

Curs 9 **Arhitecturi RISC de uniprosesare paralelă a datelor utilizând cuvinte de instrucțiuni foarte lungi tip VLIW**

1. Conceptul de paralelism

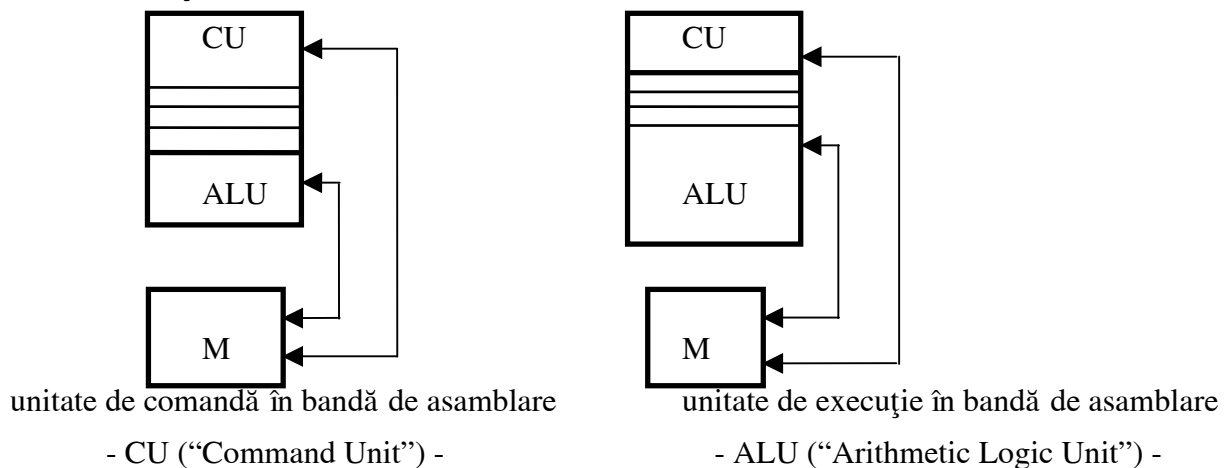
Paralelismul se referă la capacitatea de a suprapune sau executa simultan mai multe operații (citiri de la memorie (“memory fetch”, “load”), scrieri la memorie (“store”), operații aritmetice și logice, operații de intrare/ieșire etc).

Paralelismul poate fi de tip:

a) ***Bandă de asamblare*** (“Pipelining”): este o implementare similară liniei de asamblare, prin utilizarea unor tehnici de suprapunere a fazelor de execuție a mai multor instrucțiuni, ceea ce conduce la creșterea performanțelor unităților de execuție EU (“Execution Unit”) și unității de comandă CU (“Command Unit”).

Lucrul în bandă de asamblare presupune divizarea unui “task” T în “subtask”-uri T1, ..., Tk și asignarea “subtask”-urilor la un lanț de stații de procesare numite segmente de bandă de asamblare. Paralelismul se obține prin operarea segmentelor simultan. Se pot folosi benzi de asamblare atât la nivelul fluxului de instrucțiuni, cât și la nivelul fluxului de date, după cum se prezintă în figură:

Banda de asamblare poate fi implementată și combinat, atât la nivelul comenzii, cât și la nivelul execuției.

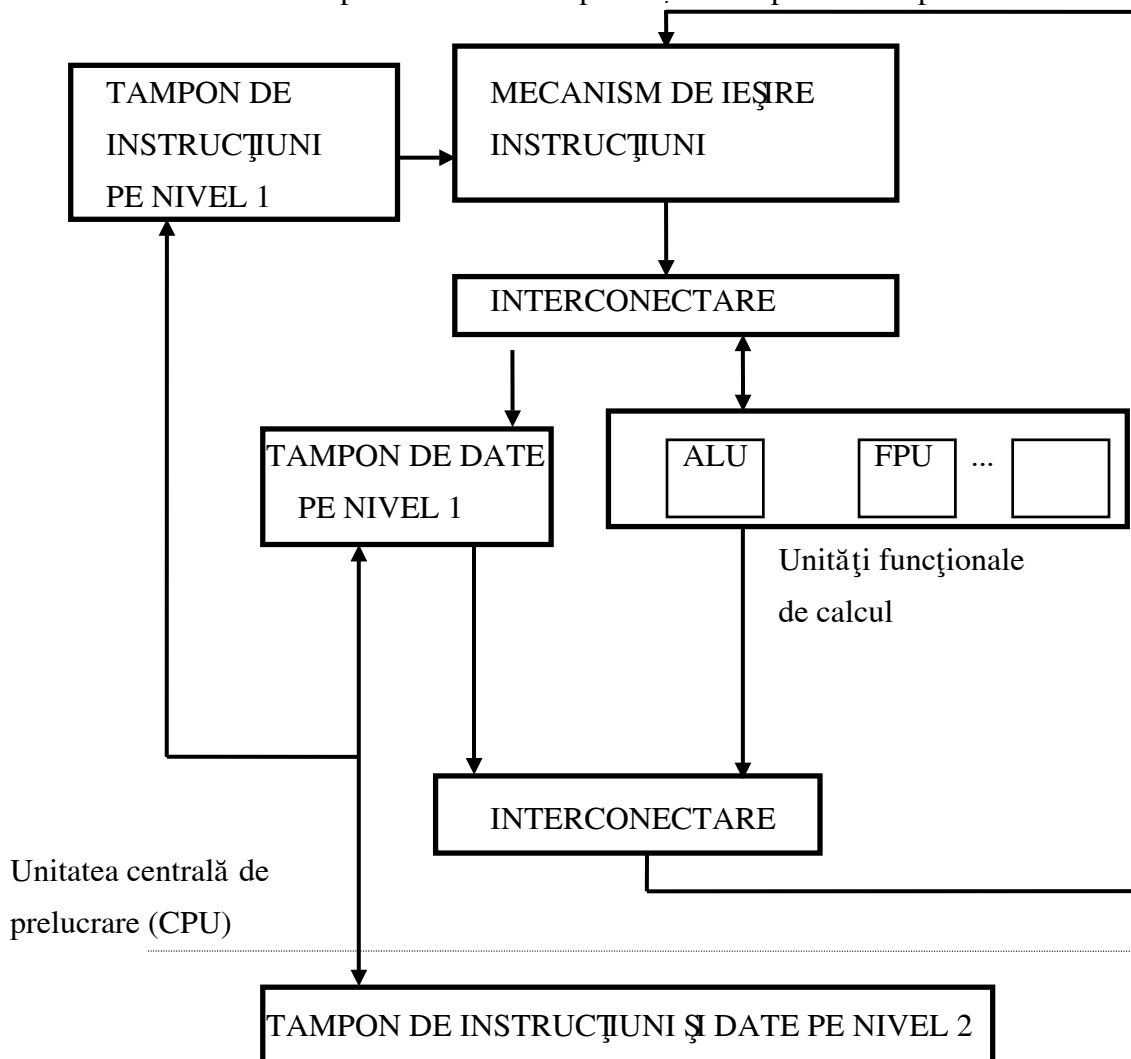


b) ***Funcțional***: prin utilizarea mai multor unități independente, ce realizează funcții (logice/aritmetice) diferite, simultan cu date diferite.

Utilizarea unităților funcționale în bandă de asamblare conferă procesoarelor proprietatea de scalare. În prezent, unitățile funcționale în bandă de asamblare sunt utilizate, atât de RISC, cât și de CISC, tendința actuală fiind de superscalare. Superscalarea [81;83;126;127] este o cale de creștere a performanței prin exploatarea paralelismului de granularitate scăzută (“fine-grain”) și execuția mai multor instrucțiuni pe ciclul de ceas. Acest lucru se realizează prin implementarea

în unitatea centrală de prelucrare, a unor unități funcționale multiple, fiecare dintre ele fiind organizată în bandă de asamblare.

În continuare este prezentat modelul posibil al unui procesor superscalar:



Un astfel de model permite execuția a patru-cinci instrucțiuni într-un ciclu de ceas. Unitățile funcționale multiple de calcul execută instrucțiunile în paralel. Mecanismul de ieșire a instrucțiunilor încarcă, într-o coadă, mai multe instrucțiuni din memoria intermediară de instrucțiuni de pe nivelul 1 și le trimite apoi în execuție la unitățile funcționale disponibile. Memoria intermediară de instrucțiuni de pe nivelul 1 furnizează instrucțiuni mecanismului de ieșire a instrucțiunilor iar memoria intermediară de date de pe nivelul 1 furnizează operanzi pentru unitățile funcționale multiple. Mecanismul de interconectare asigură transferul între componentele unității CPU.

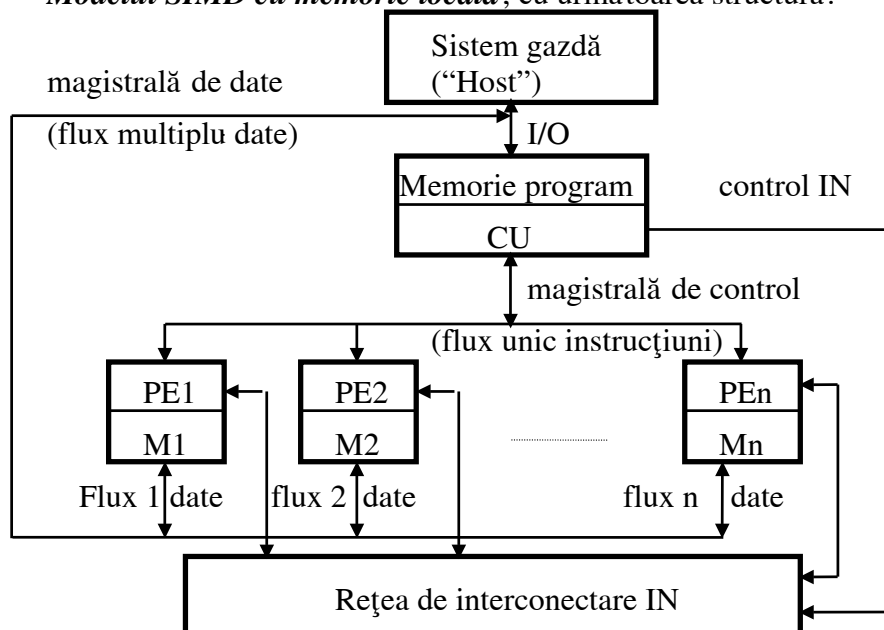
Utilizarea memoriilor intermediare pe două niveluri – nivelul 1 în același circuit cu procesorul (“on-chip”), nivelul 2 în exteriorul circuitului procesorului (“off-chip”), precum și împărțirea memoriei intermediare incorporate în tampon de instrucțiuni și tampon de date, conduc, de asemenea, la creșterea frecvenței de operare a procesorului. Un calcul al influenței memoriilor intermediare asupra creșterii performanței procesoarelor este realizat la sfârșit.

c) **Masiv de procesare SIMD** (“Single Instruction flow Multiple Data flow” – flux singular de instrucțiuni și flux multiplu de date): prin utilizarea unui masiv de elemente procesoare PE (“Processor Element”) identice, care realizează simultan aceeași operație (controlate de aceeași comandă), asupra unor date diferite stocate în memoriile locale.

Sistemele SIMD sunt constituite într-un masiv de elemente de procesoare (PE), elemente de memorie (M), o unitate de control (CU) și o rețea de interconectare (IN – “Interconnecting Network”). Acestea sunt atașate unui sistem gazdă (“Host”) care, din punct de vedere al utilizatorului, este un sistem de tip “front-end”. Rolul lui este de a executa compilări, încărcarea programelor, operații de intrare/ieșire și alte funcții ale sistemului de operare.

În funcție de modul de operare și de control există:

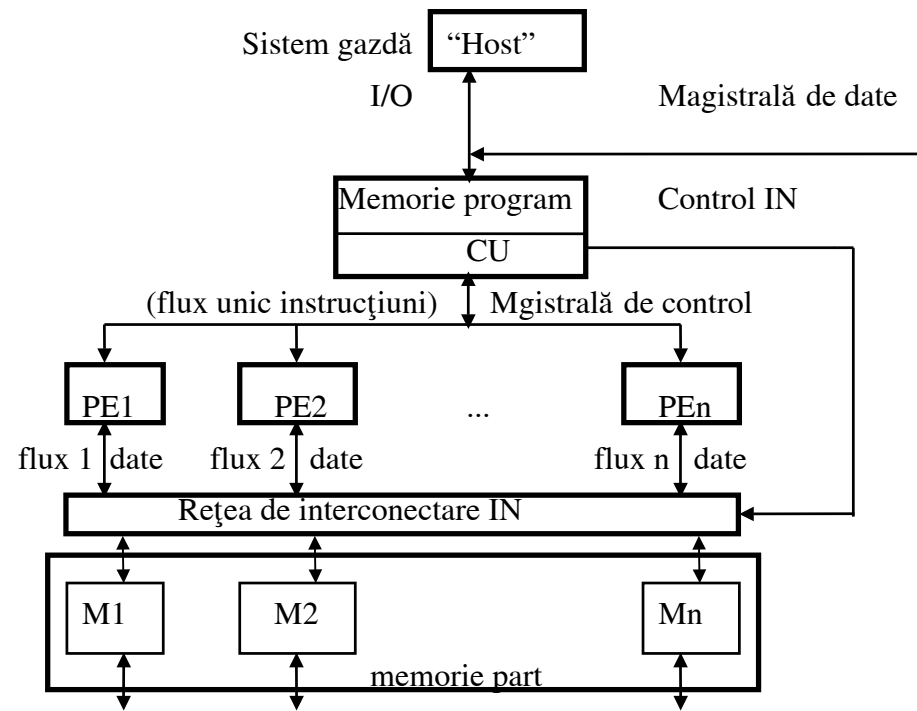
Modelul SIMD cu memorie locală, cu următoarea structură:



Conform acestui model, fiecare procesor are memoria sa locală și execută aceeași instrucțiune furnizată de unitatea de control CU. Unitatea CU citește instrucțiuni din memoria program. Instrucțiunile de control sau cele cu operanzi scalari sunt executate direct de către unitatea CU, iar instrucțiunile cu operanzi vectoriali sunt trimise spre execuție la procesoarele PE. De asemenea, se pot transmite cuvinte de date necesare pentru operațiile scalare cu vectori. Când procesoarele PE își încarcă datele de la memoria lor locală, unitatea CU transmite adresa corespunzătoare procesoarelor PE. Instrucțiunea de control de la CU furnizează aceeași adresă pentru toate PE. Comunicația de date se realizează transferând datele între procesoare prin intermediul unei rețele de interconectare. Fiecare procesor sau grup de procesoare are porturi de comunicație și registre tampon (“buffer”) de date. Rețeaua IN execută câteva funcții de mapare, depinzând de topologia rețelei și controlul cerut de unitatea CU. Procesoarele sunt echipate cu registre de stare cu indicatori pentru testări aritmetice și logice. Granularitatea, numărul de procesoare și tipul rețelei variază de la o implementare la alta.

Model SIMD cu memorie partajată (“shared”)

Conform acestui model, fiecare procesor comunică cu memoria comună partajată, prin intermediul unei rețele de interconectare și execută aceeași instrucțiune furnizată de unitatea de control CU, la fel ca în cazul modelului anterior. Modelul are următoarea structură:



Memoriile sunt separate de procesoarele PE printr-o rețea de interconectare. Un procesor poate accesa oricare memorie iar rețeaua IN permite accesul în paralel la memoria partajată.

Un avantaj este că sunt realizate mai multe funcții de mapare procesor-memorie, ceea ce înseamnă că pot fi mapați ușor în “hardware” mai mulți algoritmi. Acest lucru implică bineînțeles creșterea volumului de comutări de date, ceea ce conduce la creșterea timpului de acces la memorie.

Unitatea CU este similară cu cea din modelul SIMD cu memorie locală, însă rețeaua IN este mai complexă.

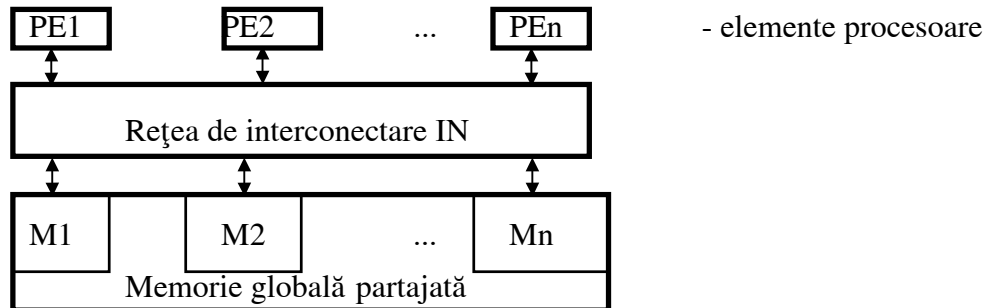
d) **Multiprocesare MIMD** (“Multiple Instruction flow Multiple Data flow” – fluxuri multiple de instrucțiuni și date): prin utilizarea mai multor procesoare, fiecare executând instrucțiuni proprii și comunicând, în general, printr-o memorie comună, sau conectate într-o rețea.

Multiprocesoarele se referă la execuția simultană de “task”-uri pe un sistem de calcul paralel asincron format din mai multe procesoare, fiecare cu CU și ALU, și module de memorie și I/O. Procesoarele cooperează strâns, dar independent.

Există două modele de bază:

Modelul MIMD cu memorie partajată (“shared-memory”)

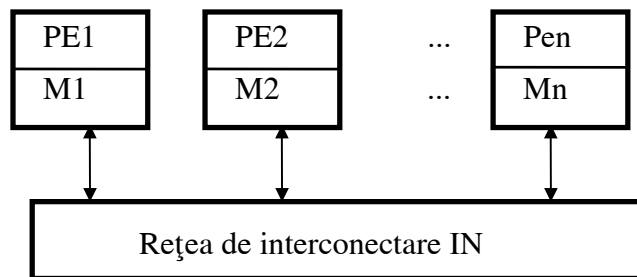
Sunt sisteme strâns cuplate, în care există o conectivitate completă între procesoare și modulele de memorie; au următoarea structură:



PE sunt procesoare de execuție nu neapărat identice, cu sau fără memorie locală. Memoria globală poate fi accesată de către toate procesoarele; creșterea lățimii de bandă la memoria globală se realizează prin implementarea local la procesoare a unor memorii intermediare (“cache”).

Modelul MIMD cu trecere de mesaje (“message-parsing”)

Structura de sistem este următoarea:

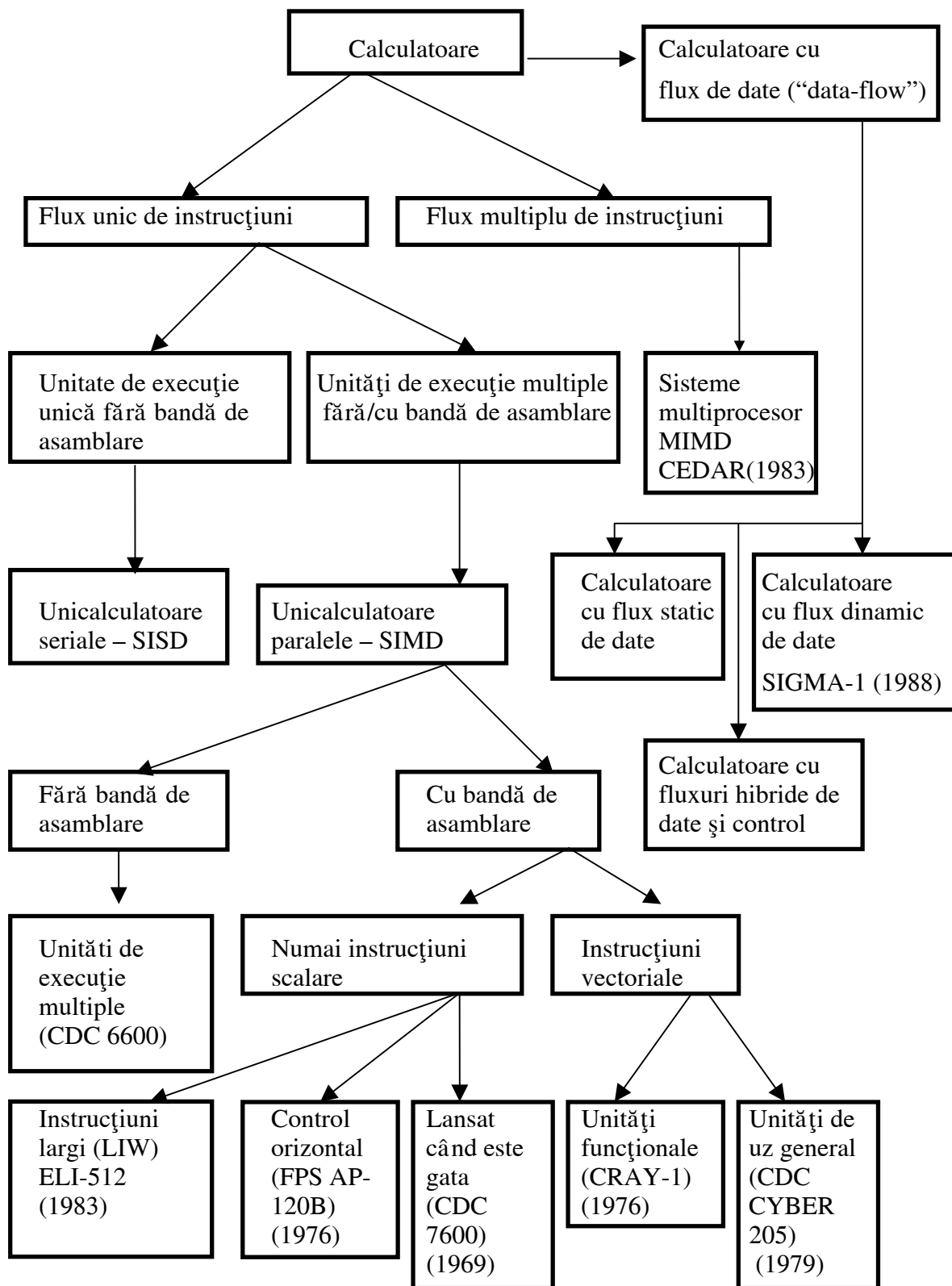


Fiecare modul sistem are un PE, memorie și interfață de I/O. Comunicația de date este dată prin mesaje și nu prin intermediul variabilelor partajate ca la modelul anterior și urmează un protocol de comunicație predeterminat. Gradul de cuplare este mai mic decât la modelul anterior.

Paralelismul este implementat într-un sistem de calcul pe mai multe niveluri, îmbinând mai multe tipuri de paralelism (cum sunt cele de mai sus). Astfel, la nivel de sistem se pot folosi mai multe procesoare, organizate într-un masiv de procesoare sau într-un multiprocesor, iar la nivelul unităților centrale de prelucrare se pot implementa mai multe unități funcționale independente, fiecare organizate în bandă de asamblare.

În acest mod s-a ajuns la obținerea unor performanțe din ce în ce mai ridicate pe sistemele de calcul, cu un cost relativ scăzut.

O clasificare structurală a sistemelor, care reflectă și gradul de paralelizare a prelucrării datelor, este următoarea:



2. Arhitectura VLIW

Prezentare generală

Primele proiecte VLIW (“Very Long Instruction Word”) au fost inspirate din arhitecturile CDC 6600, CRAY-1, IBM 360/91 și MIPS. În anii ’70, multe dintre sistemele masive de procesoare și dintre procesoarele de prelucrare dedicată de semnale foloseau instrucțiuni lungi în memoria ROM, pentru a calcula transformata Fourier rapidă (FFT) și alți algoritmi.

Precursorii adevărați ale procesoarelor VLIW pot fi considerate mașinile RISC (“Reduced Instruction Set Computer”), al căror principiu a fost expus pentru prima oară de către prof. David A. Patterson din University of California de la Berkeley în 1980. Într-un procesor RISC banda de asamblare permite execuția în paralel a fazelor unei instrucțiuni obținând execuția unei instrucțiuni pe ciclu de ceas. Începând cu proiectele RISC - I de la Universitatea din Berkeley în 1980 și MIPS de la Universitatea din Stanford în 1981 și continuând cu procesoarele SPARC de la SUN Microsystems, Transputer T800 de la Inmos, MC 88000 de la Motorola, Am 29000 de la AMD, i860 de la Intel etc., procesoarele RISC au condus la o altă idee, conform căreia să se realizeze citirea simultană a mai multor instrucțiuni și atribuirea lor mai multor unități de execuție, prin intermediul căilor multiple de date, obținând astfel o arhitectură LIW (“Long Instruction Word” – cuvânt lung de instrucțiune) sau VLIW.

Conceptul RISC a evoluat în multe direcții: sunt cunoscute proiectele PIPE (“Parallel Instruction and Pipelined Execution” – instrucțiuni paralele și execuție în bandă de asamblare) de la Universitatea din Wisconsin, RIMMS (“Reduced Instruction Set Architecture for Multi-Microprocessor” – arhitectură cu set redus de instrucțiuni pentru multiprocesoare) de la Universitatea din Reading, SIC (“Single Instruction Computer” – calculatoare cu un singur flux de instrucțiuni), WISC (“Writable Instruction Set Computer” – calculatoare cu set de instrucțiuni inscriptibile). De peste 10 ani conceptul RISC a cucerit piața procesoarelor aducând un mare câștig de performanță pentru companiile ce au investit în el (Hewlett-Packard, IBM, Mips, Motorola, Sun, etc). Una din noutățile din Silicon Valley este proiectarea, de către Intel și HP, a unui nou procesor de tip VLIW, compatibil atât Intel x86, cât și HP Precision Architecture.

Legătura cu mașinile VLIW au realizat-o mașinile RISC superscalare. Acestea sunt mașini care execută instrucțiuni multiple pe un ciclu de ceas. Totuși numai un număr mic de instrucțiuni independente (2 până la 5) pot fi executate într-o singură perioadă de ceas, deoarece ele au puține unități independente de execuție. De exemplu, un astfel de procesor poate avea o unitate de bandă de asamblare pentru prelucrare în virgulă mobilă, două unități de bandă de asamblare pentru prelucrare întregi, o unitate de bandă de asamblare pentru încărcare/memorare și o unitate de procesare ramificații. Inițial, procesoarele VLIW au fost concepute ca o alternativă la procesoarele superscalare RISC. Procesoarele VLIW utilizează, asemenea acestora, principiul microcodificării orizontale pentru a decodifica o instrucțiune VLIW. Procesorul VLIW reprezintă o extensie logică a procesorului RISC. Ca și un procesor RISC superscalar, o mașină

VLIW execută câteva operații simple la un moment dat. Diferența constă în modul de rezolvare a dependențelor la ieșire, care apar atunci când se execută mai multe operații în paralel. În cazul sistemelor VLIW, rezolvarea provine de la compilator, care este responsabil cu împachetarea mai multor instrucțiuni simple într-un cuvânt de instrucțiune lungă. Compilatoarele de VLIW sunt responsabile pentru a determina care instrucțiuni depind de altele.

Mașinile VLIW necesită construirea de unități logice multiple pentru a păstra instrucțiunile împachetate într-o instrucțiune largă. Acest lucru necesită spațiu suplimentar pe “chip”. Tehnologia RISC permite minimizarea acestui spațiu. Una din primele mașini cu instrucțiuni lungi, conținând operații multiple pe instrucțiune, a fost AP - 120B, de la Floating Point Systems din 1981. O altă mașină cu instrucțiuni lungi a fost Intel iWarp din 1988. Primul proiect VLIW a fost proiectul ELI - 512 de la Yale din 1983. Primele mașini adevărate VLIW au fost minisupercalculatoarele din anii '80 de la trei companii: Multiflow, Culler și Cydrome. Ele nu au reprezentat un succes comercial deoarece tehnologia existentă atunci nu a permis acest lucru. Totuși, experiența de a scrie compilatoare pentru VLIW este utilizată astăzi în cadrul proiectelor VLIW actuale (în acest sens trebuie subliniate eforturile companiei HP în realizarea unor compilatoare VLIW performante).

Intrarea arhitecturii VLIW pe piața comercială s-a realizat mai întâi în zona procesoarelor incorporate (“embedded processors”), cu funcționalitate integrată (cum sunt procesoarele TriMedia și TMS320C62x). În prezent, arhitectura VLIW a pătruns și pe piața comercială a calculatoarelor “desktop” (prin procesoarele IA-64 sau Itanium de la Intel și Crusoe de la Transmeta).

Procesoarele VLIW au mai multe caracteristici:

În primul rând, ele au unități funcționale independente multiple și, de aceea, pot executa simultan multe operații pe aceste unități. Operațiile sunt împachetate într-o instrucțiune foarte lungă. O astfel de instrucțiune poate include operații aritmetice (cum ar fi operații cu întregi și în virgulă mobilă), referiri la memorie, operații de control, deplasări de date și ramificații. O instrucțiune foarte lungă trebuie să posede un set de câmpuri pentru fiecare unitate funcțională; aceasta poate avea între 256 și 1024 biți sau chiar mai mult.

Aceste procesoare posedă numai o unitate de control, ce generează un cuvânt de comandă lung, care controlează explicit unitățile funcționale prin câmpuri independente. Astfel, ele pot executa numai un singur flux de cod, deoarece controlorul global selectează pentru execuție o singură instrucțiune foarte lungă, pe fiecare perioadă de ceas. În plus, fiecare operație poate fi executată în bandă de asamblare și solicită pentru execuție un număr mic de perioade de ceas.

Procesoare utilizează, pentru execuția în paralel a operațiilor, un “hardware” simplu și fără interblocare. Dependențele de date și resurse sunt rezolvate înainte de execuție și controlate

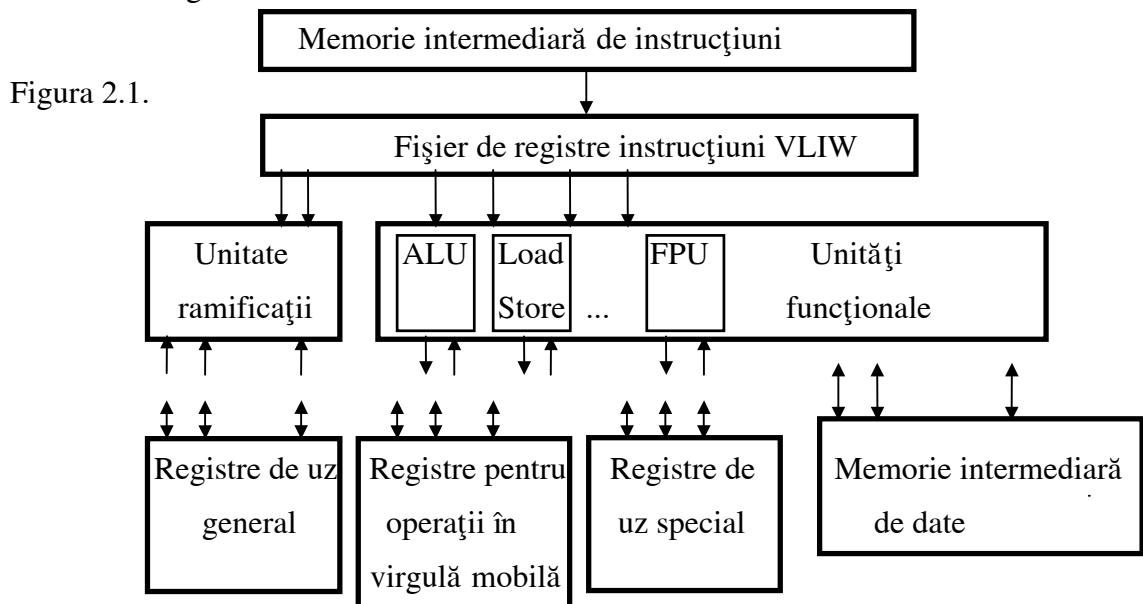
explicit de către instrucțiuni. De asemenea, procesoarele utilizează căi paralele de date pentru unitățile funcționale independente multiple și un volum mic de logică de control și sincronizare.

Prezentare funcțională

Arhitectura VLIW reprezintă una dintre cele mai eficiente soluții în proiectarea microprocesoarelor. Pentru ca un procesor să lucreze mai repede există două posibilități: creșterea frecvenței ceasului și execuția mai multor operații pe fiecare ciclu de ceas. Creșterea frecvenței ceasului necesită inventarea unor procese de fabricare și adoptarea unor arhitecturi (cum sunt benzile de asamblare mai lungi), care să mențină circuitul cât mai ocupat cu execuția „task”-urilor. Execuția mai multor operații pe un ciclu de ceas necesită integrarea de unități funcționale multiple pe același „chip”, care să permită execuția în paralel a „task”-urilor. Problema planificării „task”-urilor este astfel crucială în proiectarea procesoarelor moderne. Procesoarele superscalare actuale realizează acest lucru prin „hardware” pentru rezolvarea dependențelor de instrucțiuni. „Hardware”-ul de planificare crește însă geometric cu numărul de unități funcționale și conduce la limitări de implementare.

Alternativa este de a lăsa „software”-ul să facă planificarea „task”-urilor, ceea ce chiar realizează arhitectura VLIW. Astfel, un compilator optimizat poate examina programul, poate găsi toate instrucțiunile fără dependențe și le poate împacheta într-o instrucțiune foarte lungă, pentru ca apoi să le execute concurrent pe un număr egal de unități funcționale. O astfel de instrucțiune lungă (meta-instrucțiune) conține multe câmpuri mici, fiecare codificând direct o operație pentru o unitate funcțională particulară.

O arhitectura generală VLIW ar fi:



Un astfel de procesor VLIW constă într-o colecție de unități funcționale (sumatoare, multiplicatoare, anticipare ramificații etc), conectate printr-o magistrală, plus registre și memorii intermediare; acest lucru este foarte bun pentru că se câștigă „hardware” iar durata execuției este limitată doar de unitățile funcționale proprii.

O funcție importantă a procesoarelor VLIW este aceea că pot implementa vechile seturi de instrucțiuni CISC mai bine chiar decât procesoarele RISC. Aceasta pentru că programarea procesorului VLIW este foarte asemănătoare cu scrierea de microcod, putându-se conserva astfel programul, utilizându-se instrucțiunile CISC complexe (cum ar fi LODS și STOS de la Intel x86) implementate de procesoarele CISC prin microprograme de microcod ROM pe circuitul procesor.

Microcodul este cel mai scăzut nivel de limbaj, sincronizând porți și magistrale și transferând datele între unitățile funcționale. Procesorul RISC a eliminat microcodul în favoarea logicii cablate mai rapide, iar procesorul VLIW a scos microcodul în afara chipului procesor și l-a plasat la compilator; astfel, emularea de cod CISC se execută mai rapid pe VLIW decât pe RISC.

În schimb, compilatorul este mai complex decât la sistemele RISC; de aceea, se poate realiza o îmbinare “hardware-software” a implementării VLIW (în capitolele următoare sunt prezentate astfel de tehnici). Există o varietate de arhitecturi, care demonstrează multe sau toate caracteristicile descrise mai sus. Un model tipic pentru o arhitectură VLIW ar putea fi:

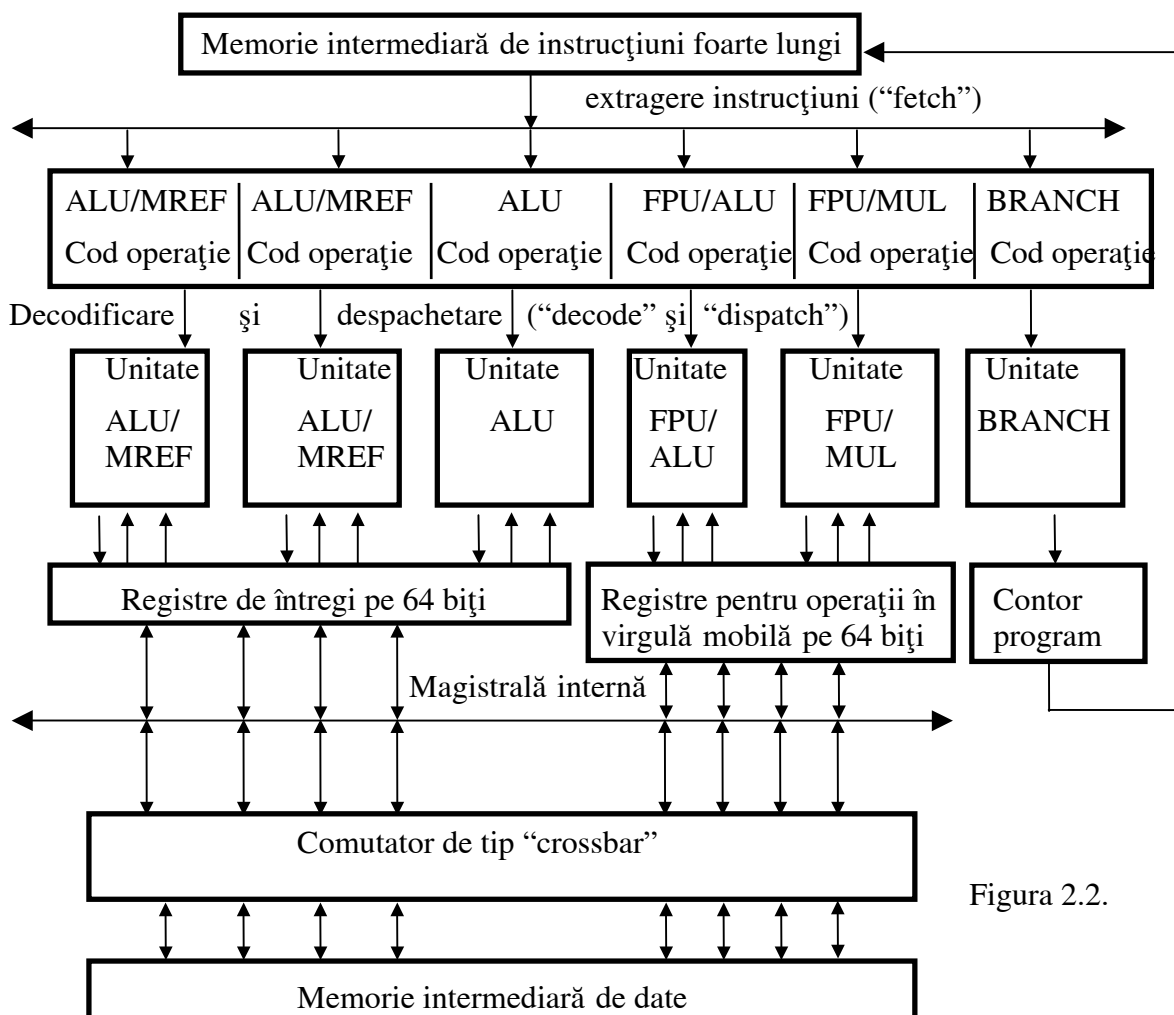


Figura 2.2.

Legendă:

ALU/MREF – operație aritmetică-logică/ referire la memorie (“load/store”)

FPU – operație în virgulă mobilă

FPU/ALU – operație de adunare/scădere în virgulă mobilă

FPU/MUL – operație de multiplicare în virgulă mobilă

BRANCH – ramificație

- codul operației conține și biții de control ai operațiilor

Căile de date constau într-un număr de unități funcționale (pentru întregi, virgulă mobilă și de salt), conectate prin registre globale și speciale multiport și prin intermediul unei rețele de comutare de tip “crossbar” la o memorie intermediară de date. Fiecare unitate funcțională este controlată de un set independent de câmpuri ale instrucțiunii, extrasă din memoria intermediară de instrucțiuni VLIW în registrul instrucțiunii, și de registre specifice de adresă pentru operanzi și rezultate. Există de asemenea o unitate de salt controlată tot de instrucțiune. Ea generează adresa instrucțiunii următoare bazându-se pe câmpurile instrucțiunii curente și condițiile de cod memorate.

Un procesor VLIW poate avea un singur fișier de registre globale (ca procesorul TRACE/7) sau fișiere multiple. De asemenea poate avea unități omogene (ca la IBM) sau heterogene (ca la TRACE/7). Procesoarele superscalare au unități funcționale multiple și fișiere mari de registre, dar ele utilizează decodificatoare de instrucțiuni complexe pentru extragerea operanzilor și procesarea lor în paralel. Există procesoare care au câteva dintre caracteristicile procesoarelor VLIW. Un exemplu este procesorul Intel i860 care are un mod dual de prelucrare instrucțiuni, în care acționează ca un mic procesor VLIW.

Datorită lungimii foarte mari a cuvântului de instrucțiune și pentru a elimina creșterea complexității “hardware”-ului, procesoarele VLIW utilizează un compilator complex care preia instrucțiunile și le pune într-o secvență de cod liniară, astfel încât ele să formeze un cuvânt de instrucțiune foarte lung, din care instrucțiunile pot fi executate în paralel, evitându-se conflictele de date. Memoriile intermediare se utilizează pentru a crește lățimea de bandă (“bandwidth”), permițând accesul rapid la instrucțiuni și operanzi. Instrucțiunile din cuvântul de instrucțiune foarte lung sunt procesate folosind benzi de asamblare multiple.

Pentru a forma și executa un cuvânt de instrucțiune foarte lung, compilatorul utilizează anumite tehnici, cum ar fi “*trace scheduling*“, “*software pipelining*” și altele, care sunt prezentate într-un capitol următor. Un compilator VLIW împachetează grupuri de operații independente într-un cuvânt de instrucțiune foarte lung, într-un mod care utilizează cât mai eficient unitățile funcționale pe timpul fiecărui ciclu de ceas. Compilatorul descoperă toate dependențele de date și apoi determină cum să rezolve aceste dependențe, reordonând întregul program prin deplasări ale blocurilor de cod; acest lucru este realizat la procesoarele RISC superscalare printr-un “hardware” adecvat de reordonare și anticipare. Pentru procesoarele

Un alt model este modelul VLIW:

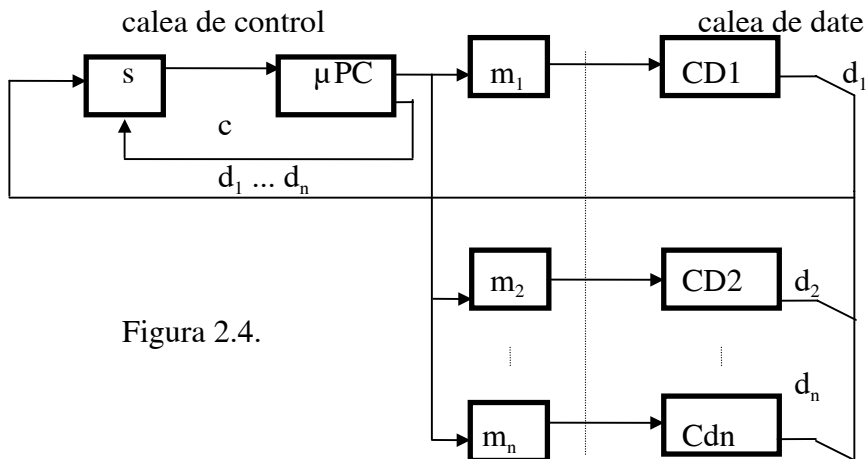


Figura 2.4.

Procesorul VLIW are unități funcționale multiple, fiecare corespunzând căii de date a unui procesor SISD. Modelul căii de control conține o ieșire separată, mapând funcțiile m_1, \dots, m_n pentru fiecare unitate funcțională din calea de date. Funcția de stare următoare s va depinde de starea fiecărei unități funcționale, d_1, \dots, d_n , starea căii de control și intrările externe.

Modelul tradițional SIMD este o simplificare a modelului VLIW. Astfel, ieșirea unei funcții singulare m va fi distribuită la fiecare unitate funcțională. Dacă funcțiile m_1, \dots, m_n ale unui model VLIW sunt identice și egale cu funcția m a unui model SIMD, atunci cele două mașini au o funcționare echivalentă; aceasta implică faptul că se poate programa un procesor VLIW să emuleze funcționarea unui procesor SIMD. De aceea se poate considera mașina VLIW ca un superset funcțional al mașinii SIMD.

Următorul model prezentat este modelul XIMD:

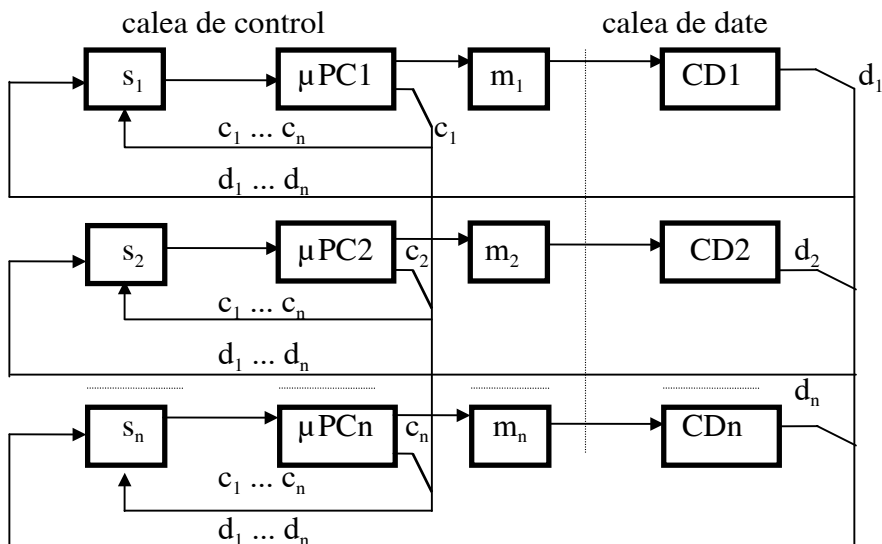


Figura 2.5.

Apare duplicarea secvențiatorului din calea de control (care și ea se duplică). Se realizează o mapare a căilor de date la fiecare secvențiator.

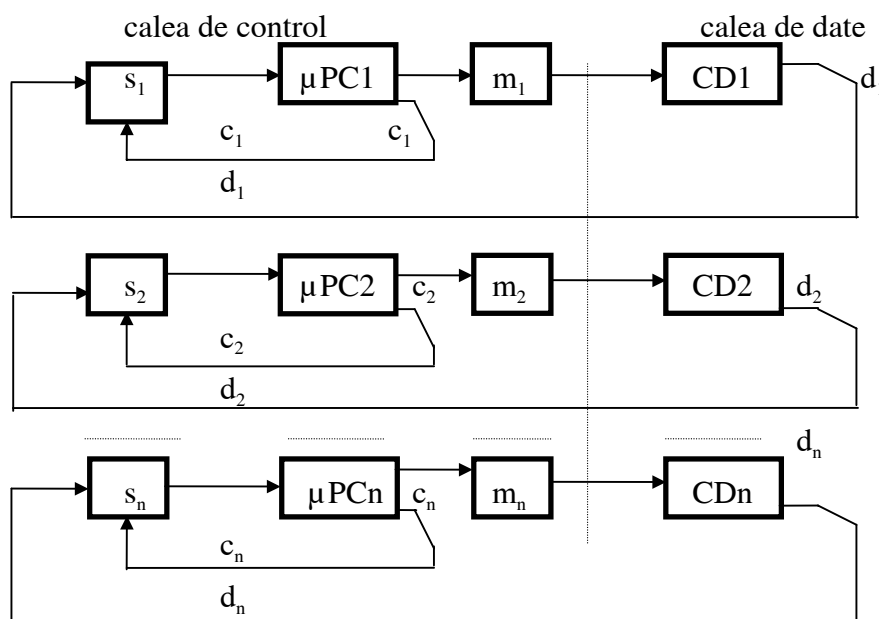
Ca și VLIW, modelul XIMD poate executa operații unice pe fiecare unitate funcțională; în plus, poate executa câte o operație unică de control pentru fiecare unitate funcțională. Se poate considera mașina XIMD ca un superset funcțional de VLIW; astfel, dacă funcțiile m_1, \dots, m_n sunt identice și $\mu PC_1, \dots, \mu PC_n$ au valori inițiale identice, atunci mașina XIMD va fi echivalentă funcțional cu mașina VLIW.

Complexitatea controlului operațiilor, care pot fi executate pe un model XIMD, este afectată de trei factori:

- numărul de valori ale fiecărei variabile de stare a căilor de control;
- numărul de valori ale fiecărei variabile de stare a căilor de date;
- complexitatea și programabilitatea fiecărei funcții de stare următoare.

Primii doi factori afectează volumul de informație care poate fi partajat de-a lungul fluxurilor de instrucțiuni, în timp ce al treilea factor determină cum un flux de instrucțiuni poate diferenția între diferite situații apărute în alte unități funcționale.

Un ultim model reliefat este modelul MIMD:



Selectând ca funcțiile m_i să nu depindă de starea altor unități funcționale, modelul XIMD poate fi funcțional echivalent cu un model MIMD.

3. Tehnici "hardware" de creștere a paralelismului pentru procesoarele VLIW

Utilizarea unui "hardware", care să monitorizeze fluxul de instrucțiuni și să grupeze mai multe instrucțiuni din flux într-o instrucțiune VLIW prin procesorul însuși, constituie o opțiune posibilă pentru furnizarea instrucțiunilor multiple. Aceasta ar permite ca fluxul secvențial de instrucțiuni să fie preluat direct la procesor, cu execuția unor părți de cod accelerate atunci când este posibil, păstrându-se în același timp compatibilitatea.

În 1988, Melvin, Shebanov și Patt au propus un “hardware” pentru compactarea microoperațiilor, generate prin citirea instrucțiunilor, într-un tampon de instrucțiuni decodificate. Implementarea a fost numită unitate de umplere (“fill unit” – FU) și funcționează după cum se descrie în continuare.

Un registru tampon, de preîncărcare instrucțiuni (“prefetch”), citește instrucțiunile din memoria principală sau memoria intermediară de instrucțiuni internă și trimite microinstrucțiunile corespunzătoare la unitatea FU. Aceasta le colectează împreună și le plasează ca un cuvânt multinod într-o memorie intermediară de instrucțiuni decodificate DIC (“Decoded Instruction Cache”). Referințele microoperațiilor au fost redenumite într-un spațiu de adrese al DIC și microoperațiile de la o singură instrucțiune pot fi separate în două cuvinte multinod. Scrierea liniilor din DIC se termină (adică un cuvânt multinod a fost finalizat) ori de câte ori a fost înregistrată o ramificație sau nu a mai rămas în cuvântul multinod nici un slot (fantă de întârziere) de microoperație gol. Microoperațiile din cuvântul multinod pot avea dependențe de date, care pot fi rezolvate mai târziu, printr-o planificare dinamică “hardware” întretesută. Scopul FU este de a se obține un cuvânt atomic larg la planificarea dinamică. Pentru recuperarea stării se utilizează un punct de test (“checkpoint”) de sistem, astfel că mașina poate executa revenire (“back up”) la punctul de test adecvat și poate rula în mod secvențial ori de câte ori apare o excepție. Implementarea conține și posibilitatea redenumirii registrelor, pentru a codifica cererile de accelerare (“forwarding”) de la cuvintele multinod, iar registrele sunt redenumite la nivelul arhitectural al setului de instrucțiuni ISA (“Instruction Set Architecture”), pentru a îmbunătăți performanța.

4. Tehnici ”software” de creștere a paralelismului pentru procesoarele VLIW

Mașinile VLIW pun mare accent pe optimizarea compilatorului, degrevând “hardware”-ul intern de mai multe funcții importante, cum ar fi:

- evitarea dependențelor dintre operații;
- compactarea mai multor operații independente într-o singură instrucțiune lungă, capabilă să se execute complet în paralel (lucru întărit și de planificarea “software” a resurselor necesare pentru executarea în paralel a tuturor operațiilor din instrucțiunea lungă);
- anticiparea cel puțin parțială a ramificațiilor condiționale, extrăgându-se astfel un paralelism ridicat la nivel de instrucțiune.

O instrucțiune VLIW (“Very Long Instruction Word”) poate include mai multe operații, cum ar fi, de exemplu: două operații întregi, două operații în virgulă mobilă, două referiri la memorie și o ramificație; aceste operații se execută în paralel pe resursele procesorului, care sunt planificate de către un compilator optimizat.

Pentru a crește paralelismul procesoarelor superscalare în general, VLIW în special, au fost dezvoltate și se folosesc diferite tehnici de compilare. Fiecare nouă proiectare a unui procesor necesită și îmbunătățirea tehnicilor de compilare.

Tehnicile dezvoltate în acest capitol ajută la creșterea vitezei de execuție (“speedup”) a arhitecturilor superscalare, în general, a celor VLIW, în special.

Scăderea performanțelor procesoarelor superscalare se datorează mai ales următoarelor probleme: limitarea “hardware” a paralelismului; neanticiparea ramificațiilor și neevitarea blocărilor benzilor de asamblare; dificultăți în menținerea unei viteze susținute de execuție în cazul apariției mai multor ramificații condiționale pe ciclul de ceas.

Tehnicile de planificare ”software”, împreună cu un “hardware” adecvat, rezolvă aceste neajunsuri și măresc viteza de execuție.

Din cadrul tehnicilor simple de compilare se pot enumera: asamblarea prin program (“software pipelining”); planificarea urmelor (“trace scheduling”); planificarea prin filtrare (“percolation scheduling”) etc.

În cadrul tehnicilor complexe de compilare se disting: deplasarea speculativă a operațiilor de încărcare/memorare în afara buclelor (“speculative load/store motion”); nespecularea (“unspeculation”); planificarea globală prin descompunerea buclelor și redenumirea registrelor (“global scheduling”); planificarea lărgită (“enhanced pipeline scheduling”); planificare cu gărzi (“sentinel scheduling”); combinarea limitată a operațiilor (“limited combining”); expandarea blocurilor de bază (“basic block expansion”); planificarea cu reacție inversă (“profiling directed feedback”) etc.

Înaintea aplicării acestor tehnici trebuie determinate dependențele programului, prin construirea grafurilor de dependențe pentru programele compilate.

În continuare se prezintă pe scurt aceste tehnici.

Tehnica de asamblare prin program (“software pipelining”)

“Software pipelining” este o tehnică utilizată atât de sistemele uniprocessor, cât și de sistemele multiprocessor. Ea constă în distribuția unei bucle pe mai multe procesoare sau unități funcționale (în cazul procesoarelor VLIW sau alte procesoare superscalare), ceea ce se realizează prin atribuirea stărilor independente la procesoare sau unități funcționale independente și, apoi, înlănțuirea continuă a operațiilor.

Tehnica “software pipelining” reprezintă o metodă pentru optimizarea buclelor, prin suprapunerea operațiilor provenind de la diferite iterații. Acest lucru este realizat prin partiționarea corpului unui program în segmente, care pot fi procesate în mod bandă de asamblare (“pipeline”); buclele sunt reorganizate astfel încât fiecare iterație din codul optimizat prin “software pipelining” să fie compusă din secvențe de instrucțiuni alese din diferite iterații

ale segmentului original de cod. Există un anumit cod de de activare, care este necesar înaintea începerii buclei, ca și un cod pentru terminare după ce bucla este compilată.

Planificatorul de “task”-uri întrețese instrucțiuni de la diferite iterații ale buclei, asamblând operațiile de încărcare de la memorie, apoi toate operațiile aritmetice-logice, iar, în final, toate memorările. O buclă asamblată întrețese instrucțiuni provenind din diferite iterații, fără a se executa desfășurarea ei. Metoda reprezintă o traducere prin program a algoritmului Tomasulo, executat prin “hardware”. Bucla asamblată prin program poate conține, de exemplu, o încărcare (“load”), o adunare (“add”) și o memorare (“store”), fiecare dintr-o iterație diferită.

Diagrama spațiu-timp pentru transformarea “pipelining” constă în următoarele faze: toate stările din buclă sunt executate simultan, pe mai multe procesoare sau unități funcționale, această fază corespunde distribuirii stărilor în spațiu (la procesoare sau unități funcționale); apoi, iterațiile sunt asamblate în timp. Timpul de execuție depinde de numărul de iterații ale buclei și numărul de stări independente pentru fiecare iterație, care furnizează adâncimea maximă a unităților în bandă de asamblare.

Tehnica “software pipelining” îmbunătățește performanța din două motive:

- se îmbunătățește execuția în bandă de asamblare prin planificarea operațiilor cu dependențe mai îndepărtate;
- crește paralelismul datorită segmentelor multiple ale buclei, care pot fi executate simultan pe o mașină superscalară.

Desfășurarea buclelor (“loop unrolling”)

“Desfășurarea buclelor” este o metodă utilizată în general pentru bucle fără dependențe de date între iterații. Prin această tehnică se expandează iterațiile unei bucle de multiple ori și se replanifică instrucțiunile din bucla desfășurată, crescând astfel distanța lexicală dintre instrucțiunile cu dependențe.

Această optimizare este disponibilă datorită compilatorului. Desfășurarea buclei este valabilă numai dacă instrucțiunile desfășurate sunt planificate în vederea minimizării dependențelor. O problemă majoră a desfășurării buclei este faptul că numărul de registre cerute crește. Pentru a evita conflictele de utilizare a registrelor la execuția concurentă a iterațiilor, se utilizează registre diferite pentru fiecare iterație.

În programele reale nu se cunoaște, de regulă, marginea superioară a unei bucle; se presupune că este n și se dorește desfășurarea unei bucle de k -ori; în loc de o singură buclă desfășurată, se va genera o pereche de bucle. Prima buclă execută de $(n \bmod k)$ -ori și are corpul ca al buclei originale. Această versiunea desfășurată a buclei este înconjurată de o a doua buclă, exterioară ei, care iterează de $(n \div k)$ -ori. Numărul de registre poate limita desfășurarea printr-un factor de desfășurare UF (“unrolling factor”); aceasta se întâmplă mai ales la operații în

virgulă mobilă. Majoritatea programelor sintetice de test pot fi desfășurate de mai multe ori, deoarece corpul buclei conține un număr mic de operații.

Un avantaj al desfășurării buclei este faptul că fiecare iterație desfășurată necesită numai o buclă de control. Câștigul de planificare pe bucla desfășurată este chiar mai mare decât pe bucla originală. Aceasta pentru că desfășurarea buclei oferă posibilitatea planificării mai multor calcule. Acest mod de planificare a buclei implică faptul că citirile și scrierile de la/la memorie sunt independente și pot fi interschimbate. Tehnica de planificare prin bucle desfășurate este simplă și utilă, datorită creșterii mărimii fragmentelor de cod, care pot fi planificate efectiv. Această transformare, realizată pe timpul compilării, este similară algoritmului Tomasulo (cu redenumirea registrelor și execuție speculativă “out-of-order”).

Dezavantajul metodei: este dificil de obținut o planificare optimă, pentru că bucla trebuie să fie desfășurată de un număr suficient de ori pentru a ține banda de asamblare complet utilizată.

Acest lucru determină ca metoda să nu fie atât de efektivă ca metoda “software pipelining”.

Tehnica planificării urmelor (“trace scheduling”)

Prin această tehnică se generează un paralelism suplimentar. Tehnica a fost dezvoltată inițial special pentru procesoarele VLIW și reprezintă o combinație de două procese separate. Primul proces, numit *selectarea urmelor* (“trace selection”), încearcă să găsească cea mai potrivită secvență de operații care vor fi asamblate împreună într-un număr mic de instrucțiuni; această secvență se numește *urme* (“trace”). Pentru a genera urme lungi, se folosește tehnica desfășurării buclei (“loop unrolling”) în timp ce ramificațiile buclei sunt executate. O dată ce o urme este selectată, următorul proces, numit *compactarea urmelor* (“trace compaction”), încearcă să comprime urma într-un număr mic de instrucțiuni lungi. Operația de compactare a urmelor deplasează operațiile dintr-o secvență, împachetându-le, pe cât este posibil, în câteva instrucțiuni lungi.

Există două limitări în compactarea urmelor:

- dependențele de date, care forțază o ordine parțială de operare;
 - punctele de salt, care creează locuri de-a lungul cărora codul nu poate fi deplasat ușor;
- în esență, codul trebuie să fie compactat în cea mai scurtă secvență posibilă, care conservă dependențele de date, iar ramificațiile sunt impedimentul principal al acestui proces.

În această tehnică, compilatorul rezolvă conflictele de resurse și utilizează planificarea pe urme pentru ca aplicațiile program să fie partiționate în căi libere de bucle (așa-numitele urme); utilizând algoritmi specifici, compilatorul determină care este cea mai bună urme de executat.

Această planificare poate include o deplasare de cod, care ar putea cauza inconsistențe de date când se execută ramificații și, pentru a compensa aceste lucruri, compilatorul adaugă

operații suplimentare de la alte urme. Procesul este repetat până când toate urmele sunt planificate. Compilatorul verifică de asemenea dacă există referiri simultane la memorie, pentru a nu le executa în paralel.

Marele avantaj al tehnicii “trace scheduling” față de tehnica “software pipelining” este faptul că tehnica “trace scheduling” include o metodă care deplasează codul de-a lungul ramificațiilor.

Tehnica de planificare prin filtrarea fluxului de instrucțiuni (“percolation scheduling”)

“Percolation scheduling” reprezintă o tehnică folosită efectiv pentru planificarea operațiilor multiple provenite din codurile scalare. Ca și tehnica “trace scheduling”, această tehnică extrage operații independente din codul scalar al programului, pe care le assemblează într-o instrucțiune VLIW de lungime corespunzând cu numărul de unități funcționale ale procesorului VLIW. Prin tehnica “percolation scheduling” deplasarea codului este filtrată la fiecare moment, pentru a furniza o instrucțiune lungă, care se poate executa fără întârziere.

Tehnica “percolation scheduling”, ca și cele prezentate anterior, încearcă să crească volumul de paralelism local al instrucțiunii, care poate fi exploatat de către o mașină care execută mai mult de o instrucțiune pe fiecare ciclu mașină.

Deplasarea speculativă de “load/store”-uri în afara buclelor

Această tehnică urmărește evitarea acceselor la memorie din bucle, în special încărcări de la memorie (“load”). Instrucțiunile de încărcare tind să introducă întârzieri într-o operație a procesorului, în special dacă datele cerute nu rezidă în memoria intermediară („cache”).

În esență, aceasta reprezintă o generalizare a ideii de a deplasa instrucțiuni invariante de buclă (“loop-invariant instructions”) în afara buclelor, cu capabilități suplimentare de deplasare a încărcărilor și memorărilor, care sunt executate condițional (adică fac parte dintr-un “if”).

Această strategie poate fi extinsă la proceduri generale, utilizând un program de analiză inter-procedurală (mai ales că nu se degradează semnificativ încărcarea (“overhead”-ul)).

Tehnica nespeculării (“unspeculation scheduling”)

Operațiile speculative (operații ale căror rezultate nu contribuie întotdeauna la rezultatul final al programului) pot apărea în codul intermediar datorită tehnicilor variate de manipulare a codului, care au fost aplicate, sau datorită faptului că astfel de operații au fost prezente în programul original. O secvență de operații devine speculativă dacă este deplasată înaintea unei ramificații condiționate, care impune ca operația să fie executată pentru a obține rezultatul dorit al programului. Astfel, în acele cazuri când ramificația condiționată conduce prin calea unde au fost localizate instrucțiunile, execuțiile instrucțiunilor pot cauza degradarea performanțelor. Obiectivul nespeculării este de a descoperi și reduce numărul operațiilor speculative, aceasta

evitând degradarea potențială a performanțelor, prin gruparea operațiilor speculative în jos pe una sau două destinații ale ramificației condiționate, transformându-le în operații nespeculative. Gruparea instrucțiunilor conduce la posibilitatea existenței unui număr de blocuri de bază cu o singură intrare și o singură ieșire, care buclează; stările imbricate “if-then-else-endif” sunt exemple de asemenea grupuri.

Tehnicile de planificare globală (“global scheduling”) și de planificare extinsă (“enhanced pipeline scheduling”)

Regiunile din programe sunt compactate prin combinarea planificării globale (“global scheduling”) și a planificării extinse (“enhanced pipeline scheduling”), începând de la regiunile cele mai interioare (bucle) și sfârșind cu cele mai exterioare (întreaga procedură).

Aceste două tehnici au fost folosite cu succes în compilatorul xl_c pentru sistemul IBM RS/6000.

Prin planificarea globală este paralelizat codul secvențial eliminându-se antidependențele și dependențele de ieșire prin redenumirea registrelor, așa cum se cere în timpul planificării. Un atribut al fluxului de date, la începutul fiecărui bloc de bază, indică ce operații sunt disponibile pentru deplasarea în sus prin acest bloc de bază. Planificarea globală constă în alegerea celei mai puțin utilizate operații din setul de operații, care se poate deplasa la un moment dat, deplasând toate operațiile în acel punct, contabilizând copiile pentru marginile care unesc căile deplasării de cod, dar nu sunt pe ele și actualizând attributele fluxului de date ale blocurilor de bază numai pe căile care au fost traversate de către operațiile deplasate. Deplasările de cod sunt globale, fără trecerea prin transformările atomice ale planificării prin filtrare (“percolation scheduling”), aceasta pentru o mai bună eficiență. Pentru executarea deplasărilor de cod, se utilizează o reprezentare intermediară care este direct executabilă mai degrabă ca un cod secvențial RISC, decât ca un cod VLIW. Algoritmul poate fi utilizat și pentru a genera cod paralelizat pentru procesoare superscalare cu unități ALU multiple.

Pentru a compacta buclele, planificarea extinsă plasează o barieră (gardă) la punctul curent de planificare, și lasă planificarea globală să caute cea mai bună operație pe toate căile care pot traversa marginile din spate ale buclei, dar nu bariera. Buclele sunt desfăcute anterior planificării și este efectuată redenumirea domeniului de lucru pentru a crește oportunitățile de planificare.

În continuare este analizată tehnica de planificare extinsă a benzii de asamblare (“enhanced pipeline scheduling”).

Această tehnică este utilizată pentru planificarea buclelor conținând constructori condiționali. Are avantajul că corpul buclei rămâne intact în timpul planificării și astfel banda de asamblare poate fi oprită după orice pas terminat.

Modelul mașini de stare este următorul:

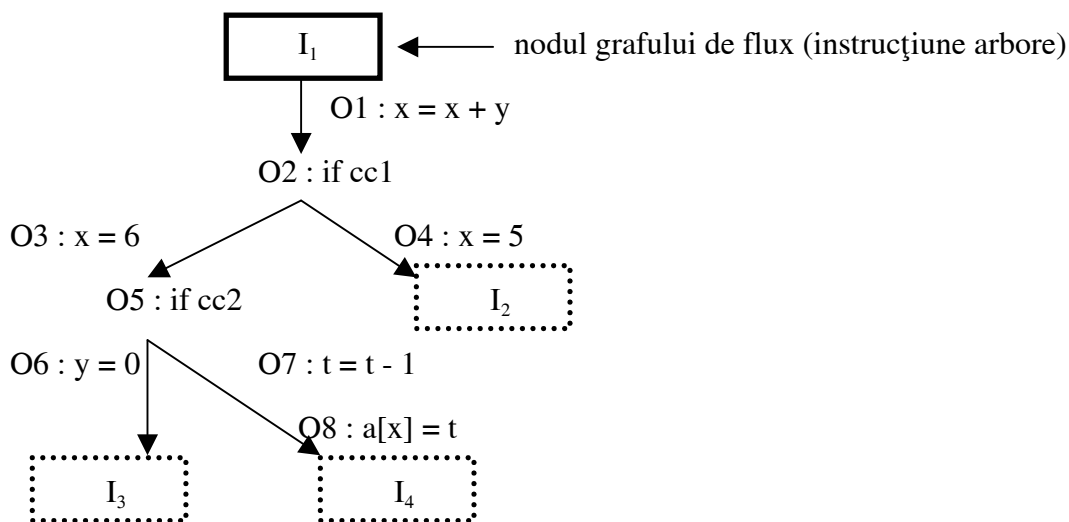


Figura 4.1. Instrucțiune arbore (“tree instruction”) utilizată în algoritmul planificării extinse

$I_2 - I_4$: etichetele altor noduri din graful de flux

cc1, cc2: coduri condiție ce sunt calculate înainte ca nodul să fi intrat în execuție și numai în acele stări, de-a lungul unei căi de la rădăcină la o frunză, care sunt executate.

Toate operațiile pe căile selectate sunt executate concurrent, utilizând vechile valori pentru toți operanzii. De exemplu, în figura 4.1, atribuirea lui t prin $O7$ nu afectează utilizarea lui t prin $O8$. După ce toate operațiile au fost executate, rezultatele sunt scrise la registre sau memorie, și controlul este transferat de la instrucțiunea corespunzătoare etichetei la frunză.

Tehnica planificării extinse crește paralelismul pe o arhitectură cu căi multiple de salt, permițând ca mai multe operații de salt condiționat să se deplaseze într-un nod al grafului de flux. La execuția condiționată, operațiile pot fi plasate pe orice ramificație din nodul grafului de flux, spre deosebire de modelul tradițional de mașină, în care toate operațiile trebuie să preceadă orice operație de salt condiționat.

Acest lucru permite migrarea codului pe arhitecturi mai performante. Migrarea aplicațiilor este asistată automat de calculator prin translatarea codului de pe o mașină pe alta. Procedul poate fi o cale de migrare și între diferite tipuri de arhitecturi (de la CISC la RISC, la arhitecturi superscalare sau la arhitecturi VLIW). Migrarea are loc la timpul de execuție.

Tehnica de planificare cu gărzi (“sentinel scheduling”)

Tehnica a fost dezvoltată de S. A. Mahlke et al. [98] și funcționează astfel:

Funcția “**reducere_dependențe_graf**” deplasează dependențele de control de la un graf de dependențe pentru a permite deplasarea speculativă de cod.

Instrucțiunile pot fi executate speculativ, dacă ele sunt permise ca speculative în arhitectură: de exemplu, ramificațiile, apelările de subrutine și instrucțiunile de intrare/ieșire (i/o) nu pot fi executate speculativ, la fel și instrucțiunile de memorare fără suport de memorie.

Funcția “**sentinel_scheduling**” utilizează graful redus de dependențe pentru a executa planificarea listei tuturor instrucțiunilor dintr-un “superblock”.

Tehnica de combinare limitată a operațiilor (“limited combining”)

Această tehnică scanează codul scalar, cautând oportunități pentru a reduce lungimea căii de cod prin combinarea operațiilor colapsabile (cum ar fi: $A\ r4 = r5$; $B\ r6 = r4@r7$, unde $r4$ este utilizat în noua instrucțiune $C\ r6=r5@r7$). Acest lucru este optimizat, când ambele operații sunt înăuntrul aceluiași bloc de bază. Tehnica de **combinare limitată** poate acoperi un număr de blocuri de bază căutând operații colapsabile, incluzând eventual puncte de intersecție, prin duplicarea codului. Această transformare va **căuta** o secvență de instrucțiuni, începând de la o instrucțiune încarcă imediat la registru (“load immediate to register”) sau copiază registru (“register copy”) (unde registru înseamnă un registru în virgulă fixă, unul în virgulă mobilă, sau o condiție de cod), prin ramificații necondiționate, până la ultima utilizare a registrului destinație pentru acele instrucțiuni. Dacă reușește **căutarea**, întreaga secvență de instrucțiuni (începând cu instrucțiunea următoare primeia și sfârșind cu ultima instrucțiune utilă) înlocuiește instrucțiunea primită inițial. Toate aparițiile registrului destinație, pentru instrucțiunea începută în aceeași secvență, sunt înlocuite prin sursa (literal sau registru) instrucțiunii primite, și este adăugat un salt necondiționat la instrucțiunea următoare ultimei utilizări. Orice cod inaccesibil lăsat de transformare poate fi șters prin tehnici cum ar fi *eliminarea codului inaccesibil* (“unreachable code elimination”).

Tehnica expansiunii blocurilor de bază (“basic block expansion”)

Blocurile de bază pot fi create de optimizări de programe sau chiar programul original; ele sfârșesc cu un salt necondiționat. Aceste ramificații nu sunt luate în considerare de către procesoarele VLIW, dar ele pot încetini procesoarele superscalare pentru că ele consumă resurse importante (de exemplu, dacă un salt condiționat neexecutat este urmat imediat de un salt necondiționat). Salturile condiționate nu pot fi evitate, dar tehnicile de expansiune a blocurilor de bază încearcă să minimizeze apariția salturilor necondiționate la execuție, copiind codul de la destinația ramificației. Pentru a se evita blocarea benzii de asamblare, această tehnică determină numărul de instrucțiuni de neramificație care pot fi deplasate înaintea saltului necondiționat, examinând codul precedent al ramificației. Căutarea se oprește când s-au înlocuit suficiente instrucțiuni de neramificație, când o instrucțiune (dintr-o buclă) este revizitată, sau când se execută o revenire dintr-o procedură. Există o limită a ferestrei de lucru și depășirea ei oprește căutarea; în acest caz, punctul de oprire este cel care determină blocarea minimă. Odată oprită căutarea, codul este copiat, începând de la destinația ramificației necondiționate până la punctul de oprire. Ramificația necondiționată originală este ștearsă și este inserat un nou salt necondiționat după codul copiat, pointând către instrucțiunea imediat următoare punctului de oprire.

Tehnica de planificare cu reacție inversă (“feedback”) – PDF (“Profiling Directed Feedback”)

Informația de “profile” este utilizată pentru a minimiza numărul de execuții ale blocurilor de bază. Din setul de blocuri de bază se vor număra execuțiile unui subset care acoperă toate blocurile (scade astfel “overhead”-ul). Pentru a genera și utiliza codul de “profile”, se compilează un program de două ori și se rulează fiecare versiune separat. În timpul primei treceri a PDF, compilatorul inserează un cod de numărare la timpul de execuție dintr-un subset de blocuri de bază. Când programul astfel compilat este executat, se creează un fișier care indică de câte ori au fost executate blocurile de bază conținând codul de numărare. Numărarea rulărilor multiple ale aceluiași program poate fi acumulată. În timpul celei de-a doua treceri a PDF, compilatorul citește înapoi (“feedback”) din fișierul creat anterior, din același loc în care compilatorul a inserat la prima trecere codul de numărare, calculează setul complet de numărări de blocuri de bază din setul de blocuri de bază de numărare citit din fișier și-l utilizează pentru optimizări.

Nu toate blocurile de bază au nevoie însă de cod de numărare. Este posibil să se micșoreze volumul de instrucțiuni numărate într-o buclă (reducând astfel și “overhead”-ul PDF) prin deplasarea încărcărilor și memorărilor invariante în afara buclei.

5. Compatibilitatea arhitecturală VLIW

În cazul arhitecturilor VLIW pot apărea probleme de compatibilitate a programelor executate pe mașini de generații diferite. Problema compatibilității arhitecturilor VLIW poate fi rezolvată utilizând diferite scheme “hardware” sau “software” pentru păstrarea compatibilității codului compilat. Astfel de scheme utilizează tehnici “hardware”, cum ar fi *divizarea execuției* (“split-issue”), *unitatea de umplere* (“fill unit”), sau tehnici “software” de recompilare a codului obiect al unui program VLIW dintr-o generație în alta. Înainte de a se prezenta câteva dintre aceste tehnici, este realizată o scurtă caracterizare a mașinilor VLIW din punctul de vedere al problemelor de compatibilitate. Astfel, arhitecturile VLIW sunt mașini orizontale cu cuvânt lung de instrucțiune (format dintr-o multioperație) constând în câteva operații executate simultan, în același ciclu de ceas.

Programele VLIW sunt “latency-cognizant”, adică se cunoaște latența (întârzierea) tuturor unităților funcționale. Arhitecturile VLIW se mai numesc arhitecturi cu latență cunoscută UAL (“Unit Assumed Latency”), spre deosebire de arhitecturile cu latență necunoscută NUAL (“Non-Unit Assumed Latency”), care atribuie întârzieri pentru toate unitățile funcționale; majoritatea arhitecturilor superscalare sunt UAL.

Există două modele de planificare pentru programele cu latență cunoscută:

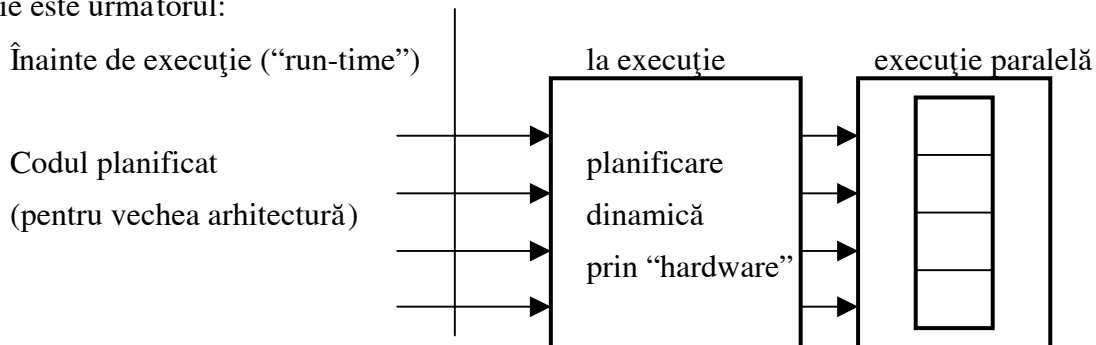
- a) modelul “Equals” (E – egal): oricare operație are exact latența de execuție specificată pentru ea;

- b) modelul “Less-Than-or-Equals” (LTE – mai mic sau egal): latența poate fi mai mică sau egală cu cea specificată.

Modelul E produce o planificare puțin mai scurtă decât LTE în principal datorită reutilizării registrelor. Totuși, LTE simplifică implementarea întreruperilor și furnizează compatibilități liniare atunci când latențele sunt reduse.

Tehnici “hardware”

Aceste tehnici folosesc o planificare dinamică “hardware” a codului. Un model de execuție este următorul:



a) Metoda “split-issue”

Tehnica de planificare dinamică “hardware” prin divizarea execuției (“split-issue”), elaborată de B. Rau, se caracterizează prin următoarele:

Fiecare operație este divizată într-o pereche de operații de tipul “read_and_execute” (RE – citește și execută), “destination_writeback” (DW – scrie destinație). RE utilizează ca destinație un registru anonim iar DW scrie la destinația specificată în operația originală. RE este executată în următorul ciclu de ceas disponibil, în care nu sunt dependențe sau constrângeri de resurse, pe când DW este executată în ultimul ciclu de ceas aflat după întârzierea dată de (ciclul de ieșire + originea – operația – latența – 1). Pentru a se asigura că DW nu se execută mai devreme sunt prevăzuți indicatori.

b) Metoda “fill-unit”

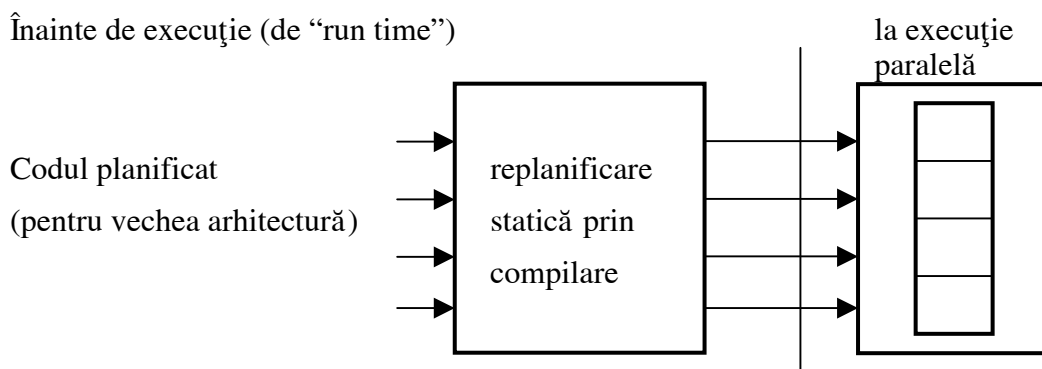
Această metodă a fost descrisă anterior, de aceea se va prezenta doar pe scurt modul de lucru:

- Procesorul execută un flux de operații UAL sau NUAL.
- Concurrent cu execuția, unitățile de umplere compactează operațiile într-o multioperație, care se memorează într-o linie de memorie intermediară ascunsă (“shadow instruction cache”).
- Formarea unei noi multioperații de către unitatea de umplere este terminată când se înregistrează o instrucțiune de ramificație.

Tehnici “software”

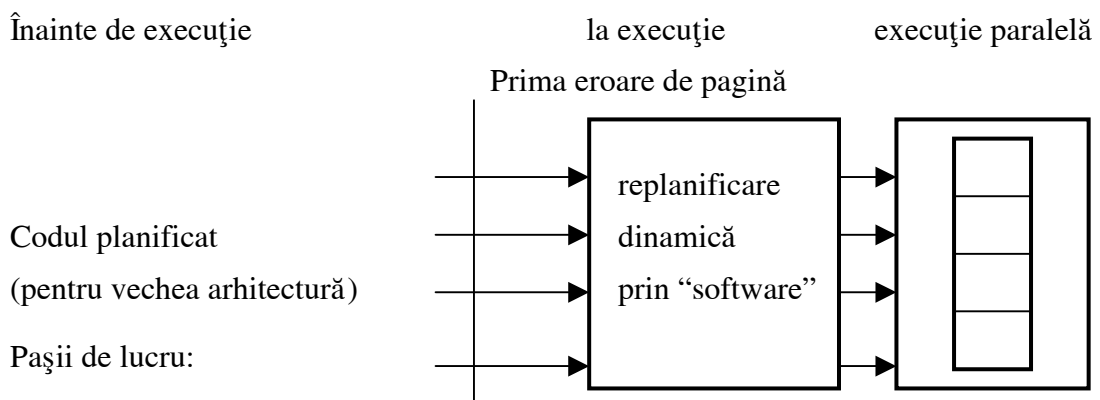
a) *Recompilarea statică*

Această tehnică recompilează întregul program în mod static (“off-line”) și poate lua avantajele compilatoarelor optimizate. Recompilarea completă poate fi evitată prin menținerea copiilor multiple de programe pentru diferitele arhitecturi destinație într-un fișier obiect partiționat; un modul adecvat acestei cerințe poate fi planificat la instalare. Această metodă introduce devierea de la procesul normal pentru dezvoltare și de la procesul rutinei de instalare pentru utilizator, deoarece copiile multiple pot conduce la inconsistențe de date și pentru că spațiul de memorare cerut de copii poate deveni excesiv de mare. Modelul de execuție este următorul:



b) *Replanificarea dinamică*

Tehnica de replanificare dinamică recurge la o versiune limitată de planificare prin program (“software scheduling”) fără a necesita resurse “hardware” suplimentare. Modelul de execuție este următorul:



- Încărcătorul (“loader”-ul) sistemului de operare citește antetul (“header”-ul) programului binar și detectează versiunea generației arhitecturii.
- După ce prima pagină a programului este încărcată pentru execuție, agentul de gestiune a erorii de pagină invocă modulul de replanificare dinamică, care replanifică execuția pe mașina curentă; procesul este repetat pentru fiecare nouă eroare de pagină.
- Paginile traduse sunt salvate pe un spațiu de înlocuire (“swap”); sunt replanificate numai paginile care sunt executate în timpul duratei de viață a programului.

6. Modele arhitecturale VLIW

În acest capitol sunt prezentate diferite modele arhitecturale VLIW, implementate în cadrul unor mașini de uniprocésare paralelă performantă. În primul rând este descris un model VLIW pentru un procesor de semnal. Principiul VLIW este foarte potrivit pentru prelucrări multimedia, acestea presupunând un număr mare de operații necesare a se efectua în paralel. Ca realizare practică sunt cunoscute procesoarele TriMedia, realizat de firma Philips și TMS320C62x, realizat de firma Texas Instruments. Capitolul se încheie cu prezentarea procesoarelor VLIW actuale, IA-64 sau Itanium de la firma Intel și Crusoe (compatibil cu un procesor Mobile Intel Pentium) de la firma Transmeta.

Modelul VLIW pentru un procesor de semnal

Un procesor de semnale video poate fi proiectat conținând mai multe elemente de procesare (PE- "Processing Elements") lucrând în paralel. Semnalul de prelucrat este digitalizat și introdus la procesare prin intermediul unui comutator de tip "crossbar" (cu bare încrucișate). Elementele de procesare execută operațiile în paralel pe baza unor grafuri de flux de semnal SFG ("Signal Flow Graph").

Arhitectura unui astfel de procesor, numit VSP ("Video Signal Processor"), este prezentată în figura următoare:

CS – comutator "crossbar" pentru cele 10 unități:

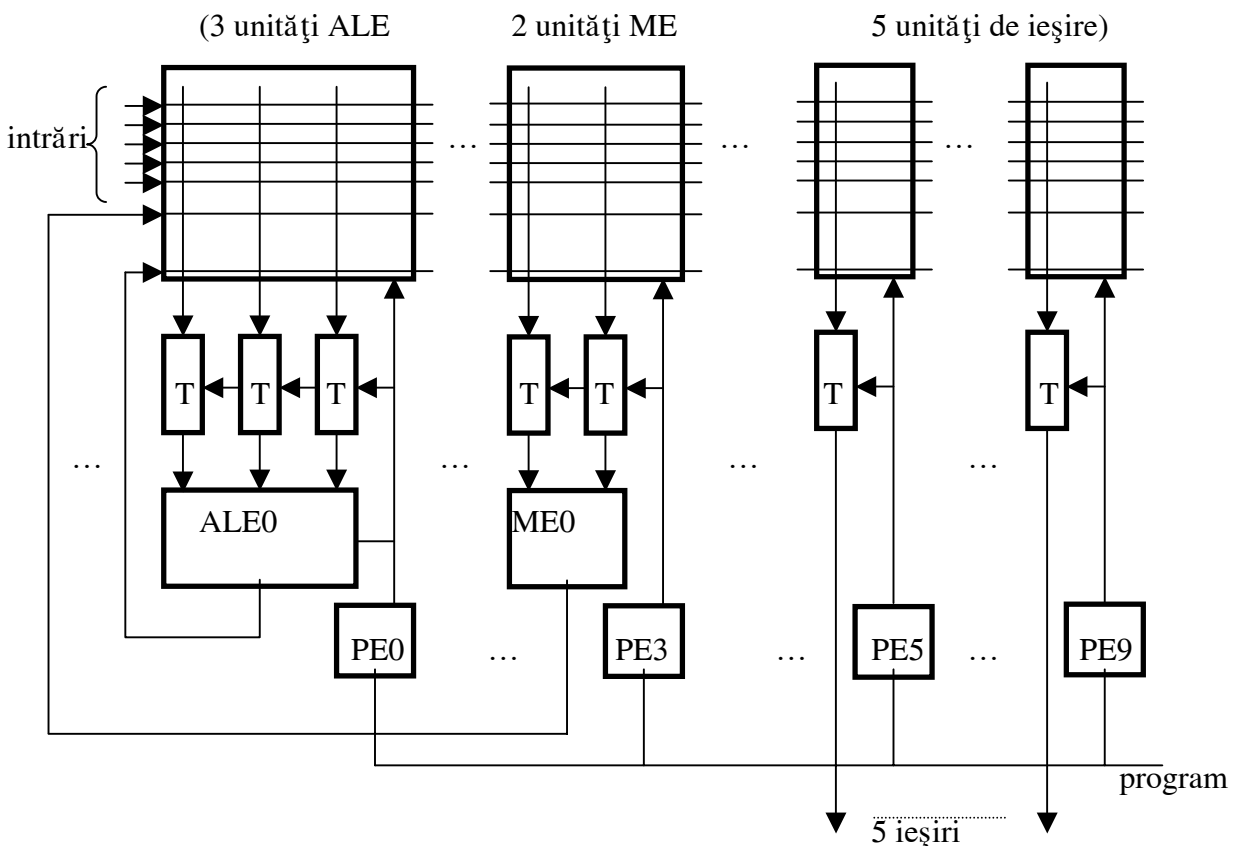


Figura 6.1.

T – tamponare;

CS – comutator cu bare încrucișate ("Crossbar Switch")

Cele 10 PE corespund următoarelor unități:

- 3 elemente logice aritmetice (ALE – “Arithmetic Logic Element) în bandă de asamblare;
- 2 elemente de memorare (ME – “Memory Element);
- 5 elemente de ieșire.

Fiecare dintre ALE, ME și elementele de ieșire are programul său individual de lungime până la 16 instrucțiuni. Un PE execută programul periodic iar frecvența cea mai scăzută este:

$$\frac{27 \text{ MHz}}{16} = 1,6875 \text{ MHz}$$

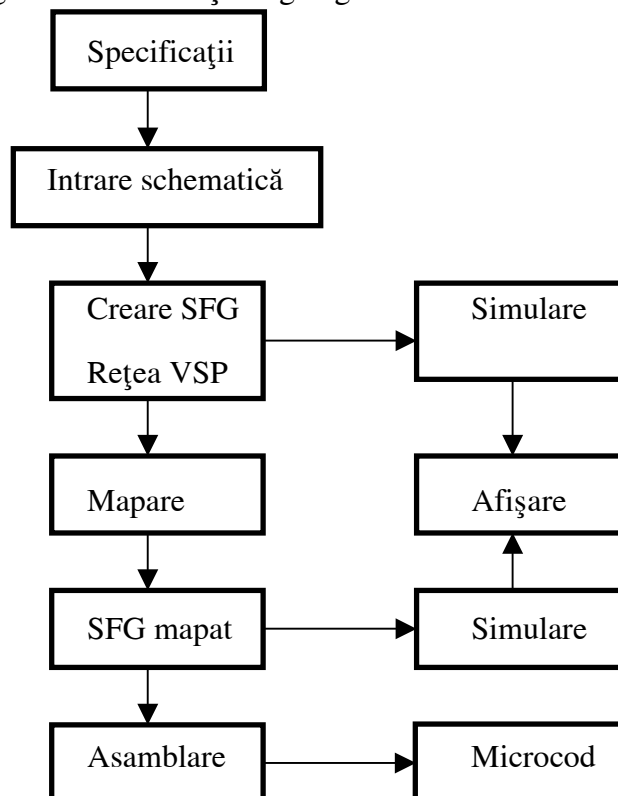
Instrucțiunile pentru ALE au 56 biți, cele pentru ME au 36 biți iar cele pentru elementele de ieșire au 12 biți. O instrucțiune VLIW va avea: $3 \times 56 + 2 \times 36 + 5 \times 12 = 300$ biți

Programele PE sunt încărcate serial printr-o interfață. Sistemele pot avea 1, 3, 5, 8 sau 16 procesoare de semnale video (VSP).

Grafurile de flux de semnal (SFG) descriu algoritmiile programelor de prelucrare a semnalelor. Pașii de mapare a SFG pe rețeaua VSP sunt:

- asignare PE: atribuie fiecărei operații elementul PE pe care se execută
- asignare timp: atribuie fiecărei operații timpul de start al primei execuții (în cicli de ceas)
- asignare cale: atribuie fiecărei precedente etichetate de date calea de urmat prin VSP

Traectoria programelor urmărește organigrama:



Exemple de implementări actuale de procesoare de semnal, având nucleul de procesare de tip VLIW:

Procesorul Texas Instruments TMS320C62x

Acest procesor conține 8 unități de execuție independente împărțite pe două căi de date, fiecare având aceleași patru tipuri de unități funcționale: L (adunare întregă, operații logice, numărare pe biți, adunare în virgulă mobilă, conversie din virgulă mobilă); S (adunare întregă, operații logice, operații pe biți, deplasări, constante, ramificații/control, comparări în virgulă mobilă, conversii din virgulă mobilă, împărțiri în virgulă mobilă); D (adunare întregă, încărcare – memorare); M (înmulțire întregă și în virgulă mobilă).

Frecvența de ceas a procesorului este de 200 MHz, suficientă pentru procesări de semnale.

Fiecare cale de date are un fișier de 16 registre pe 32 biți cu 10 porturi de citire și 6 porturi de scriere; există posibilitatea de a transfera valori între cele două fișiere de registre de pe cele două căi de date, prin intermediul unei magistrale interne globale.

Schema de execuție este următoarea:

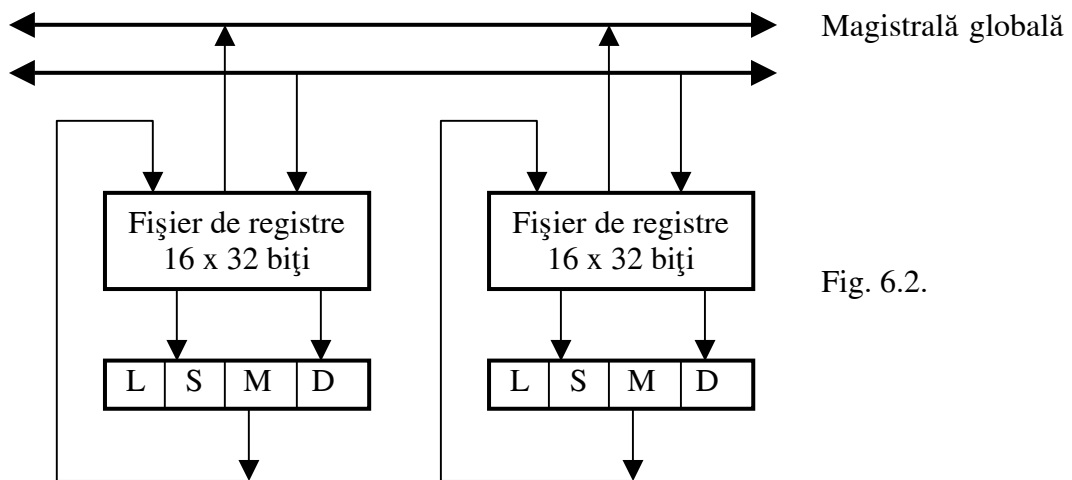


Fig. 6.2.

Calea internă de execuție are 256 biți și se pot executa 8 operații pe fiecare perioadă de ceas, cu operanzi pe 32 biți.

Din memoria internă de instrucțiuni se citesc pachete conținând până la 8 instrucțiuni, care pot fi citite simultan. Fiecare pachet citit se poate expanda în până la 8 pachete de execuție, funcție de biții de paralelizare asociați. Un pachet de execuție are 8 instrucțiuni care se pot executa în paralel, cum s-a specificat mai sus.

Schematic, un pachet citit se poate extinde în pachete de execuție ca în exemplul următor:

Fie următorul pachetul citit:

A B C D E F G H

1 1 1 1 0 0 1 1 1 0

Se obține o secvență de 3 pachete de execuție A || B || C || D, E, F || G || H, în care s-a ținut cont de faptul că un șir de biți 1 urmat de un 0 conduce la o execuție paralelă, în timp ce un șir de biți începând cu un 0 indică o execuție secvențială a instrucțiunilor.

Cele 3 pachete de execuție sunt:

x A B x x C x D
E x x x x x x x
x F x G x x H x

rezultă 24 de instrucțiuni

Execuția instrucțiunilor se realizează în trei faze, fiecare în bandă de asamblare, conform următoarelor trei etape:

- Citire: are patru etaje (PG – generarea adresei program; PS – trimiterea adresei program; PW – accesarea programului din memoria intermediară; PR – primirea pachetului de program citit);
- Decodificare: are două etaje (DP – despachetarea instrucțiunilor, conversia pachetului citit în pachete de execuție și rutarea lor către unitățile funcționale; DC – decodificarea instrucțiunilor);
- Execuție: are maxim 10 etaje; 90% din instrucțiuni utilizează doar primele cinci etaje, doar operațiile de adunare/înmulțire în virgulă mobilă în dublă precizie utilizează și celelalte cinci etaje.

Procesorul Phillips Trimedia TM 1000

TM 1000 este un procesor de semnale multimedia, având nucleul de procesare CPU în arhitectură VLIW. Procesorul posedă cinci unități multifuncționale de execuție, având un total de 27 de funcții și poate executa 5 operații pe perioada de ceas; are de asemenea un fișier de 128 de registre cu 15 porturi de citire (deoarece o operație necesită doi operanzi și o informație pentru anticiparea ramificațiilor) și 5 porturi de scriere.

Se pot citi 64 de octeți simultan din memoria intermediară de instrucțiuni, împachetați într-o instrucțiune VLIW.

Execuția se realizează pe cinci etaje: citire instrucțiuni (FI), decompresie (DC), citire registre (RR), execuție (EX) și scriere rezultat (WB).

Schema microarhitecturală internă a procesorului este următoarea:

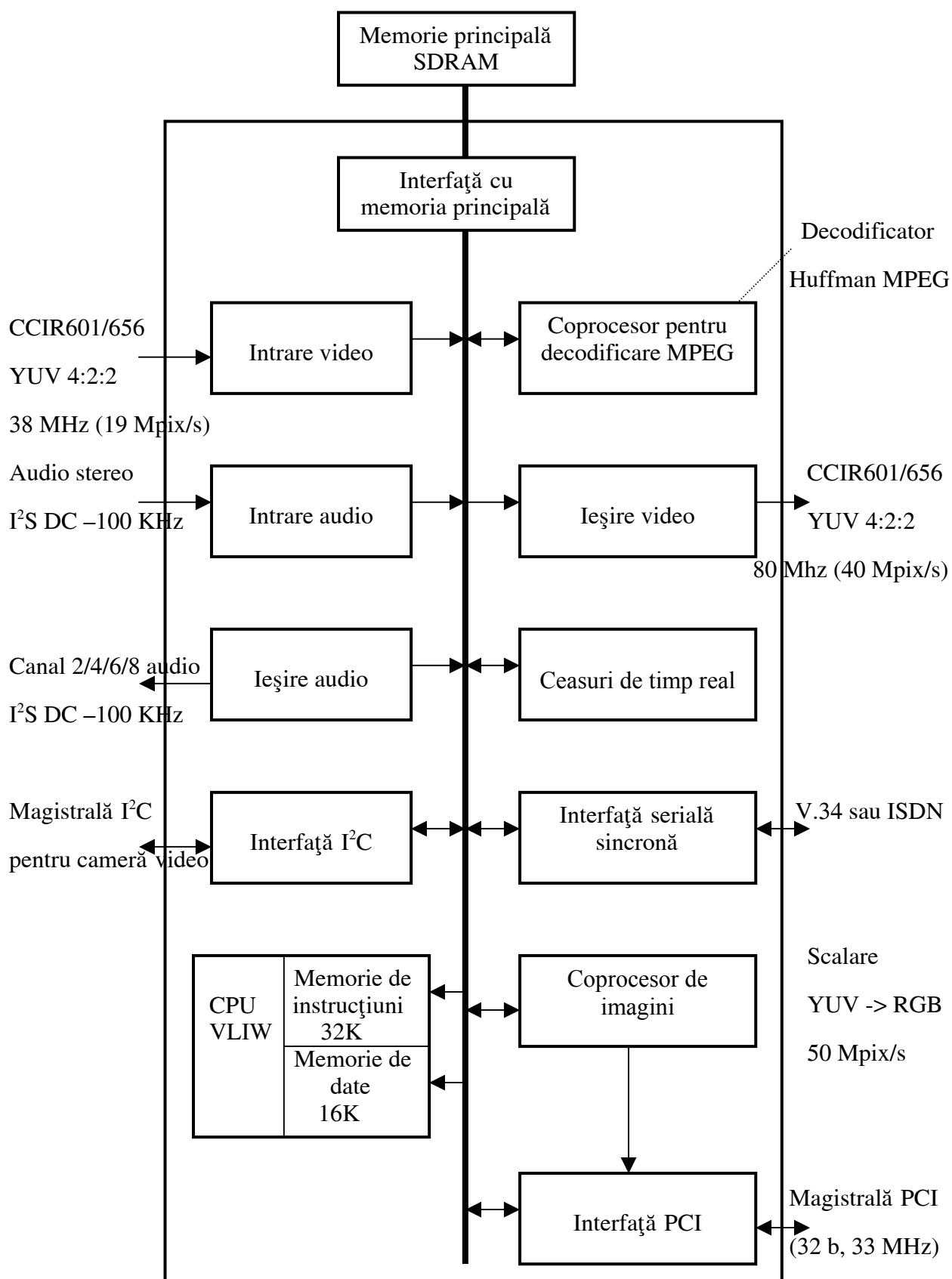


Fig. 6.3.

Implementări actuale de procesoare VLIW

Procesorul Transmeta Crusoe

În mod tradițional, procesoarele VLIW au fost proiectate cu scopul de a maximiza performanța prin execuția în paralel a cât mai multor instrucțiuni posibile. Firma Transmeta și-a propus, și a reușit, să realizeze un procesor cu performanțe moderate, comparativ cu procesoarele curente, dar cu un consum foarte mic de tensiune și fiind astfel potrivit pentru emularea procesoarelor din calculatoarele mobile, din familia Mobile Intel Pentium, sau a mașinilor virtuale Java.

Procesorul Crusoe are instrucțiuni pe 128 biți, numite molecule, fiecare moleculă codificând 4 operații, numite atomi. Formatul moleculei determină direct cum sunt rutate operațiile la cele 5 unități funcționale (2 unități de prelucrare întregi, o unitate de prelucrare în virgulă mobilă, o unitate de încărcare/memorare și o unitate de procesare a ramificațiilor).

Emularea codului pentru execuție se realizează cu ajutorul unui translator binar ce rulează la timpul de execuție. Astfel, numai acest "software" trebuie modificat atunci când se schimbă arhitectura procesorului. În urma execuției codului emulat se colectează informații care conduc la auto-optimizări ale translatorului.

Figura următoare prezintă schema bloc a procesorului Crusoe:

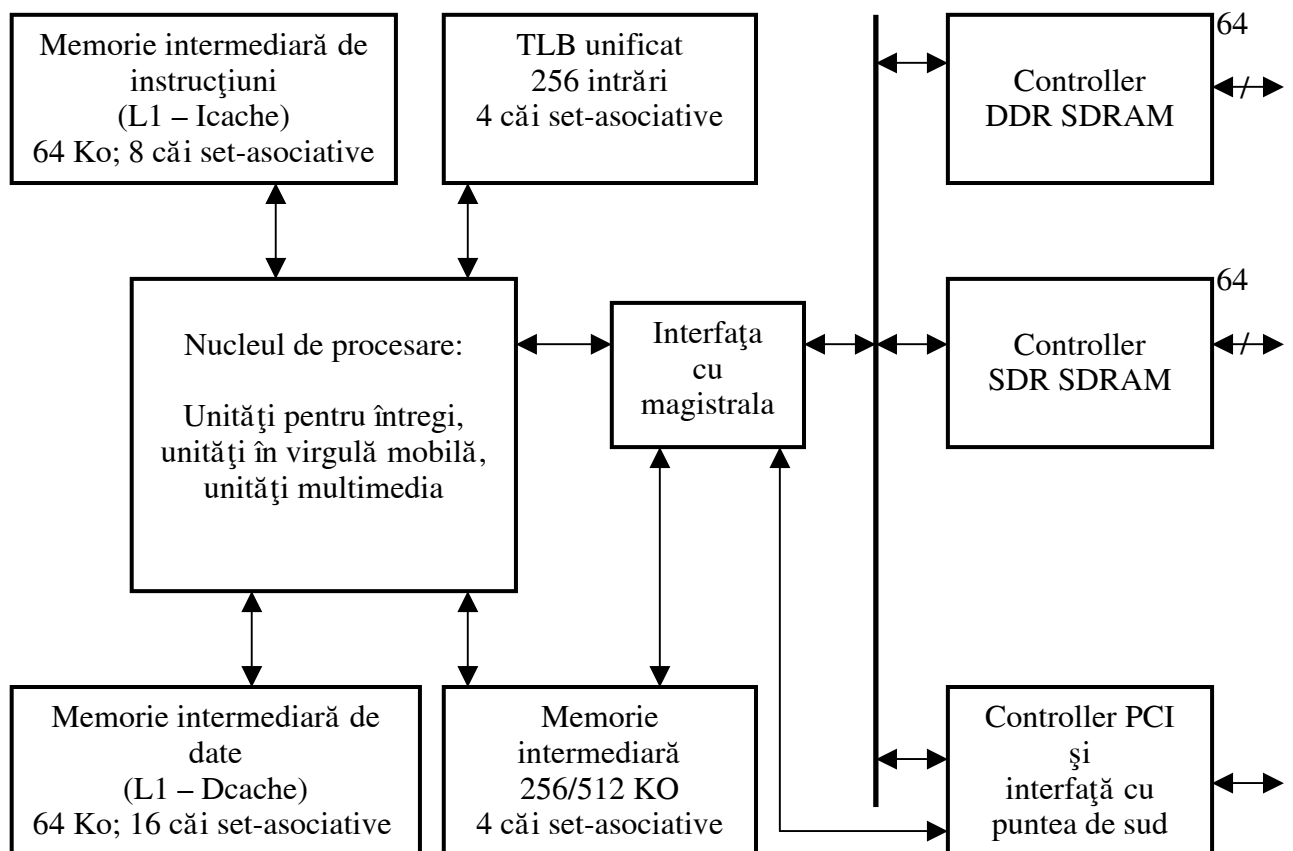


Fig. 6.4.

Procesorul Intel Itanium IA-64

Acest procesor este prima implementare pe 64 biți pentru a procesa instrucțiuni de tip VLIW, dezvoltat împreună de firmele Intel și Hewlett Packard (HP).

Procesorul conține 4 unități de prelucrare întregi, 4 unități multimedia, 2 unități de încărcare/memorare, 2 unități de prelucrare în virgulă mobilă cu dublă precizie și 2 unități de prelucrare în virgulă mobilă cu precizie simplă.

Acesta lucrează la o frecvență de 800 de MHz, este realizat într-o tehnologie de 0,18 microni și are banda de asamblare cu 10 etaje.

Arhitectura IA-64 utilizează un format fix de instrucțiuni împachetate pe 128 biți. Fiecare instrucțiune VLIW, numită multioperație, constă într-unul sau mai multe pachete de 128 biți, fiecare pachet având 3 operații și un câmp de control. IA-64 utilizează un tampon de decuplare și suport speculativ pentru a crește viteza de execuție.

Itanium posedă un set de 128 de registre generale și un alt set de 128 de registre în virgulă mobilă pe 82 biți.

Arhitectura IA-64 implementează suport pentru algoritmul de planificare pe bucle ("software pipelining"), dezvoltat de Bob Rau de la HP în cadrul unui proiect anterior (Cydra-5); pentru aceasta s-a implementat posibilitatea de rotație a registrelor astfel încât nu mai este necesară desfacerea buclei și redenumirea registrelor utilizate în iterații succesive.

IA-64 suportă execuție și control speculativ utilizând un mecanism "software/hardware" propriu. Ea suportă de asemenea planificarea statică și dinamică a ramificațiilor.

Setul de instrucțiuni conține și instrucțiuni SIMD potrivite pentru procesare multimedia.

Procesorul Itanium este probabil cea mai complexă arhitectură VLIW implementată comercial, ceea ce justifică preocupările multor alte firme de a include în proiectele lor modelul arhitectural VLIW.

Schema bloc a procesorului Intel Itanium este următoarea:

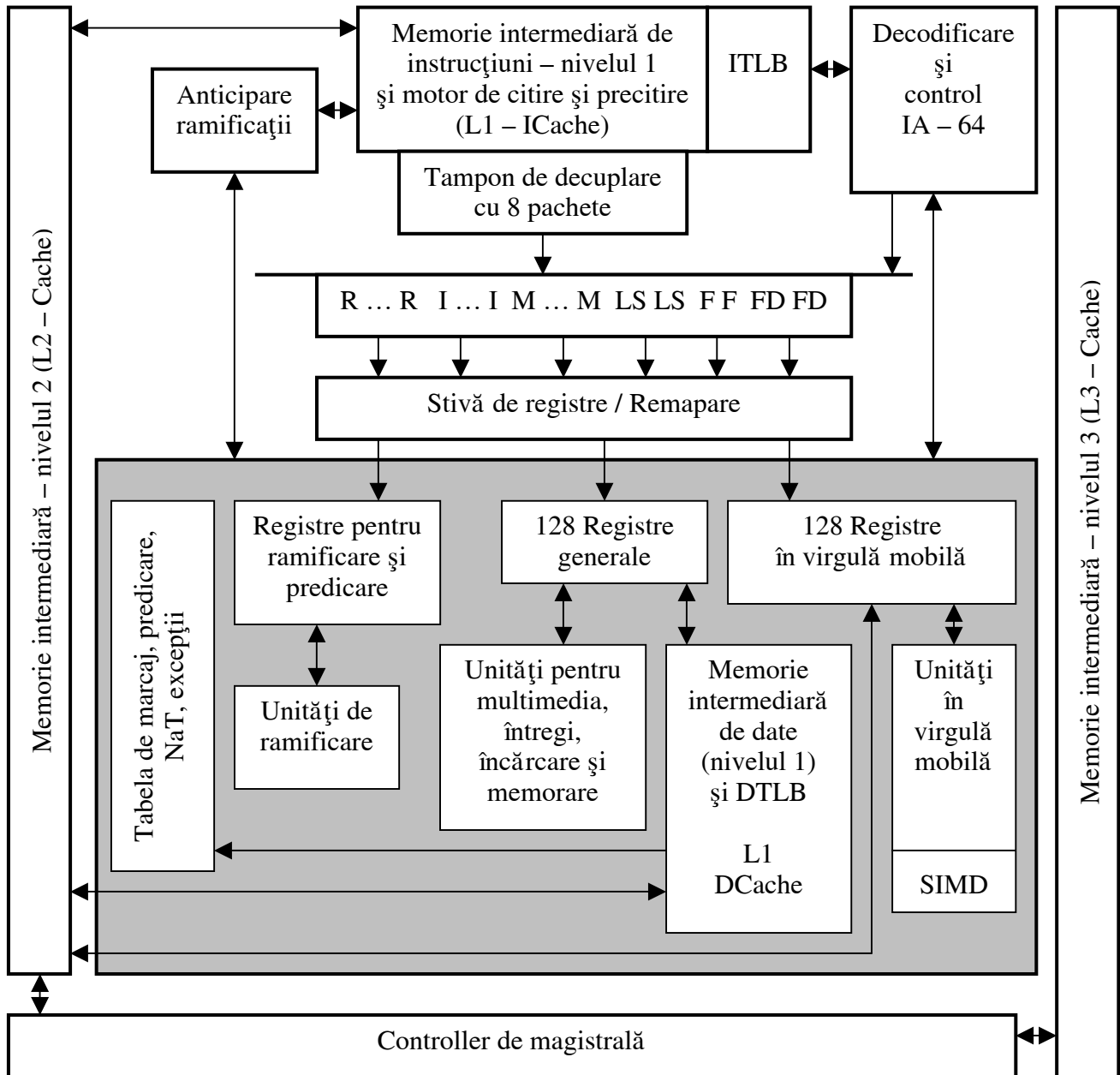
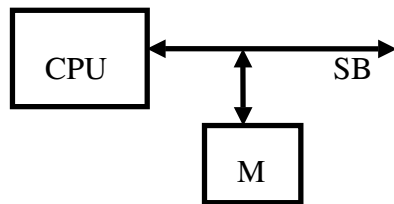


Fig. 6.5.

7. Creșterea performanței procesoarelor superscalare utilizând memorii intermediare

Utilizarea memoriilor intermediare (numite și memorii tampon sau “cache”) înăuntrul și în afara procesorului a condus la o creștere semnificativă a vitezei de execuție, atât în cazul procesoarelor CISC, cât și în cazul celor RISC, deci și pentru procesoarele VLIW.

Se consideră un sistem tipic constând într-un procesor (CPU – “Central Processing Unit”), o magistrală sistem (SB – “System Bus”) și o ierarhie de memorii M:



Se poate calcula rata de execuție a instrucțiunii (IR – “Instruction Rate”) în MIPS utilizând lățimea de bandă B (“Bus Bandwidth”). Lățimea de bandă în Mocteți/sec este definită prin volumul de lucru transferat la și de la memorie în unitatea de timp.

Se poate calcula B prin relația:

$$B = u_b \times \frac{n_b}{n_c \times t_c} \quad (1)$$

unde: u_b este utilizarea magistralei, n_b este numărul de octeți transferat pe un acces la memorie, n_c este numărul de cicli de ceas per transfer cu memoria și t_c este durata ciclului de ceas.

Se consideră că n_c este format din numărul de cicli de ceas pentru magistrala sistem și numărul de cicli de ceas pentru acces la memoria sistem. De asemenea, se consideră că un grad de utilizare a magistralei de minim 75% este întâlnit în majoritatea sistemelor de calcul.

Acum, IR poate fi calculată prin relația:

$$IR = \frac{B}{n_t} = \frac{B}{n_i + n_d} \quad (2)$$

unde: n_t este numărul total mediu de octeți per instrucțiune executată, n_i este numărul de octeți instrucțiune și n_d este numărul de octeți de date.

Există patru căi diferite de a utiliza o memorie intermediară într-un sistem de calcul:

Sistem cu memorie intermediară de instrucțiuni/date în afara procesorului

Când apare un “cache miss” (adică o eroare de pagina lipsă în memoria intermediară), atunci se vor adăuga un număr suplimentar de cicli de ceas (n_{ca}) necesari pentru a accesa memoria sistem. Deci, se poate calcula noul n_c^* prin relația:

$$n_c^* = n_c + n_{ca} \times \text{miss ratio extern} \quad (3)$$

Sistem cu memorie intermediară de instrucțiuni în interiorul procesorului

Când apare o eroare de pagină, se modifică atât n_c , cât și n_i .

Se pot calcula astfel:

$$n_c^* = n_{c1} + n_c \times \text{miss ratio} \quad (4)$$

unde: n_{c1} este n_c pentru a accesa memoria intermediară internă

$$n_i^* = \text{miss ratio intern} \times \frac{n_i \times \text{numărul de octeți de instrucțiune pe miss}}{n_b} \quad (5)$$

Sistem cu memorii intermediare separate pentru instrucțiuni și pentru date în interiorul procesorului

Când apare o eroare de pagină, se modifică atât n_c , cât și n_i , dar și n_d .

n_c și n_i se calculează ca mai sus, iar n_d astfel:

$$n_d^* = \text{miss ratio} \times \frac{n_d \times \text{numărul de octeți de date pe miss}}{n_b} \quad (6)$$

Sistem cu memorii intermediare interne separate pentru instrucțiuni și pentru date în conjuncție cu o memorie intermediară externă rapidă pentru instrucțiuni/date

În acest caz :

$$n_c^* = n_{c1i} + (n_{c2} + n_{ca} \times \text{miss ratio extern}) \times \text{miss ratio intern} \quad (7)$$

unde: n_{c1i} este n_c pentru memoria intermediară de instrucțiuni internă și

n_{c2} este n_c pentru accesul la memoria intermediară

n_i^* și n_d^* se calculează ca mai sus.

Evaluarea performanței unui procesor VLIW

Se vor compara performanțele obținute la un procesor superscalar performant (fie acesta procesorul PowerPC 620) și la un procesor VLIW.

Parametrii comuni cunoscuți sunt:

Frecvența ceasului: 500 MHz; gradul de utilizare a magistralelor: 75% (rezonabil de altfel).

Memoria intermediară de date și instrucțiuni externă = 128 Kocteți; miss ratio = 2.5%; linii de 128 octeți.

Memoria intermediară de date internă = 64 Kcuvinte; miss ratio = 2.5%; linii de 32 octeți.

Memoria intermediară de instrucțiuni internă = 64 Kcuvinte; miss ratio = 2%; linii de 64 octeți.

Magistrala externă de date este pe 128 biți iar magistrala internă de date este pe 64 biți.

$n_b = 128/8 = 16$ octeți; $n_{c_{li}} = 8$ cicli de ceas pentru un acces la memoria intermediară de instrucțiuni internă; $n_{c_2} = 10$ cicli de ceas pentru un acces la memoria intermediară externă; $n_{ca} = 5$ cicli de ceas pentru un acces la memoria sistem.

$$t_c = 1/500 \text{ MHz} = 0.2 \times 10^{-8} \text{ secunde};$$

Se aplică formulele (7), (1), (5), (6) și (2) din paragraful anterior.

Se obține:

$$n_c^* = 8 + (10 + 5 \times 0.025) \times 0.02 \cong 8.203 \text{ cicli de ceas.}$$

$$B = 0.75 \times [16/(8.203 \times 0.2 \times 10^{-8})] = 731.44 \text{ Mocteți/secundă.}$$

Pentru procesorul RISC, având magistrala internă de instrucțiuni pe 256 biți, rezultă:

$$n_i = 256/(8 \times 4 \text{ instrucțiuni}) = 8 \text{ octeți}; \quad n_d = 64/8 = 8 \text{ octeți}$$

$$n_i^* = 0.02 \times [(8 \times 32)/16] = 0.32 \text{ octeți} \quad n_d^* = 0.025 \times [(8 \times 32)/16] = 0.4 \text{ octeți}$$

Se obține o frecvență de execuție pentru procesorul RISC:

$$I_{R\text{PowerPC}} = 731.44/(0.32 + 0.4) \cong 1015.9 \text{ MIPS.}$$

Se va considera acum un procesor VLIW cu aceiași parametri de lucru ca și procesorul PowerPC, mai puțin magistrala internă de instrucțiuni, care este pe 432 biți.

Se obține:

$$n_i = 432/(8 \times 13 \text{ instrucțiuni}) = 4.154 \text{ octeți}; \quad n_d = 64/(8 \times 12 \text{ operații}) = 0.667 \text{ octeți.}$$

Rezultă:

$$n_i^* = 0.02 \times [(4.154 \times 32)/16] = 0.166 \text{ octeți}$$

$$n_d^* = 0.025 \times [(0.667 \times 32)/16] = 0.034 \text{ octeți}$$

Se obține o frecvență de execuție pentru procesorul VLIW:

$$I_{R\text{VLIW}} = 731.44/(0.166 + 0.034) \cong 3657.2 \text{ MIPS.}$$

Cum $I_{R\text{VLIW}} / I_{R\text{PowerPC}} \cong 3.6$, rezultă o creștere teoretică de 3.6 ori a ratei de execuție a instrucțiunilor pe un procesor VLIW față de execuția pe un procesor superscalar RISC.

Prin eliminarea dependențelor dintre instrucțiuni încă din faza de compilare a programului, înainte de execuția lui pe procesor, se obține pentru procesorul VLIW o creștere suplimentară a performanței de aproximativ 25%.

$$\text{Se obține } I_{R\text{VLIW}}^* = 1.25 \times I_{R\text{VLIW}} = 4571.5 \text{ MIPS.}$$

Rezultă o creștere de performanță de $I_{R\text{VLIW}}^* / I_{R\text{PowerPC}} \cong 4.5$ ori față de cea a procesorului RISC.

ÎNTREBĂRI

1. Să se prezinte tipurile de paralelism implementate în sistemele de calcul numeric.
2. Pe baza unei scheme bloc de funcționare, să se prezinte caracteristicile și modul de operare ale unui procesor VLIW.
3. Să se definească o arhitectură VLIW pentru procesorul DLX (MIPS) prezentat la curs.
4. Să se prezinte o modalitate “hardware” de creștere a performanței unui procesor VLIW.
5. Să se prezinte o modalitate “software” de creștere a performanței unui procesor VLIW.
6. Cum se poate păstra compatibilitatea între diferite generații de arhitecturi VLIW?
7. Evaluați performanța procesorului DLX (MIPS sau VLIW) cu și fără utilizarea de memorii intermediare interne separate pentru date și instrucțiuni și memorie intermediară externă pentru date/instrucțiuni; se va calcula frecvența de execuție a instrucțiunilor (în MIPS), utilizând unele din formulele prezentate în curs.