și apoi se încearcă depistarea unei potriviri de adrese dintre această adresă și o adresă din memoria cache. Dacă procesul se încheie cu HIT, atunci datele de 12 biți sunt citite și transmise către CPU.

Dacă procesul se încheie cu MISS, atunci perechea adresă-dată este citită din memoria principală, trimisă la procesor și memorată în memoria cache. În cazul cel mai

defavorabil, memoria este complet ocupată și noua pereche adresă-dată trebuie suprascrisă în locul altei perechi stabilită de către algoritmul de replasare utilizat (uzual se folosesc algoritmi LRU și FIFO).

Mapare directă

O versiune mai ieftină de memorie cache este prezentată în figura 4.5. În cadrul acestei implementări, adresa CPU se împarte în două câmpuri: INDEX (numărul biților acestui câmp este egal cu numărul biților necesari pentru a accesa memoria cache) și ETICHETĂ. În concluzie, cei 15 biți ai adresei CPU se împart după cum urmează: 6 biți pentru câmpul ETICHETĂ și 9 biți pentru câmpul INDEX.

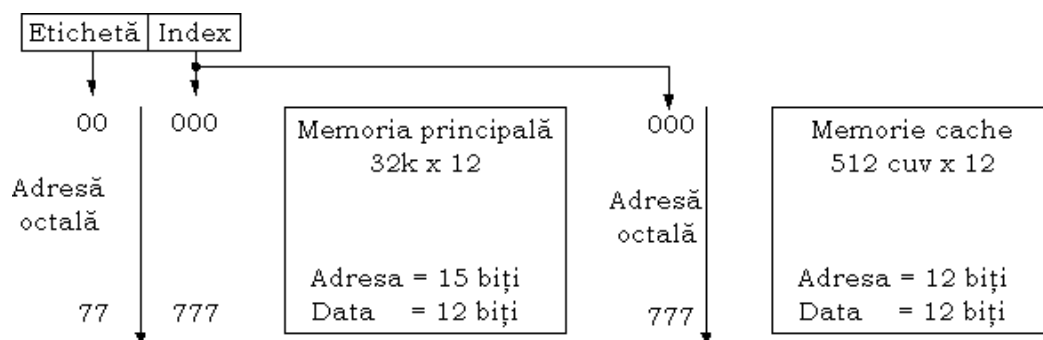


Figura 4.5: Maparea directă.

Fiecare cuvânt memorat în memoria cache este format din câmpul de date și eticheta asociată. Atunci când o cerere este lansată, câmpul INDEX este utilizat pentru adresa de acces a memoriei cache. Câmpul ETICHETĂ al adresei CPU este comparat cu "eticheta" cuvântului citit din memoria cache.

Dacă procesul returnează MISS, cuvântul cerut se citește din memoria principală și este memorat în memoria cache împreună cu o nouă etichetă, înlocuindu-se astfel valoarea anterioară.

Dezavantajul major constă în scăderea valorii HIT RATIO dacă două sau mai multe cuvinte ale căror adrese au același index dar etichete diferite sunt accesate în mod repetat.

Un exemplu numeric este prezența în figura 4.6.

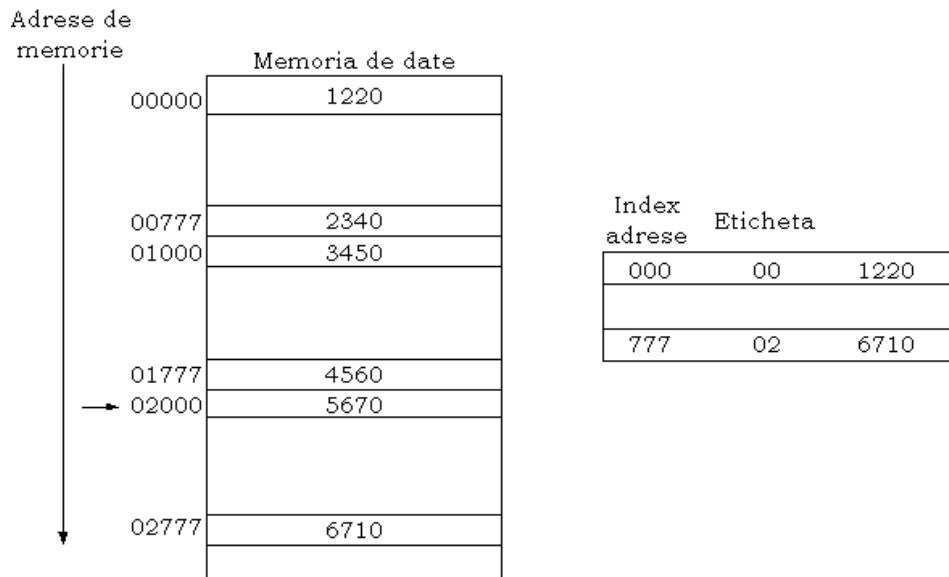


Figura 4.6: Exemplu mapare directă.

Cuvântul aflat la adresa 00000 este memorat în memoria cache. Procesorul dorește să acceseze cuvântul aflat la adresa 02000. Indexul de adrese este 000 și el va fi utilizat pentru accesarea memoriei cache. În urma comparației dintre cele două etichete, se constată că rezultatul întors este MISS, deoarece eticheta memoriei cache este 00, iar eticheta adresă este 02.

Se accesează memoria principală și cuvântul 5670 este transferat către CPU. Cuvântul din memoria cache de la 000 este replasat cu eticheta 02 și data 5670.

Maparea asociativă pe mai multe căi

În cadrul acestei metode de organizare a memoriei cache, fiecare cuvânt al memoriei cache poate memora unul sau mai multe cuvinte din memorie sub același index de adresă. În figura 4.7 este prezentat un exemplu de memorie cache cu o organizare asociativă pe mai multe căi.

Index adrese	Eticheta	Date	Eticheta	Date
000	01	3450	02	5670
777	02	6710	00	2340

Figura 4.7: Maparea asociativă pe mai multe căi.

Fiecare index de adresă referă două cuvinte de date și etichetele asociate lor. Fiecare etichetă necesită 6 biți, iar fiecare cuvânt 12 biți, rezultând lungimea cuvântului de

$2 \cdot (6 + 12) = 36$ biți. Un index de adrese de 9 biți poate acoperi 512 cuvinte deci dimensiunea memoriei cache este de 512×36 . Această memorie cache poate acoperi un număr de 1024 cuvinte ale memoriei principale deoarece un cuvânt din memoria cache conține două cuvinte dată.

Când CPU generează o cerere de memorie, valoarea index a adresei este utilizată pentru accesarea memoriei cache. Câmpul etichetă din adresa CPU este comparat cu cele două etichete din memoria cache în vederea stabilirii unei potriviri. Compararea logică este realizată printr-o căutare asociativă de etichete într-o mulțime. În caz de MISS, algoritmi de înlocuire cei mai frecvent utilizați sunt FIFO și LRU.

Scrierea în memoria cache

În cazul în care CPU generează semnal de scriere în memoria principală, există două metode care se pot aplica:

1. reactualizarea informației din memoria principală în paralel cu reactualizarea informației din memoria cache dacă cuvântul se găsește în memoria cache, la adresa specificată. Această procedură se numește WRITE-THROUGH.
2. WRITE-BACK. În cadrul acestei metode, doar locația memoriei cache este reactualizată în cadrul operației WRITE. Locația reactualizată este marcată și doar în caz de suprascriere a informației se va face reactualizarea locației memoriei principale.

2. Desfășurarea lucrării

1. Se va utiliza un simulator de memorie cache pentru a simula diferite organizări ale memoriei cache pentru primele 1 milion de referințe în monitorizarea execuției programului *gcc*. Sunt disponibile (<http://www.csit-sun.pub.ro/resources>) atât *dinero* (simulatorul de memorie cache), cât și *gcc*. Se presupune o memorie cache de instrucțiuni de 32kB și o memorie cache de date de 32kB folosind aceeași organizare. Să se aleagă cel puțin două tipuri de asociativități și două dimensiuni de bloc. Să se deseneze o diagramă ce prezintă organizarea unei memorii cache cu cea mai bună rată de eșec.
2. Se va studia influența folosirii unui nivel secundar de memorie cache asupra performanței unui procesor. Se presupune existența unui procesor cu un CPI (numărul de cicluri/instrucțiune) de bază de 1.0 și o frecvență de ceas de 500MHz. Se consideră că toate referințele la memorie vor avea succes în memoria cache primară. Memoria principală are un timp de acces de 200ns, incluzând aici și timpul necesar tratării cazurilor de eșec. Frecvența de eșec pe instrucțiune în memoria cache primară este de 5%. Care va fi creșterea de viteză

a mașinii dacă este adăugat un nivel suplimentar de memorie cache cu un timp de acces de 20ns, atât pentru eșec, cât și pentru succes și suficient de mare pentru a reduce frecvența de eșec la memoria principală la 2%?

3. Probleme propuse

1. Pe baza hărții adreselor de memorie stabilită în cadrul exemplului 1, să se proiecteze circuitul care realizează conexiunea memoriei cu CPU.
2. Pe baza organizării hardware a unei celule de memorie asociativă prezentată în figura 4.3, să se proiecteze o memorie asociativă cu $m = 4$ și $n = 3$. Să se realizeze simularea operațiilor de scriere/citire în/din memoria asociativă proiectată.
3. Se reconsideră exemplul prezentat în figura 4.6. În cadrul exemplului mărimea blocului de memorie are dimensiunea de 1 cuvânt. Cum arată aceeași proiectare dacă se utilizează un bloc de capacitate 8 cuvinte ?
4. Enunțați cel puțin un avantaj al folosirii metodei WRITE-THROUGH de scriere în memoria cache.
5. Folosind limbajul de asamblare pentru MIPS și simulatorul SPIM exemplificați toate modurile de adresare la memorie cunoscute. Spre exemplu, pentru a exemplifica modul de adresare indirectă cu autodecrementare se poate scrie programul MIPS din figura 4.8.

```
# adresare indirectă cu autodecrementare
.text
la      $a1, 0x0040001c
sub     $a1, $a1, 0x00000004
lw      $t1, ($a1)
lw      $t2, ($t1)
pc:    not      $t1, $t2
adr:    .word   0x00400020
li      $v0, 1
op:     .word   7
        move   $a0, $t1
        syscall
end:    mop
```

Figura 4.8: Program MIPS pentru exemplificarea modului de adresare indirectă cu autodecrementare a memoriei principale.