

Teme de Proiect An IV Calculatoare

la disciplina **Simularea și Optimizarea Arhitecturilor de Calcul**

A. CARACTERISTICI COMUNE

Scrieți în simulator pentru o arhitectură superscalară parametrizabilă. Scopul principal este acela de a determina diferiți parametri de performanță, configurații optime, pentru o arhitectură Harvard de memorie (cache-uri de instrucțiuni și date separate).

Principalii parametri ai arhitecturii sunt:

- **FR (rata de fetch)** - specifică numărul de instrucțiuni citite simultan din cache sau memorie într-un ciclu de tact; poate lua valori de 4, 8 sau 16 instrucțiuni.
- **IBS (instruction buffer size)** - dimensiunea buffer-ului de prefetch, măsurată în număr de instrucțiuni; plaja de valori: 4 (minim FR), 8, 16, 32; buffer-ul de prefetch este o coadă ce lucrează după principiul FIFO (first in first out). Vor fi citite FR instrucțiuni simultan de la adresa specificată de PC (program counter) și depuse în partea superioară a buffer-ului. În același ciclu de execuție, instrucțiuni din partea inferioară sunt expediate spre unitățile de decodificare și execuție. O intrare în buffer va conține câmpurile:
 - OPCODE - codul operației executată de instrucțiunea respectivă;
 - PC_crt - adresa (Program Counter-ul) instrucțiunii curente;
 - DATE / INSTR - adresa din / la care se citesc / se scriu date din sau în memorie, în cazul instrucțiunilor cu referire la memorie, respectiv adresa instrucțiunii țintă în cazul instrucțiunilor de salt.
- **IRmax (issue rate maxim)** - numărul maxim de instrucțiuni, lansate în execuție simultan într-un ciclu de execuție, din buffer-ul de prefetch. Poate lua valorile: 2, 4, 8, 16 (maxim FR) instrucțiuni. Dacă rata de fetch este mai mică decât numărul maxim de instrucțiuni executate concurent într-un ciclu, atunci performanța este limitată de procesul de *fetch instrucțiune*. Considerăm execuția instrucțiunilor “*in order*” - ordinea inițială a instrucțiunilor. O instrucțiune va fi executată abia după ce toate celelalte instrucțiuni de care ea depinde au fost executate.
- **Latența** - numărul de cicli necesari execuției instrucțiunilor aritmetice, de salt și cele cu referire la memorie (în cazul în care accesele pentru obținerea datei sunt cu hit în cache). Inițial are valoarea 1.
- **Cache-ul de instrucțiuni (IC) și Cache-ul de date (DC)** - sunt cache-uri mapate direct, organizate în blocuri de capacități parametrizabile [4, 8, 16 (maxim IBS) locații]. Încărcarea și evacuarea datelor în cache se face la nivel de bloc și nu la nivel de locație.
 - V – bit de validare (0 – nu e validă data; 1 – validă;). Inițial are valoarea 0. Este necesar numai pentru programe automodificabile la cache-urile de instrucțiuni. La prima înscriere în cache este setat pe 1.
 - D – bit Dirty. Este necesar la scrierea în cache-ul de date (vezi pct.4).

SIZE_IC, SIZE_DC - dimensiunea cache-urilor de instrucțiuni respectiv de date au plaja de valori de la 64 locații (128, 256, ...) până la 8192 locații.

$TAG = data \div SIZE_D(IC)$

$BLOC_OFFSET = [data \bmod SIZE_D(IC)] \div BLOC_SIZE (FR)$

BLOC_SIZE, FR - dimensiunea în locații a blocului din cache-ul de date respectiv instrucțiuni;
data – data emisă din program;

TAG - câmp de identificare al datei;

BLOC_OFFSET - offset-ul de bloc din cache;

TagE, TagC – Tag emis din program respectiv Tag citit din cache

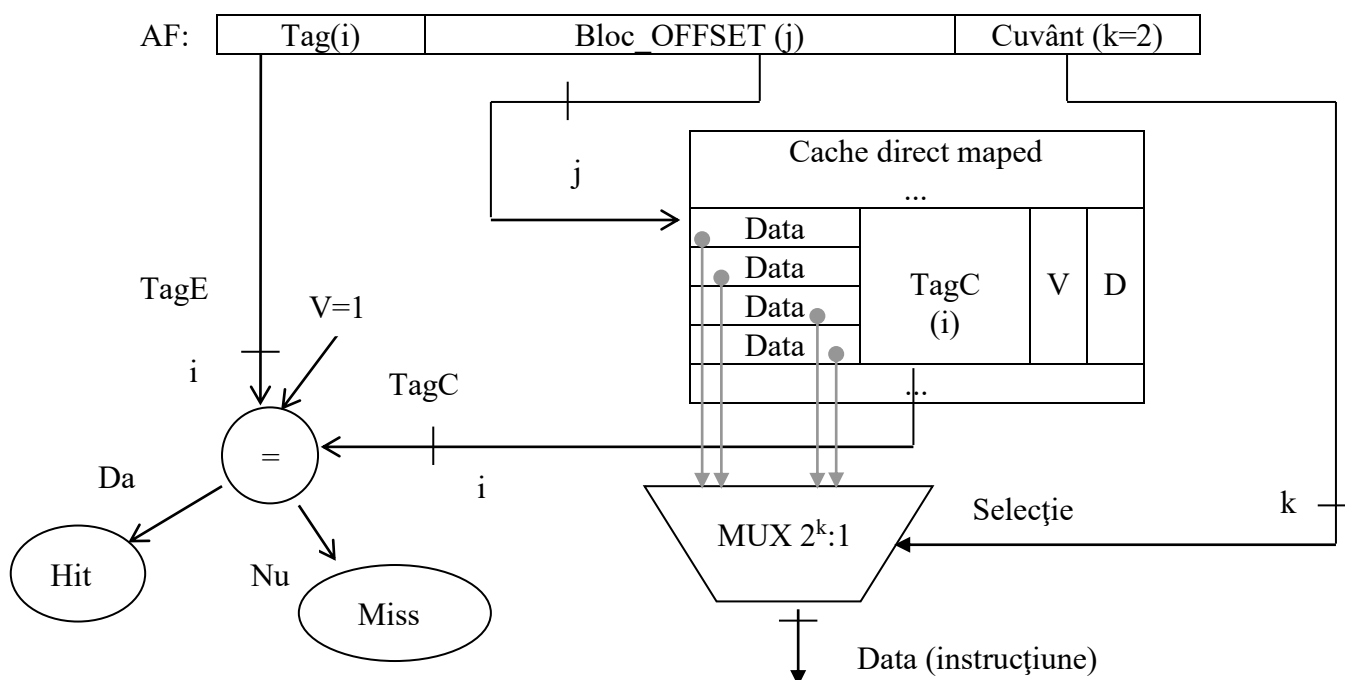


Figura 1. Structura cache-urilor de Instrucțiuni/Date

În cazul cache-urilor mapate direct, datele vor fi memorate în același loc de fiecare dată când sunt accesate. Din acest motiv vom ști la fiecare acces ce dată va fi evacuată din cache.

La o căutare în cache (IC sau DC) se ia tag-ul valorii căutate și se verifică dacă ea există la indexul sau la offset-ul de bloc respectiv. În caz afirmativ spunem că avem acces cu HIT, altfel MISS în cache și trebuie actualizat cache-ul.

- **Memoria principală** - (care se accesează numai la *miss* în cache) va avea o latență parametrizabilă de N_PEN (10, 15, 20) tați procesor. În cazul acceselor de date, sunt introduse penalizări numai pentru instrucțiunile LOAD (la STORE nu e nevoie din cauza procesorului de ieșire care pipelineizează scrierea în memorie, făcând-o transparentă pentru procesor). Este posibilă execuția a două instrucțiuni cu referire la memorie de genul: Load + Load sau Load + Store.
- Presupunem existența unui număr suficient de mare (maxim IRmax) de **seturi de regiștrii generali**: un set de regiștrii generali este necesar pentru execuția unei instrucțiuni de tip aritmetico-logic sau cu referire la memorie.
- Presupunem un *branch prediction* perfect, adică cunoașterea în permanență a adresei corecte a următoarei instrucțiuni ce se va executa.

Programul va simula fișiere *trace* (*.trc), rulate pe arhitectura HSA (Hartfield Superscalar Architecture). Este vorba de 8 benchmark-uri Stanford, care cuprind probleme clasice de sortare, problema turnurilor din Hanoi, problema damelor, generare de permutări și înmulțiri de matrici.

Rezultatele generate sunt sunt rata de procesare (număr de instrucțiuni raportat la număr de cicli de execuție), rate de miss în cache-uri (IC, DC), procentul din timpul total cât buffer-ul de prefetch este gol. Se vor determina parametri optimi și factorii de limitare în fiecare din cazuri.

CARACTERISTICI DISTINCTE

- 1.1. Generați rezultate urmate de grafice privind influența ratei de fetch (FR) asupra ratei de procesare IR(FR), asupra ratei de miss în cache-urile de date și instrucțiuni $R_{missDC}(FR)$ și $R_{missIC}(FR)$.
- 1.2. Studiați influența dimensiunii buffer-ului de prefetch asupra ratei de procesare IR(IFS) și asupra procentului din timpul total cât buffer-ul de prefetch este gol IBE (instruction buffer empty) - procesorul stă, nu are instrucțiuni de executat.
- 1.3. Studiați influența capacității cache-ului de instrucțiuni asupra ratei de procesare IR(SIZE_IC) și asupra ratei de miss la cache-ul de instrucțiuni $R_{missIC}(SIZE_IC)$.

- 2.1. Determinați influența numărului maxim de instrucțiuni ce pot fi trimise simultan în execuție asupra ratei de procesare IR(IRmax).
- 2.2. La acest punct nu se va mai considera număr nelimitat de seturi de regiștrii generali. Se va determina numărul optim de seturi de regiștrii (2, 3, 4, ...IRmax) în variantele cu cache de date uniport (o singură instrucțiune cu referire la memorie se poate executa) sau biport (două instrucțiuni cu referire la memorie se pot executa: L+L sau L+S).
- 2.3. Pentru valoarea optimă determinată la punctul 2.2. a numărului de seturi de regiștrii, studiați comparativ performanța (rata de procesare) pe două tipuri de cache de date (uniport sau biport).

- 3.1. Se va modela și simula un cache de instrucțiuni și date perfect din punct de vedere al ratei de hit (100% hit). Se va dovedi o metrică ideală cu care vom compara rezultatele obținute folosind cache-uri de instrucțiuni și date reale (rată de hit diferită de 100%). Grafic IR(Rata de hit).
- 3.2. Simulați execuția realistă a branch-urilor penalizând 0 respectiv 1 sau 2 cicli, la fiecare branch trimis spre execuție. La 0 cicli penalizare avem varianta optimă de predicție perfectă. Reprezentați grafic comparativ rata de procesare în cele trei situații.

4. Modelați și simulați procesul de scriere în cache prin prisma celor două strategii posibile: *write back* și *write through*. Write back e preferată în majoritatea implementărilor actuale deoarece scrierea are loc la viteza memoriei cache iar multiplele scrieri în bloc necesită doar o scriere în nivelul cel mai de jos al memoriei. Cu write through, informația e scrisă în ambele locuri: în blocul din cache și în blocul din memoria principală. Prin write back informația e scrisă doar în blocul din cache. Write back implică evacuare efectivă a blocului - cu penalitățile de rigoare - în memoria principală. Rezultă este necesar un bit **Dirty** asociat fiecărui bloc. Starea acestui bit indică dacă blocul e *Dirty* (modificat cât timp a stat în cache), sau *Clean* (nemodificat). Dacă bitul este curat, blocul nu e scris la miss, deoarece nivelul inferior – memoria principală - are copia fidelă a informației din cache. Dacă avem *citire din cache cu miss* și *Dirty* setat pe 1 atunci vom avea o penalizare egală în timp cu timpul necesar evacuării blocului - existent în cache dar nu cel care îl solicit - la care se adaugă timpul necesar încărcării din memorie în cache a blocului necesar în continuare. La **write through** nu există evacuare de bloc la cache-urile mapate direct, dar există penalități la fiecare scriere în memorie. Se vor genera graficele IR(BLOC_SIZE) și $R_{missDC}(BLOC_SIZE)$ în cele două ipostaze: scriere în cache prin write back și scriere în cache prin write through. Se va studia comparativ realismul, prin rata de procesare, introdus prin cele două tehnici de scriere față de situația când nu se folosește nici una din aceste tehnici.

5. Considerăm următoarea tehnică de îmbunătățire a performanței unui procesor. Aceasta constă în: presupunând că știm adresele de acces ale tuturor instrucțiunilor Load și Store din buffer-ul de prefetch, atunci în momentul execuției unui anumit Load, am putea considera executate toate

instrucțiunile Load din buffer-ul de prefetch care accesează același bloc, cu condiția ca locația respectivă să nu fie alterată de un Store intercalat. Acest lucru poate fi realizat într-un mecanism Out of Order unde există buffere pentru Load-uri și Store-uri sau într-un Trace Procesor. Folosind această tehnică, rata de procesare poate fi influențată de doi factori: dimensiunea cache-ului de date (la un cache de capacitate redusă - mai puține blocuri - probabilitatea de acces în același bloc fiind mai mare), și de capacitatea buffer-ului de prefetch - într-un buffer de prefetch mai mare fiind mai probabil să găsească mai multe instrucțiuni Load care accesează același bloc de date.

- 5.1. Studiați influența tehnicii prezentate, numită *multiple load*, coroborată cu dimensiunea cache-ului de date și tehnică de scriere write back, asupra ratei de procesare.
 - 5.2. Studiați influența tehnicii *multiple load*, coroborată cu dimensiunea buffer-ului de prefetch și tehnică de scriere write back, asupra ratei de procesare.
 - 5.3. Realizați un grafic comparativ privind rata de procesare cu și fără *multiple load*, cu tehnică de scriere write back.
6. Implementarea unei memorii Trace Cache în cadrul arhitecturii HSA. Determinarea câștigului de performanță față de un sistem care integrează doar un cache de instrucțiuni. Se va varia lungimea liniei din Trace Cache, gradul de asociativitate, dimensiunea buffer-ului de prefetch. În primă fază, se consideră predicția perfectă a salturilor, ulterior se va simula inclusiv folosirea unui predictor realist de branch-uri (adaptiv pe două niveluri, neuronal, etc).
 7. Implementarea unor algoritmi de evacuare a blocurilor din cache în vederea îmbunătățirii performanței arhitecturilor *single* și *multicore* (<http://www.jilp.org/jwac-1/JWAC-1%20Program.htm>). Dezvoltarea de simulatoare trace-driven pe framework-ul JWAC-1 (ISCA-2010 - <http://www.jilp.org/jwac-1/>) folosind trace-urile propuse. Metrici folosite (Speed-up, Rata de hit / miss).
 8. Implementarea unor algoritmi de evacuare a blocurilor din cache în vederea îmbunătățirii performanței arhitecturilor *single* și *multicore* (<http://crc2.ece.tamu.edu/>). Dezvoltarea de simulatoare trace-driven pe framework-ul **ChampSim** (ISCA-2017 - <https://github.com/ChampSim/ChampSim>) folosind trace-urile propuse. Metrici folosite (Speed-up, Rata de hit / miss).

B. CARACTERISTICI COMUNE

▼ Metodologia de simulare:

Simularea este efectuată pe benchmark-urile Stanford, 8 fișiere trace cu extensia (*.tra). Aceste fișiere sunt o prelucrare a programelor scrise în mnemonică de asamblare (*.ins) și a trace-urilor originale (*.trc), cu scopul de a evidenția toate salturile (inclusiv cele care nu se fac).

Întregul fișier trace (*.trc) este o înlanțuire de triplete **<TipInstr AdrCrt AdrDest>**, unde **TipInstr** poate lua una din valorile 'B' – branch, 'L' – load, 'S' – store; **AdrCrt** reprezintă valoarea registrului PC – adresa instrucțiunii curente, iar **AdrDest** reprezintă adresa de memorie a datei accesate – în cazul instrucțiunilor cu referire la memorie ('L' / 'S') sau adresa destinație a saltului – în cazul instrucțiunilor de salt și ramificație ('B').

Există totuși o deficiență a trace-urilor cu extensia *.trc: nu evidențiază salturile care nu se fac. Din acest motiv s-au generat noile trace-uri (fișierele *.tra). Acestea conțin doar branch-urile (atât cele care se fac cât și cele care nu se fac) și exclude instrucțiunile Load / Store. Forma acestor fișiere este următoarea:

BT	12	30
BS	32	98
BM	100	33
NT	36	37
BRA	2	100
MOV	PC, RA	(RETURN)

Reprezentarea salturilor se face tot sub forma unor triplete <**TipBr AdrCrt AdrDest**>, unde **TipBr** se prezintă sub forma unei codificări pe două caractere, primul dintre ele indică dacă saltul se face ('B' – saltul se face, 'N' – saltul nu se face), iar al doilea caracter indică tipul saltului: 'T' sau 'F' – salturi condiționate, 'S' – apeluri de tip Call, 'M' – apeluri de tip Return, 'R' – salturi necondiționate. Alegerea acestei codificări a fost inspirată de mnemonicile întâlnite în sursa în limbaj de asamblare (BT, BF – salt condiționat, BSR – instrucțiune de tip Call, BRA – salt necondiționat și MOV PC, RA – instrucțiune de tip Return).

▼ Automatul de predicție:

Este descris printr-un șir de caractere cu un format mai special, ce prezintă atât numărul de stări, tranzițiile între stări cât și predicția aferentă fiecărei stări. Pentru o mai bună înțelegere a funcționării automatului exemplificăm:

Automatul **BCBAADCD:12** - pe 2 biți

Tabelul care descrie funcționarea primului automat este următorul:

Stare Curentă	Stare următoare		Predicție
	Pentru intrare = 0	Pentru intrare = 1	
A	B	C	0
B	B	A	0
C	A	D	1
D	C	D	1

Prima parte a șirului până la caracterul ':' reprezintă tranzițiile pentru starea 'A' cu intrare 0, apoi cu intrare 1, apoi tranzițiile din 'B' pentru aceleași intrări, etc. Numărul din a doua parte este "văzut" în binar sub forma "0010" și reprezintă ieșirile asociate fiecărei stări în parte (stării 'A' îi este asociat cel mai puțin semnificativ bit, în acest caz bitul '0', următorul bit lui 'B', bitul '0', următorul bit lui 'C', bitul '1', etc).

CARACTERISTICI DISTINCTE

1. Să se proiecteze o structură de predicție a salturilor de tip BTB care poate fi **complet-asociativă** (vezi figura 2) sau **mapată direct**.

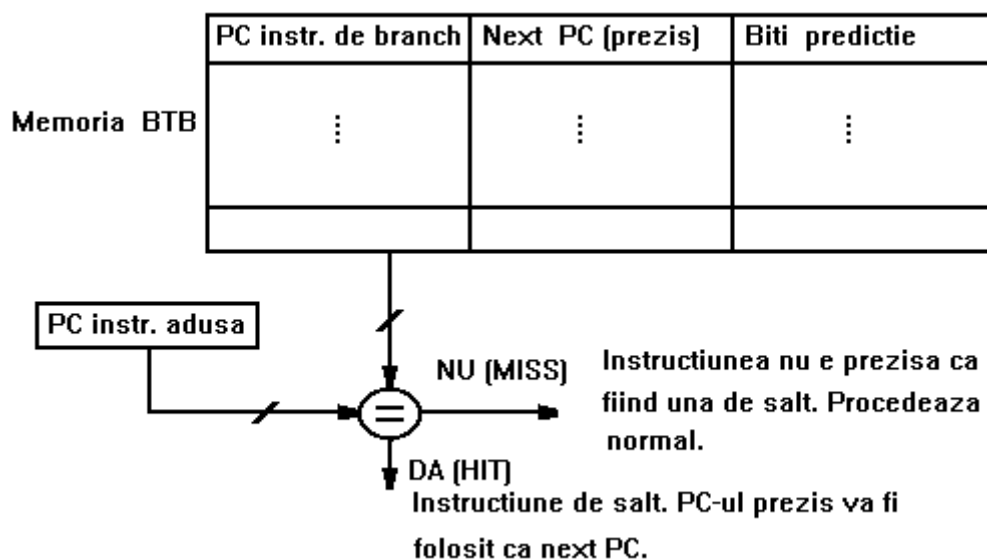


Figura 2. Structură de predicție BTB asociativă

În cazul proiectării unui BTB complet asociativ evacuarea se face prin metoda LRU (cel mai puțin recent folosit).

În principiu, simulatorul dezvoltat va funcționa astfel:

a) Inițializare simulator:

- Dimensiunea tabeli (număr de intrări)
- Inițializarea cu 0 a adreselor destinație
- Inițializare stare automat de predicție
- Tipul de arhitectură ales (mapat direct / complet asociativ)

b) Simularea propriu-zisă

Se stabilește de către utilizator benchmark-ul de tip trace care va fi utilizat. Din acest benchmark, se citesc secvențial instrucțiunile de ramificație și se compară predicția reală din trace cu cea propusă din tabelă. Aici pot să apară 3 cazuri distincte: predicție corectă, predicție incorectă, predicție incorectă datorată exclusiv adresei de salt incorecte din tabelă. Acest ultim caz se poate datora faptului că adresa de salt din tabelă a fost modificată de exemplu de către un alt salt, având astfel un **fenomen de interferență al salturilor** dar pot fi și alte cauze posibile. În final se actualizează locația folosită din tabela de predicții (starea automatului de predicție, câmpurile LRU și nextPC).

c) Generarea de rezultate

La finele simulării propriu-zise se generează rezultate semnificative precum numărul total de salturi executate, procentajul de predicții corecte, incorecte și respectiv afectate de interferențe ale salturilor. Se cere să se prezinte sub formă grafică funcțiile:

- $A_p = f(\text{tip_arhitectură})$
- $A_p = f(\text{dim_tabeli})$
- Să se repete rezultatele utilizând automatul de predicție pe 1 bit: **ABAB:2**. Realizați grafic: $A_p = f(\text{tip_automat})$

2. Să se implementeze o schemă de predicție corelată pe două niveluri de tip GAg (Global History Register, Global Prediction History Table) – vezi figura 3.

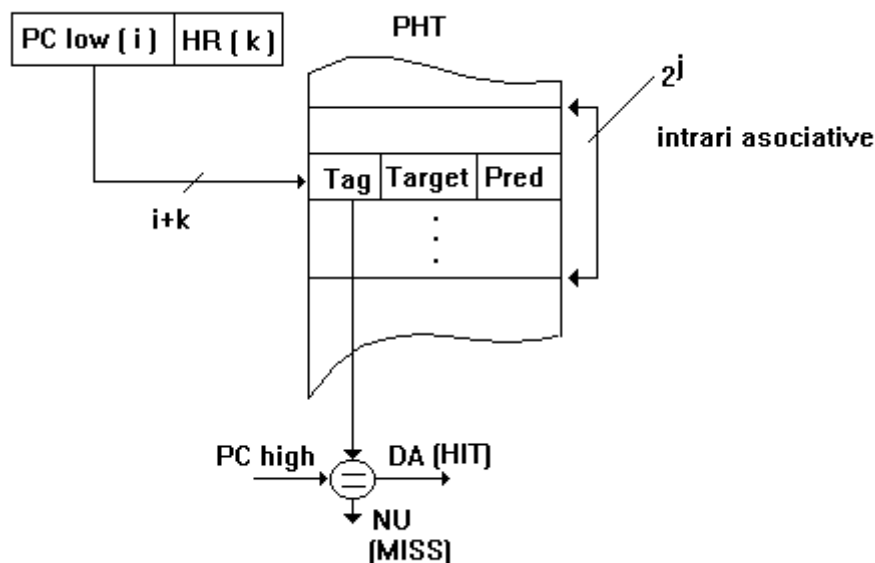


Figura 3. Structură de predicție de tip GAg

Tabela de predicții PHT (Prediction History Table) este adresată cu un index rezultat din concatenarea a **două informații ortogonale**: PClow (i biți), semnificând gradul de localizare al saltului, respectiv registrul de predicție (HR- History Register pe k biți), semnificând "contextul" în care se situează saltul în program. Ambele contribuții s-au făcut cu scopul eliminării interferențelor branch-urilor în tabela de predicție. Adresarea PHT exclusiv cu HR duce la serioase interferențe (mai multe salturi pot accesa aceleași automat de predicție din PHT), cu influențe evident defavorabile asupra performanțelor. PHT poate avea diferite grade de asociativitate (inițial considerăm tabela PHT complet asociativă). Un cuvânt din această tabelă are un format similar cu cel al cuvântului dintr-un BTB.

În cazul proiectării tabelii PHT complet asociativă evacuarea se face prin metoda LRU (cel mai puțin recent folosit).

Simulatorul dezvoltat va funcționa astfel:

a) Inițializare simulator:

- Inițializarea număr biți PClow (i = gradul de localizare al saltului)
- Inițializarea dimensiunii registrului de deplasare HRglobal (k = gradul de corelare al saltului)
- Inițializarea cu 0 a adreselor destinație
- Inițializare stare automat de predicție
- Tipul de arhitectură ales (mapat direct / complet asociativ)

b) Simularea propriu-zisă

Se stabilește de către utilizator benchmark-ul de tip trace care va fi utilizat. Din acest benchmark, se citesc secvențial instrucțiunile de ramificație și se compară predicția reală din trace cu cea propusă din tabelă. Aici pot să apară 3 cazuri distincte: predicție corectă, predicție incorectă, predicție incorectă datorată exclusiv adresei de salt incorecte din tabelă. Acest ultim caz se poate datora faptului că adresa de salt din tabelă a fost modificată de exemplu de către un alt salt, având astfel un **fenomen de interferență al salturilor** dar pot fi și alte cauze posibile. În final se actualizează corespunzător locația

folosită din tabela de predicții (starea automatului de predicție, câmpurile LRU și nextPC) precum și registrul de istorie al salturilor HRglobal.

c) Generarea de rezultate

La finele simulării propriu-zise se generează rezultate semnificative precum numărul total de salturi executate, procentajul de predicții corecte, incorecte și respectiv afectate de interferențe ale salturilor. Se cere să se prezinte sub formă grafică funcțiile:

- $Ap = f(HR_{global})$
- $Ap = f(i)$; $i = \text{size of PC_low}$
- $Ap = f(\text{tip_arhitectură})$

3. Să se implementeze un modul care detectează *salturile dificil de prezis într-un anumit context*. Realizați un predictor adaptiv pe 2 niveluri (GAg, PAg, gshare) în care structura de predicție va fi indexată cu informație de corelație extinsă (*path*). Determinați cantitativ și calitativ în ce măsură este redus numărul de salturi dificile. Determinați câștigul în acuratețe față de schema originală (GAg nemodificată) în condiții echivalente hardware (tabele de predicție de aceeași capacitate). Funcționarea simulatorului respectă subpunctele anterioare 2a÷2c.
4. Implementați un predictor contextual de tip PPM complet pentru predicția salturilor din benchmark-urile Stanford. Comparați acuratețea de predicție obținută cu cea a predictorului GAg de la punctul 2. Studiați fezabilitatea unui predictor simplificat bazat pe context compus dintr-un predictor Markov de ordin maxim și unul de ordin 0. Funcționarea simulatorului respectă subpunctele anterioare 2a÷2c.
5. Implementarea unui predictor de salturi condiționate de tip **perceptron simplu, fast path-based și idealized piecewise**. Analiză comparativă asupra costurilor. Modelare și implementare sub platformă academică SPEC 2000 și industrială standardizată CBP – campionatul mondial de predicția salturilor. Simularea prin programare distribuită a benchmark-urilor.
6. Implementarea unor predictoare de salturi condiționate de tip *state-of-the-art*: **O-GEHL, L-TAGE, GTL**. Analiză comparativă asupra costurilor. Modelare și implementare sub platformă academică SPEC 2000 și industrială standardizată CBP – campionatul mondial de predicția salturilor. Simularea prin programare distribuită a benchmark-urilor.
7. Dezvoltarea de simulatoare trace-driven pe framework-ul CBP (ISCA-2011) folosind trace-urile propuse. Metrice folosite (Acuratetea predicției) Simularea unor predictoare de salturi de tip „*state-of-the-art*” pentru platforma Campionatului Mondial al Predictoarelor de Salturi (http://www.jilp.org/jwac-2/cbp3_framework_instructions.html, <http://www.jilp.org/jwac-2/program/JWAC-2-program.htm>, http://www.jilp.org/jwac-2/program/07_seznec.tgz, http://www.jilp.org/jwac-2/program/cbp3_07_seznec.pdf, http://www.jilp.org/jwac-2/program/cbp3_07_seznec.pptx, etc).

C.

1. Realizați un scheduler minimal care să cuprindă următoarele tehnici de rearanjare a instrucțiunilor în scopul eliminării statice a dependențelor reale de date (RAW) și a dezambiguizării referințelor la memorie. Simularea se face pe trace-urile benchmark-urilor SPEC folosite de simulatorul PSATSIM sau pe fisierele (.ins, .trc) aferente benchmark-urilor Stanford.

a) Fuzionarea instrucțiunilor (*merging*)

Tehnica “*merging*” implică combinarea a două instrucțiuni într-una singură. Există trei categorii de astfel de instrucțiuni. Prima categorie, numită **MOV Merging** implică o pereche de instrucțiuni în care prima din ele este MOV. A doua categorie numită **Immediate Merging** se caracterizează prin faptul că ambele instrucțiuni au ca operanzi sursă valori imediate. A treia categorie se numește **MOV Reabsorption** și are ca a doua instrucțiune o instrucțiune MOV sau instrucțiunea ce se va infiltra va fi convertită la tipul primei instrucțiuni.

➤ MOV Merging

Când apare o dependență reală de date între o instrucțiune MOV și o instrucțiune care încearcă să promoveze în sus în structura de cod a programului, se verifică faptul dacă cele două instrucțiuni pot fi procesate în paralel. În caz afirmativ instrucțiunea ce se infiltrează își continuă drumul ascendent prin basic block. În continuare vom prezenta tipurile de situații ce pot să apară în cadrul tehnicii MOV merging.

– Combinarea cu o **instrucțiune de adunare**.

a) Secvența inițială:

```
MOV R6, R7
ADD R3, R6, R5    /* instrucțiunea care se infiltrează */
```

Secvența modificată:

```
MOV R6, R7; ADD R3, R7, R5
```

b) Secvența inițială:

```
MOV R6, #4
ADD R3, R6, #5    /* instrucțiunea care se infiltrează */
```

Secvența modificată:

```
MOV R6, #4; MOV R3, #9
```

Prin înlocuirea registrului R6 cu valoarea imediată respectivă instrucțiunea de adunare devine MOV.

– Combinarea cu o **instrucțiune store**

Secvența inițială:

```
MOV R3, #0
ST (R1, R2), R3    /* instrucțiunea care se infiltrează */
```

Secvența modificată:

```
MOV R3, #0; ST (R1, R2), R0
```

Registrul R3 este înlocuit cu R0 deoarece toate procesoarele RISC au registrul R0 cablat la 0.

– Combinarea cu o **instrucțiune relațională**

Secvența inițială:

```
MOV R4, #4
```

GT B1, R4, R3 /* instrucțiunea care se infiltrează */

Secvența modificată:

MOV R4, #4; LTE B1 R3, #4

Registrul R4 este înlocuit cu valoarea imediată memorată de el și instrucțiunea GT devine LTE pentru a permite operanzilor instrucțiunilor interschimbarea.

– Combinarea **instrucțiunilor gardate**

a) Secvența inițială:

EQ B3, R0, R0 /* B3 := true */

TB3 ADD R10, R11, R12 /* instrucțiunea care se infiltrează */

Secvența modificată:

EQ B3, R0, R0; ADD R10, R11, R12

Instrucțiuni de tipul EQ Bi, R0, R0 și NE Bi, R0, R0 sunt folosite pentru a înlocui instrucțiunile MOV Bi, #true sau MOV Bi, #false pe care arhitectura HSA nu le pune la dispoziție. Deoarece B3 este întotdeauna *true* garda TB3 aferentă instrucțiunii de adunare poate fi înlăturată. Dacă B3 ar fi evaluată întotdeauna la *false* atunci instrucțiunea care se infiltrează va fi înlocuită cu un NOP.

b) Secvența inițială:

MOV B1, B2

TB1 LD R4, (R0, R6) /* instrucțiunea care se infiltrează */

Secvența combinată:

MOV B1, B2; TB2 LD R4, (R0, R6)

Pentru eliminarea dependențelor de date, garda instrucțiunii LD devine acum B2.

c) Secvența inițială:

MOV B1, B2

BT B1, label /* instrucțiunea care se infiltrează */

Secvența modificată:

MOV B1, B2; BT B2, label

d) Secvența inițială:

EQ B1, R0, R0

BT B1, label /* instrucțiunea care se infiltrează */

Secvența modificată:

BRA label

Dacă registrul boolean care constituie gardă pentru instrucțiunea care se infiltrează are valoare constantă saltul fie va fi eliminat fie va fi transformat într-unul necondiționat (BRA).

➤ Immediate Merging

Această tehnică implică orice pereche de instrucțiuni care au valori imediate pe post de al doilea operand sursă.

Secvența inițială:

SUB R3, R6, #3

ADD R4, R3, #1 /* instrucțiunea care se infiltrează */

Secvența modificată:

SUB R3, R6, #3; ADD R4, R6, #-2

➤ MOV Reabsorption

Acest tip de combinare implică transformarea instrucțiunii MOV într-o instrucțiune de același tip cu prima.

Secvența inițială:

```
ADD R3, R4, R5
MOV R6, R3          /* instrucțiunea care se infiltrează */
```

Secvența combinată:

```
ADD R3, R4, R5; ADD R6, R4, R5
```

Ideea aflată la baza acestei metode este de a absorbi instrucțiunea MOV prin *renaming* aplicat registrului destinație al primei instrucțiuni reducând astfel expansiunea codului. În cazul în care prima instrucțiune este una cu referire la memorie (Load sau Store) duplicarea instrucțiunilor poate duce la reducerea performanței datorită numărului limitat de porturi de date ale cache-ului.

b) Analiza anti-alias a referințelor la memorie

Dependențele de date nu apar doar între regiștri, ci și între locații de memorie referite în instrucțiunile LD și ST. La fel ca celelalte dependențe ele provoacă deseori degradarea performanței procesoarelor. Pentru a face distincție între locațiile de memorie referite de cele două tipuri de instrucțiuni, folosim tehnica numită analiză anti-alias statică a memoriei (static memory disambiguation). Pentru a decide dacă o instrucțiune LD poate fi inserată în fața unei instrucțiuni ST, lucru ce se poate face în siguranță doar dacă cele două adrese diferă, adresele locațiilor de memorie sunt comparate și este returnată una din valorile:

- **Diferit:** Adresele sunt întotdeauna diferite.
- **Identic:** Adresele sunt întotdeauna identice.
- **Eșuează:** Adresele nu se pot distinge.

Dacă valoarea returnată este “Diferit” instrucțiunea LD poate fi inserată în fața instrucțiunii ST. De asemenea, dacă valoarea returnată este “Identic” instrucțiunea LD poate fi înlocuită cu o instrucțiune MOV ca în exemplul următor:

Secvența inițială:

```
ST (R0, R5), R6
LD R10, (R0, R5)
```

Secvența devine:

```
MOV R10, R6
ST (R0, R5), R6
```

Pe de altă parte, în cazul în care instrucțiunea LD este urmată de o instrucțiune ST, pentru ca aceasta din urmă să poată fi infiltrată în fața primeia **trebuie** ca cele două adrese să fie distincte întotdeauna.

Secvența de cod ce va fi analizată este următoarea:

MOV R6, R5	LES B1, R18, #25
DIV R5, R6, #120	BT B1, L10 (#0)
ASL R13, R5, #4	ADD R19, R19, #104
SUB R13, R13, R5	LES B1, R19, #2600
ADD R13, R13, #-60	BT B1, L11 (#0)
ST (R0, R17), R13	LD R18, 12(SP)
ADD R17, R17, #4	LD RA, 0(SP)
ADD R18, R18, #1	ADD SP, SP, #128

_Innerproduct:	MOV PC, RA (#0) SUB SP, SP, #128 MOV R10, R5 MOV R13, R8 ST (R0, R10), R0 MOV R8, #1 ASL R5, R13, #1 ASL R5, R5, #2 ADD R5, R5, R13 ADD R13, R5, R6 ASL R9, R9, #2 ADD R6, R13, #4 ADD R7, R7, #104	L16: ADD R13, R7, R9 LD R5, (R0, R6) LD R13, (R0, R13) MULT R5, R5, R13 LD R13, (R0, R10) ADD R13, R5, R13 ST (R0, R10), R13 ADD R6, R6, #4 ADD R8, R8, #1 LES B1, R8, #25 BT B1, L16 (#0) ADD SP, SP, #128 MOV PC, RA (#0)
----------------	---	--

D.

1. Generarea fișierelor *trace* pentru aplicațiile paralele din suita **PARSEC** și **SPLASH-2** folosindu-vă de simulatorul multicore **multi2sim**. Practic, în urma simulării unui benchmark Parsec cu simulatorul multi2sim se va genera pentru fiecare CORE în parte pe care a rulat benchmark-ul un fișier trace. Trace-urile rezultate vor trebui simulate pe arhitectura **SMPCache** (configurația multicore). Realizați o documentație aferentă evoluției trace-urilor considerând protocoalele de coerență **MSI / MESI** și exemplificați diferențele.
2. Implementarea unui **Automatic Design Space Explorer** aferent unei arhitecturi superscalare (Hatfield – vezi punctul A), pentru determinarea optimului din punct de vedere a ratei de procesare folosind algoritmi de căutare locali (în vecinătatea unei configurații microarhitecturale) de tip *Hill Climbing*, *Simulated Annealing*.
3. Implementarea unui **Automatic Design Space Explorer** aferent unei arhitecturi superscalare (vezi simulatorul PSATSim – arhitectura PowerPC), pentru analiza multiobiectiv (performanță, consum de energie) folosind tehnici de optimizare non-Pareto (suma ponderată a obiectivelor, algoritmul VEGA (*Vector Evaluated Genetic Algorithm*), ordonare lexicografică).
4. Implementarea unui **Automatic Design Space Explorer** aferent unei arhitecturi superscalare (vezi simulatorul PSATSim – arhitectura PowerPC), pentru analiza multiobiectiv (performanță, consum de energie) folosind tehnici de optimizare Pareto: NSGA-II (*Non-dominated Sorting Genetic Algorithm*) [27], SPEA2 (*Improving the strength Pareto evolutionary algorithm*) [28].
5. Implementarea unui **Automatic Design Space Explorer** aferent unei arhitecturi superscalare (vezi simulatorul PSATSim – arhitectura PowerPC) pentru identificarea arhitecturii optime din perspectivă multi-obiectiv (CPI, Energie) și care să implementeze mecanismul de Meta-Optimizare (nivelul de meta-optimizare poate rula mai mulți algoritmi de optimizare în paralel: NSGA-II și SPEA – abordare random și ponderată).



Bibliografie

- [1] **Chen I.K., Coffey J.T., Mudge T.** – *Analysis of Branch Prediction via Data Compression*, Proceedings of the 7th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS VII), Cambridge, MA, USA, October 1996, pp. 128-137.
- [2] **Florea A., Vințan L.** – *Simularea și optimizarea arhitecturilor de calcul în aplicații practice*, Editura MatrixRom, București, 2003, ISBN 973-685-605-4.

- [3] **Florea A.** – *Predicția dinamică a valorilor în microprocesoarele generației următoare*, Editura MatrixRom, București, 2005, ISBN 973-685-980-0.
- [4] **Jimenez D., Lin C.** – *Neural methods for dynamic branch prediction*. *ACM Transactions on Computer Systems*, 20(4):369–397, November 2002.
- [5] **Jimenez D.** – *Fast Path-Based Neural Branch Prediction*, in the Proceedings of the 35th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO-35), December 2003.
- [6] **Mudge T. Chen, I., Coffey, J.** – *Limits to Branch Prediction*, Technical Report, University of Michigan, January 1996.
- [7] **Nair, R.** – *Optimal 2-Bit Branch Predictors*, IEEE Transaction on Computers, No. 5, 1995.
- [8] **Vințan N. L., Florea A.** – *Microarhitecturi de procesare a informației*, Editura Tehnică, București, ISBN 973-31-1551-7, 2000 (312 pg.).
- [9] **Vintan L., Gellert A., Florea A., Oancea M., Egan C.** – *Understanding Prediction Limits through Unbiased Branches*, “Lecture Notes in Computer Science”, vol. 4186-0480, pp. 483-489, Springer-Verlag, ISSN 0302-9743, Berlin Heidelberg, 2006.
- [10] **Florea A., Gellert A., Radu C., Calborean H., Crapciu A.**, [An Interactive Graphical Trace-Driven Simulator for Teaching Branch Prediction in Computer Architecture](#), The 6th EUROSIM Congress on Modelling and Simulation, (EUROSIM 2007), ISBN 978-3-901608-32-2, September 9-13, Ljubljana, Slovenia 2007, (special session: *Education in Simulation / Simulation in Education I*).
- [11] **Gellert A., Florea A., Vintan M., Egan C., Vintan L.**, [Unbiased Branches: An Open Problem](#), Proceedings of the Twelfth Asia-Pacific Computer Systems Architecture Conference (ACSAC 2007), Korea, Seoul, August 2007.
- [12] **Vințan L., Egan, C.** – *Extending Correlation in Branch Prediction Schemes*, International Euromicro’99 Conference, Milano, Italy, September 1999.
- [13] **Vințan N. L., Florea A.** – *Branch Prediction: a Criticism and a Novel Scheme* - SINTES 10, Craiova, 2000.
- [14] **Vințan L., Florea A.** - *Investigating New Branch Prediction through Quantitative Approach – Beyond 2000: Engineering Research Strategies*, Sibiu, 1999.
- [15] <http://webspaces.ulbsibiu.ro/adrian.florea/html/>
- [16] <http://webspaces.ulbsibiu.ro/lucian.vintan/html/>
- [17] Sez nec A. – *The Idealistic GTL Predictor*, Journal of Instruction Level Parallelism, no. 9, May, 2007
- [18] Sez nec A. – *The L-TAGE Branch Predictor*, Journal of Instruction Level Parallelism, no. 9, May, 2007
- [19] Sez nec A. – *Analysis of the OGEHL predictor*, Proceedings of the 32th International Symposium on Computer Architecture (IEEE-ACM), Madison, June 2005.
- [20] <http://arco.unex.es/smpcache/#MemoryTraces>
- [21] <http://parsec.cs.princeton.edu/download.htm>
- [22] http://www.multi2sim.org/wiki/index.php5/Main_Page
- [23] <http://www.jilp.org/jwac-1/JWAC-1%20Program.htm>
- [24] <http://www.jilp.org/jwac-2/program/JWAC-2-program.htm>
- [25] <http://www.crhc.illinois.edu/TechReports/2007reports/magellan.pdf>
- [26] <http://www.autonlab.org/tutorials/hillclimb02.pdf>
- [27] K. Deb, A. Pratap, S. Agarwal, and T. Meyarivan, “A fast and elitist multiobjective genetic algorithm: NSGA-II,” *Evolutionary Computation, IEEE Transactions on*, vol. 6, no. 2, pp. 182-197, 2002.
- [28] E. Zitzler, M. Laumanns, L. Thiele, and others, “SPEA2: Improving the strength Pareto evolutionary algorithm,” in *Eurogen*, pp. 95–100, 2001.
- [29] Chiș Radu, *Developing Effective Multi-Objective Optimization Methods for Complex Computing Systems*, PhD Thesis, Lucian Blaga University of Sibiu, September 2017.