

# ARHITECTURA CALCULATOARELOR 2003/2004

## CURSUL 8

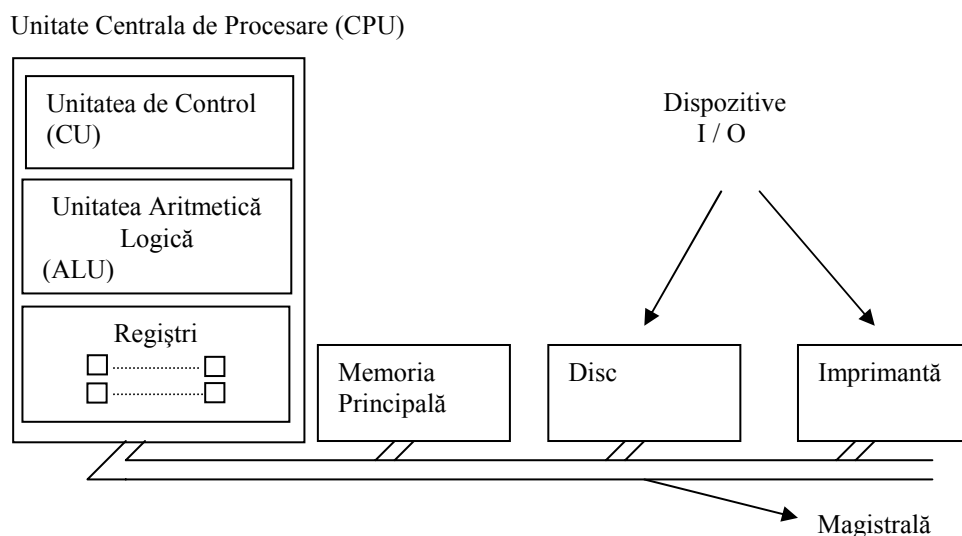
### Capitolul 3: Organizarea sistemelor de calcul

Un calculator numeric este un sistem interconectat de **procesoare, memorii, și dispozitive intrare/ieșire**. Aceste trei componente sunt conceptele cheie ale arhitecturii calculatorului.

#### 3.1 Procesoare

##### 3.1.1 Organizarea CPU

Organizarea unei structuri de calcul simple, bazată pe o magistrală, este prezentată în figura 3.1. **CPU - ul (Unitatea Centrală de Procesare)** este “creierul” calculatorului. Funcțiunea sa este de a executa programele stocate în memoria principală prin extragerea instrucțiunilor, decodificarea lor și apoi executarea lor una după alta. Componentele sunt conectate printr-o magistrală (bus), care este, din punct de vedere fizic, o colecție de trasee electrice pentru transmiterea adreselor, datelor și a semnalelor de control. Magistrala poate fi externă CPU-ului, conectând dispozitivele de I/O și memoria, sau poate fi internă CPU-ului.



**Figura 3.1** Organizarea unui calculator simplu cu CPU și două dispozitive I/O

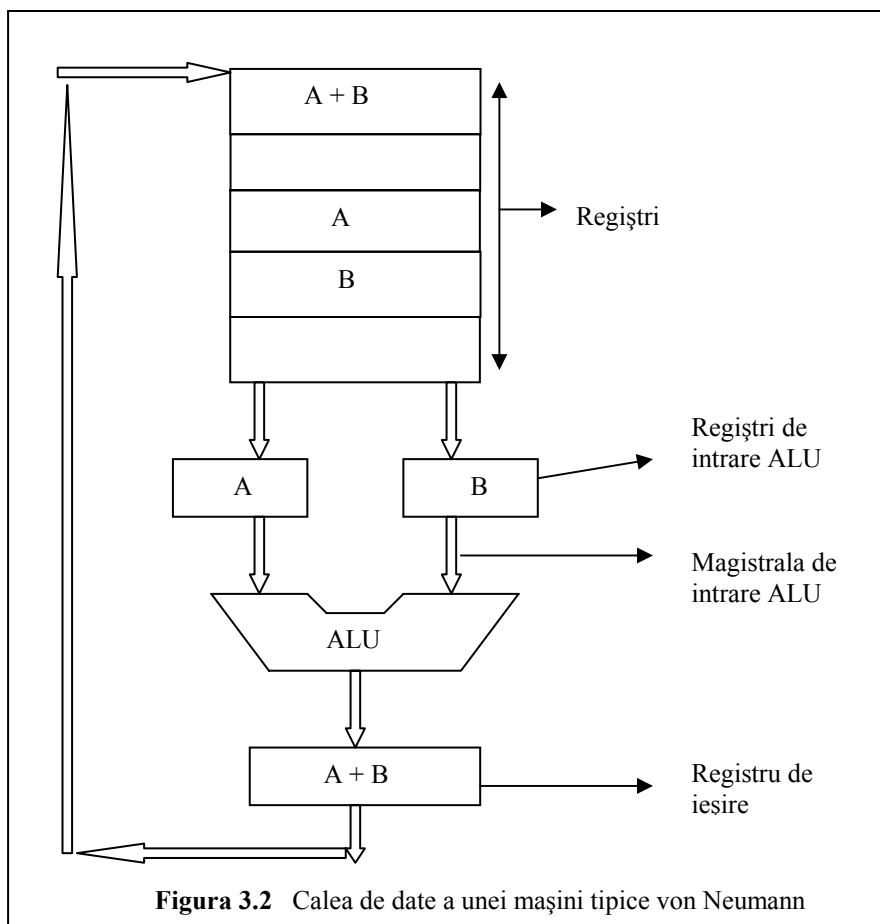
CPU este compus din mai multe părți distincte. **CU (Unitatea de control)** este responsabilă pentru extragerea instrucțiunilor din memoria principală și determinarea tipului lor. **ALU (Unitatea Aritmetică – Logică)** execută operații aritmetice și logice, cum ar fi adunarea sau AND logic, care sunt necesare pentru realizarea instrucțiunilor.

CPU conține deasemenea o memorie mică de înaltă viteză folosită pentru a stoca rezultate temporare și anumite informații de control. Memoria este alcătuită dintr-un număr de regiștri,

din care fiecare are o anumită mărime și funcțiune. De obicei, toți regiștrii au aceeași mărime (lungime). Fiecare registru poate conține un număr, până la un maximum determinat de mărimea registrului. Regiștrii pot fi citiți și scriși cu o mare viteză pentru că sunt interni CPU-ului.

Regiștrul cel mai important este **PC (Numărătorul de program)**, care indică adresa următoarei instrucțiuni care urmează să fie executată. Numele “Numărător de program” este într-un fel greșit deoarece acest registru nu are nimic de a face cu numărarea, dar termenul este foarte folosit. De asemenea important este **IR (Regiștrul de instrucțiune)**, care memorează codul instrucțiunii curente care se execută. Majoritatea calculatoarelor au numeroase alte registre, unele de uz general ca și altele pentru scopuri specifice.

Un exemplu de organizare tipică von Neumann pentru CPU este prezentată în Figura 3.2. Această structură este numită **cale de date (data path)** și este alcătuită din regiștri (tipic de la 1 la 32), ALU și mai multe magistrale conectând componentele. Regiștrii sunt conectați la doi regiștri de intrare ai ALU, etichetați A și B în figura 3.1. Acești regiștri memorează intrarea ALU în timp ce ALU calculează. Calea de date este foarte importantă în toate sistemele de calcul și vom reveni asupra acestui subiect.



Chiar ALU realizează adunarea, scăderea, și alte operații asupra intrarilor sale și oferă rezultatul în registrul de ieșire. Acest registru de ieșire poate fi stocat înapoi într-un registru. Mai tarziu, registru poate fi scris (adică stocat) în memorie, dacă se dorește. Nu toate

structurile de calcul au A, B, și registru de ieșire. În exemplu este presupusă o operație de adunare.

Majoritatea instrucțiunilor pot fi împărțite în două categorii: registru - memorie sau registru - registru. Unele instrucțiuni registru - memorie permit aducerea cuvintelor de memorie în regiștri, unde ele pot fi folosite ca intrări ALU pentru instrucțiunile corespunzătoare. Cuvintele (word) sunt unități de date care se mută între memorie și regiștrii. Un cuvânt poate fi și un întreg. Alte instrucțiuni registru - memorie permit ca registrele să fie stocate înapoi în memorie.

Celalalt tip de instrucțiune este registru - registru. O instrucțiune registru - registru tipică aduce cei doi operanzi din registre, îi depune în regiștrii de intrare ALU, realizează unele operații cu ele (de exemplu, adunarea sau AND logic) și stochează rezultatul într-unul dintre regiștri. Procesul trecerii a doi operanzii prin ALU și stocarea rezultatului este numit **ciclu căii de date** și este partea cea mai importantă a majorității CPU - urilor. Prin extensie, el definește ceea ce poate face mașina. Cu cât ciclul căii de date este mai rapid cu atât mașina respectivă este mai rapidă.

### 3.1.2 Executarea instrucțiunilor

CPU execută fiecare instrucțiune într-o secvență de pași. În linii mari pașii sunt următorii:

1. Extrage instrucțiunea următoare din memorie în registrul de intrare,
2. Schimbă numărătorul de programe pentru a indica instrucțiunea următoare,
3. Determină tipul instrucțiunii care tocmai a fost încărcată,
4. Dacă instrucțiunea folosește un cuvânt din memorie determină unde este el,
5. Dacă este necesar, aduce cuvântul într-un registru CPU,
6. Execută instrucțiunea,
7. Reia de la pasul 1 pentru a începe executarea următoarei instrucțiuni.

Această succesiune de pași este frecvent referită drept ciclul **extragere (fetch) - decodificare - execuție**. Ea este activitatea centrală pentru CPU în toate calculatoarele.

Această descriere a modului cum lucrează CPU - urile seamănă foarte mult cu un program scris în limbaj natural. În figura 3.3 se prezintă un astfel de program rescris într-o manieră procedurală (Java) și numit interpretor.

Mașina care este interpretată are doi regiștri vizibili din punctul de vedere al utilizatorului:

- numărător de programe (PC), pentru memorarea suitei de adrese ale următoarei instrucțiuni care trebuie extrase și
- acumulatorul (AC), pentru memorarea rezultateor aritmetice.

Ea are deasemenea regiștri interni pentru memorarea:

- instrucțiunii curente (instr),
- tipului instrucțiunii curente (instr\_type),
- adresei operandului instrucțiunii (data\_loc),
- a operandului curent însuși (data).

Se presupune că instrucțiunile conțin o singură adresă de memorie. Locația de memorie adresată conține operandul, de exemplu valoarea care să fie adăugată la acumulator.

```

public class Interp {
    static int PC;           // registru pentru adresa instrucțiunii următoare
    static int AC;           // acumulatorul, un registru pentru operații aritmetice
    static int instr;        // registru pentru instrucțiunea curentă
    static int instr_type;   // tipul de instrucțiune (opcode)
    static int data_loc;     // adresa operandului, sau -1 dacă nu este niciunul
    static int data;         // registru pentru operandul curent
    static boolean run_bit = true; // un comutator care poate fi trecut pe off pentru a
                                // opri calculatorul

    public static void interpret(int memory[ ], int starting_address ) {
// Acest program interpretează programe pentru un calculator simplu cu instrucțiuni care au un
//singur operand din memorie. Această mașină are un registru AC (acumulator), folosit pentru
//operații aritmetice. Spre exemplu, instrucțiunea ADD adaugă un întreg din memoria la AC.
//Interpreterul continuă rularea până ce o instrucțiune HALT aduce run_bit la valoarea 0.
//Starea rulării proceselor pe această mașină constă din memorie, numărătorul de program,
//run_bit și AC. Parametrii de intrare sunt imaginea memoriei și adresa de început.
        PC = starting_address;
        While (run_bit)      {
            instr = memory[PC];           // aduce instrucțiunea următoare în instr
            PC = PC + 1;                   // incrementează numărătorul de program
            instr_type = get_instr_type(instr); // determină tipul de instrucțiune
            data_loc=find_data(instr, instr_type); // localizează data, -1 dacă nu e niciuna
            if ( data_loc >= 0)           // dacă -1 atunci nu este nici un operand
                data = memory [data_loc]; // încarcă operand
            execute(instr_type, data);    // execută instrucțiunea
        }
    }

    private static int get_instr_type(int addr)      { ... }
    private static int find_data(int instr, int type) { ... }
    private static void execute(int type, int data)  { ... }
}

```

**Figura 3.3** Un interpretor pentru un calculator simplu (scris în Java).

Faptul că este posibil să se scrie un program care să simuleze funcționarea CPU arată că nu este necesar ca un program să fie executat de către un CPU hardware constând dintr-o cutie plină de electronice. În schimb, un program poate fi executat de alt program care îi extrage, decodifică și execută instrucțiunile. Un program (ca în figura 3.3) care extrage, decodifică și execută instrucțiunile altui program este numit **interpretor**.

Această echivalență între procesoarele hardware și interpretoarele software are implicații importante pentru organizarea și proiectarea sistemelor de calcul. După ce a specificat limbajul mașinii, L, pentru un calculator nou, echipa de proiectare decide dacă vor să construiască un procesor hardware pentru a executa direct programe în L sau, alternativ, dacă vor să scrie un interpretor pentru a interpreta programe pentru L. Dacă ei aleg să scrie un interpretator, ei trebuie deasemenea să ofere o altă mașină hardware care să ruleze interpretatorul. Anumite construcții hibrid sunt deasemenea posibile cu unele execuții hardware, dar și cu interpretare software.

Un interpretor sparge instrucțiunile unei mașini țintă în pași mici. Ca o consecință, mașina pe care rulează interpretorul poate fi mult mai simplă și mai puțin scumpă decât un procesor hardware pentru mașina țintă. Această economie este semnificativă dacă mașina țintă are un număr mare de instrucțiuni și instrucțiunile sunt complicate, cu multe opțiuni. Economia provine în esență din faptul că un sistem hardware este înlocuit de un sistem software (interpretator).

Primele computere aveau seturi simple de instrucțiuni. Dar goana pentru computere mai puternice conduce, printre alte lucruri, la instrucțiuni individuale mai puternice. Foarte repede s-a descoperit că instrucțiunile mai complexe conduc la execuții mai rapide ale programelor, chiar dacă instrucțiunile individuale durează mai mult. O instrucțiune în virgulă mobilă este un exemplu de instrucțiune complexă. Instrucțiunile pentru accesarea blocurilor de elemente este alt exemplu. Tot așa de simplă este observația că pentru două instrucțiuni care apar frecvent consecutiv este preferabilă o singură instrucțiune care poate să realizeze mai rapid efectul celor două.

Instrucțiunile mai complexe sunt preferabile deoarece executarea operațiilor individuale poate fi executată în paralel folosind resurse hardware diferite. Pentru computerele scumpe de înaltă performanță costul acesta suplimentar ar putea fi justificat. De aceea computerele scumpe de înaltă performanță tind să aibă mult mai multe instrucțiuni decât calculatoarele cu un preț mai scăzut. Totuși creșterea costului dezvoltării de software și cerințele de compatibilitate a instrucțiunilor au creat necesitatea de a implementa instrucțiuni complexe chiar și pe calculatoarele de preț scăzut, unde prețul este mai important decât viteza.

O implementare hardware directă, adică neinterpretată, a fost folosită numai în implementarea modelelor scumpe. Calculatoarele simple cu instrucțiuni interpretate au avut deasemenea alte beneficii. Printre cele mai importante importante au fost:

- Abilitatea de a repara instrucțiunile implementate incorect sau a unor deficiențe de proiectare în hardware – ul de baza,
- Oportunitatea de a adăuga unele noi instrucțiuni, cu cost minimal, chiar după livrarea mașinii,
- Proiectarea structurată care a permis dezvoltarea eficientă, testarea și documentarea instrucțiunilor complexe.

Prin 1950 IBM a sesizat că dezvoltarea unei singure familii de mașini, toate executând aceleași instrucțiuni, are multe avantaje atât pentru IBM cât și pentru clienții ei. IBM a introdus termenul de **arhitectură** pentru a descrie acest nivel de compatibilitate. O nouă familie de calculatoare ar avea o aceeași arhitectură, dar multe implementări diferite care să poată toate să execute același program, diferind numai prețul și viteza. Dar cum să construim un computer de preț scăzut care ar putea executa toate instrucțiunile complicate ale unei mașini scumpe de înaltă performanță? Raspunsul sta în **interpretare**. Această tehnică, care a fost prima dată sugerată de Wilkes (1951), a permis proiectarea calculatoarelor simple, de preț mai scăzut, care totuși au putut să execute un mare număr de instrucțiuni. Rezultatul a fost arhitectura sistemului IBM/ 360, o familie de calculatoare compatibile, traversând aproape două ordine de magnitudine, atât în preț cât și în capacitate.

Pentru că piața de calculatoare a explodat în anii 1970 și capacitățile de calcul au crescut rapid, cererea de calculatoare cu preț scăzut a favorizat proiectarea calculatoarelor folosind interpretoare. Abilitatea de a imbrina hardware și interpretoare pentru un set de instrucțiuni

particulare s-a dovedit un factor important pentru proiectarea eficientă. Pentru că tehnologia de semiconductori a avansat rapid avantajele costului au depășit ponderea acordată performanțelor mai înalte și **astfel arhitecturile bazate pe interpretor au devenit calea convențională de a proiecta calculatoarele**. Aproape toate noile calculatoare proiectate începând din 1970, de la microcalculatoare la calculatoare mainframes, s-au bazat pe o interpretoare.

Prin anii '70 folosirea procesoarelor simple care folosesc interpretoare a devenit foarte larg răspândită exceptând cele mai scumpe modele de înaltă performanță, ca și Cray - 1 și seriile Control Data Cyber. Folosirea unui interpretor a eliminat inevitabilele limitări de cost ale instrucțiunilor complexe și arhitecturile au început să includă mult mai mult instrucțiunile complexe, în particular cu căi multiple de a specifica operanzii. Curentul a atins zenitul cu calculatorul VAX de la DEC care au avut mai multe sute de instrucțiuni și mai multe de 200 de căi diferite de specificare a operanzilor care să fie folosiți în fiecare instrucțiune. Din nefericire arhitectura VAX a fost concepută de la început să fie implementată cu un interpretor, fără a lua în considerare posibilitatea de a implementa un model de performanță. Ca urmare s-au inclus un număr foarte mare de instrucțiuni de valoare marginală care erau dificil de executat direct. Această omisiune s-a dovedit a fi fatală lui VAX. DEC a fost cumparat de Compaq în 1998.

Deși la început procesoarele cu 8 biți erau mașini foarte simple cu seturi de instrucțiuni foarte simple, prin anii 1970, chiar microprocesoarele au fost transformate în proiecte bazate pe interpretoare. În timpul acestei perioade una din provocările majore a fost creșterea complexității circuitelor integrate. Un avantaj major al abordării bazate pe interpretor a fost abilitatea de a proiecta un procesor simplu, care transfera complexitatea la memoria care găzduiește interpretorul. De aceea o proiectare complexă hardware a putut fi transformată într-o proiectare complexă software. Succesul Motorola 68000 care a avut un mare set de instrucțiuni interpretate și eșecul concurentului Zilog Z8000, care a avut un set de instrucțiuni egal de mare, dar fără un interpretator, a demonstrat avantajele unui interpretor prin aducerea rapidă pe piață a noului microprocesor. Succesul a fost tot mai mare surprinzând conducerea lui Zilog, mai ales că predecesorul lui Z8000, Z80, a fost mai popular decât predecesorul lui 68000, 6800). Bineînțeles alți factori au fost deasemenea implicați aici, nu în ultimul rând a fost lunga istorie a firmei Motorola ca fabricant de cipuri, în timp ce Exxon (detinatorul lui Zilog) este o companie de petrol nu un fabricant cipuri.

**Alt actor lucrând în favoarea interpretării a fost existența memoriilor ROM rapide (control stores) capabile să stocheze interpretorul.** Să presupunem că instrucțiunea tipică 68000 interpretată corespunde la 10 instrucțiuni interpretor, numite **microinstrucțiuni**, a câte 100 nanosecunde fiecare, și două accesuri la memoria principală, a câte 500 nanosecunde fiecare. Timpul total de execuție a fost atunci 2000 nanosecunde, de numai două ori mai rău decât ar fi putut să ofere cea mai bună execuție directă. Dacă memoria rapidă nu ar fi fost disponibilă, atunci instrucțiunea ar fi luat 6000 de nanosecunde.

### 3.2 Memoria cache

De-a lungul timpului, **microprocesoarele au fost întotdeauna mai rapide decât memoriile**. Odată cu îmbunătățirea memoriilor, se perfecționau și procesoarele, păstrând raportul. De fapt, cum devine posibil să se pună mai multe circuite pe un chip, proiectanții de procesoare folosesc aceste noi facilități pentru operații de tip pipeline și superscalare, făcând procesoarele

și mai rapide. Proiectanții memoriilor au folosit de obicei noile tehnologii pentru creșterea capacității chipurilor, nu a vitezei, așa că problema a avut o tendință de înrăutățire în timp. În practică, acest dezacord înseamnă ca după ce procesorul emite o cerere către memorie, nu va primi cuvântul necesar timp de multe impulsuri de tact. **Cu cât memoria este mai lentă, cu atât mai multe tacturi va trebui să aștepte procesorul.**

Sunt două moduri de a trata problema. Cel mai simplu mod este de a începe citirile READ de memorie când sunt întâlnite, dar continuând execuția și **oprirea procesorului dacă o instrucțiune încearcă să folosească cuvântul de memorie înainte ca acesta să fie disponibil.** Cu cât este mai lentă memoria, cu atât mai des va apărea această problemă, și cu atât mai grave vor fi sancțiunile atunci când se întâmplă. De exemplu, dacă memoria are o întârziere de 10 perioade de tact, este foarte probabil ca una dintre următoarele 10 instrucțiuni să încerce folosirea cuvântului citit.

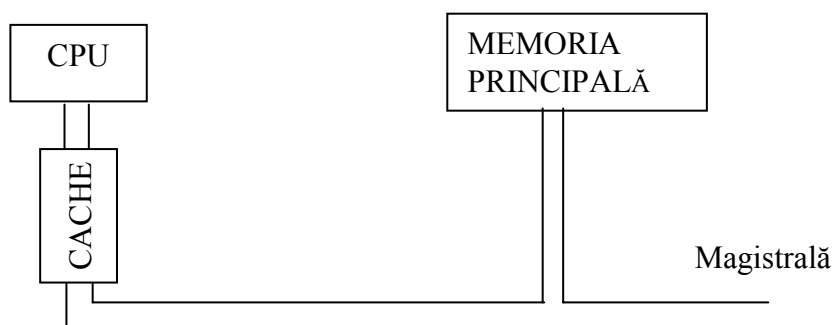
Cealaltă soluție este să se instruiască mașinile să nu se oprească, dar în schimb să ceară **compilatoarelor să nu genereze cod care să folosească cuvintele înainte de a sosi din memorie.** Problema este că această abordare este mult mai greu de implementat. După o instrucțiune de încărcare, nu mai este nimic de făcut, și compilatorul este forțat să adauge instrucțiunea NOP (no operation), care nu face decât să ocupe o pauză și să piardă timpul. Efectul este că această abordare conduce la o întârziere software, nu hardware, dar pierderea de performanță este aceeași.

De fapt, problema este mai mult de economie (bani) decât de tehnică. Inginerii știu cum să construiască memorii la fel de rapide ca procesorul, dar pentru a funcționa la viteza maximă, acestea trebuie să fie așezate pe chipul procesorului (deoarece transferul prin magistrală este foarte lent). Adăugarea unei memorii mari pe chipul procesorului îl face mai mare, deci și mai scump, și chiar dacă nu costul ar fi problema, există limitări tehnologice ale dimensiunii chipului procesorului. **Astfel, alegerea devine să avem o cantitate mare de memorie lentă, sau o cantitate mică de memorie rapidă.** Bineînțeles că noi am prefera o cantitate mare de memorie rapidă la un preț mic.

Se cunosc tehnici de combinare a unei cantități mari de memorie lentă cu o cantitate mică de memorie rapidă, pentru a obține viteza memoriei rapide (aproape) și capacitatea unei memorii mari, la un preț moderat. **Memoria mică și rapidă se numește cache** (de la francezul "cacher", care înseamnă a ascunde). În continuare vom descrie pe scurt cum funcționează și cum se folosesc memoriile cache.

Ideea de bază a cache-ului este simplă: **cuvintele de memorie cele mai utilizate se pastrează în cache.** Când procesorul are nevoie de un cuvânt, caută întâi în cache. Numai în cazul când cuvântul nu este acolo, îl va cauta în memoria principală. Dacă o parte substanțială a cuvintelor este în cache, timpul mediu de acces poate fi redus foarte mult. Succesul sau eșecul depinde de ce fracțiune din cuvinte utile sunt în cache. De muți ani, oamenii știu că programele nu accesează memoria complet aleator. Dacă se face o referință la locația de memorie de adresă A, **este probabil ca următoarea referință de memorie să fie în vecinătatea** lui A. Un exemplu simplu este chiar programul. În afară de salturi și apeluri de procedură, **instrucțiunile se aduc din locații de memorie consecutive.** Mai mult, majoritatea timpului de execuție al programului este petrecut în **bucle**, unde un număr limitat de instrucțiuni se execută de mai multe ori. Asemănător, un program de manipulare a unei matrice este probabil să acceseze de multe ori aceeași matrice înainte de a trece la altceva.

Observația că referințele de memorie într-un interval de timp scurt tind să folosească numai o zonă mică din memoria totală se numește "**principiul localizării**" și este baza tuturor sistemelor de cache. Ideea generală este că atunci când un cuvânt este accesat, acesta și câțiva dintre vecinii săi sunt aduși din memoria mare și lentă în cache, astfel că data viitoare când se folosește, va fi accesat foarte rapid. O așezare obișnuită a procesorului, memoriei cache și memoriei principale este ilustrată în figura 3.4. Memoria cache se situează logic între CPU și memoria principală. Fizic, sunt mai multe localizări posibile.



**Figura 3.4** Poziția relativă a memoriei cache față de restul componentelor sistemului

Dacă un cuvânt este citit sau scris de  $k$  ori într-un interval scurt, calculatorul va folosi 1 referire către memoria lentă și  $k-1$  referiri către memoria rapidă. Cu cât mai mare este  $k$ , cu atât mai bună este performanța globală.

Putem formaliza aceste operații prin introducerea următoarelor valori:

- $c$  timpul de acces al cache,
- $m$  timpul de acces al memoriei principale,
- $h$  raportul de reușită, care este raportul din toate referirile care pot fi satisfăcute de cache și numărul total de referiri.

În exemplul din paragraful precedent,  $h = \frac{k-1}{k}$ .

Unii autori definesc deasemenea raportul de ratare, care este  $1-h = \frac{1}{k}$ .

Cu aceste definiții, putem calcula timpul mediu de acces,  $t_{med}$ , astfel:

$$T_{med} = c + (1-h)m$$

Dacă  $h \rightarrow 1$ , toate referirile pot fi satisfăcute din cache, și timpul de acces se apropie de  $c$ . Din contră, dacă  $h \rightarrow 0$ , o referire la memoria principală este necesară de fiecare dată, așa că timpul de acces se apropie de  $(c + m)$ ,  $c$  pentru a verifica memoria cache (fără succes), și  $m$  pentru a efectua referirea la memorie.

La unele sisteme, referirea (accesul) la memorie poate fi pornită în paralel cu căutarea în cache, așa că dacă nu se găsește în cache, ciclul de memorie este deja început. Oricum, această strategie necesită ca memoria să poată fi oprită pe parcurs la o găsimă a cuvântului în cache, ceea ce face implementarea mai complicată.



Funcționând pe principiul localizării, memoriile principale și cache sunt divizate în blocuri de dimensiune fixă. Când se vorbește de aceste blocuri în cache, se folosește de obicei denumirea "linii de cache". Când se produce o ratare, toată linia de cache se încarcă din memorie, nu numai cuvântul necesar. De exemplu, pentru o linie de 64 biți, o referință către adresa de memorie 256 va încărca în cache linia formată din octeții de la 256 la 319. Cu puțin noroc, unii dinte octeții din linia de cache respectivă vor fi apelați în scurt timp. **Operarea în acest fel este mult mai rapidă decât extragerea câte unui cuvânt, deoarece este mai rapid să se extragă k cuvinte o dată, decât câte un cuvânt de k ori.**

Proiectarea memoriei cache este un subiect de importanță tot mai mare pentru procesoarele de înaltă performanță.

#### **Observații:**

1. O problemă este dimensiunea memoriei cache. Cu cât este mai mare, cu atât mai bune sunt rezultatele.
2. O altă problemă este dimensiunea liniei de cache. 16 k cache pot fi împărțiți în 1k linii de 16 octeți, 2k linii de 8 octeți și alte combinații.
3. A treia problemă este organizarea, sau cum știe memoria cache ce cuvinte păstrează la un moment dat, deci cum se selectează informațiile de păstrat în cache. Ideal ar fi să se poată face o predicție asupra informațiilor care urmează să fie folosite.
4. O a patra problema este dacă datele și instrucțiunile sunt păstrate în același cache sau nu. Un cache unificat (instrucțiunile și datele folosesc același cache), este mai simplu de proiectat și favorizează automat extragerile de instrucțiuni. Oricum, astăzi tendința este către un cache separat, cu instrucțiunile într-o parte și datele în alta. Aceasta se mai numește arhitectura Harvard. Forța care îi ghidează pe proiectanți în această direcție este răspândirea procesoarelor "pipeline". Zona de instrucțiuni trebuie să acceseze instrucțiunile în același timp în care în zona de operanzi se accesează datele. Un cache separat permite accesul simultan, pe când cel unificat nu. De asemenea, cum instrucțiunile nu se modifică în timpul execuției, conținutul memoriei cache de cod nu va fi niciodată rescris în memorie.
5. În final, o a cincea problemă este numărul memoriilor cache. Este obișnuit astăzi să avem chipuri cu cache primar pe chip, cache secundar în afara chipului, dar în aceeași capsulă, și un al treilea nivel de cache mai departe.