

Probleme propuse pentru examenul de la
disciplina Structura și Organizarea
Calculatoarelor

== MOODLE TYPE ==

Prof. dr. ing. Vasile Gheorghiuță GĂITAN
Ș. I. dr. ing. Ionel ZAGAN

Cap. 1 Introducere

1. Fie două implementări diferite M1 și M2 ale aceluiași set de instrucțiuni. Există trei clase de instrucțiuni (A, B și C) în setul de instrucțiuni. M1 are un ceas de 80 de MHz și M2 are un ceas de 100 MHz. Numărul mediu de cicli pentru fiecare clasă de instrucțiuni și frecvența lor (pentru un program tipic) sunt după cum urmează:

Instruction Class	Machine M1 – Cycles/Instruction Class	Machine M2 – Cycles/Instruction Class	Frequency
A	1	2	60%
B	2	3	30%
C	4	4	10%

- Calculează CPI mediu pentru fiecare mașină M1 și M2.
 - Calculează MIPS mediu pentru fiecare mașină M1 și M2.
 - Care mașină are raportul MIPS mai mic? Pentru care clasă individuală de instrucțiuni trebuie să schimbi CPI-ul, și cât de mult, pentru ca această mașină să aibă aceleași performanțe sau mai bune față de mașina cu MIPS-ul mai mare (poți schimba doar CPI-ul pentru una dintre clasele de instrucțiuni de pe mașina mai lentă)?
2. Presupuneți că aveți un procesor care execută un program constând din 50% înmulțiri în virgulă flotantă, 20% împărțiri în virgulă flotantă, iar restul de 30% sunt alte instrucțiuni (legea lui Amdahl).
- Managementul dorește ca acest procesor să ruleze de 4 ori mai repede. Se poate face ca împărțirea să meargă de 3 ori mai repede și înmulțirea de 8 ori mai repede. Puteți îndeplini cerința managementului doar cu o singură îmbunătățire, și care este aceasta?
 - Dogbert a preluat compania înlăturând toți managerii precedenți. Dacă realizați ambele îmbunătățiri, la înmulțire și la împărțire, care este viteza procesorului astfel îmbunătățit relativ la procesorul original?
3. Presupunem că putem îmbunătăți performanțele în virgulă flotantă ale unui procesor cu un factor având valoarea 15 (aceleași instrucțiuni în virgulă flotantă sunt de 15 ori mai rapide pe noul procesor). Ce procentaj dintre instrucțiuni trebuie să fie în virgulă mobilă pentru a atinge o îmbunătățire a vitezei de cel puțin 4 ori?
4. Așa cum este definit MIPS-ul se poate defini MFLOPS-ul care înseamnă milioane de instrucțiuni în virgulă flotantă per secundă. Dacă procesorul A are un MIPS mai mare decât procesorul B, atunci în mod obligatoriu are și un MFLOPS mai mare decât B?
5. Se consideră benchmark-ul SPEC. Numiți doi factori care influențează performanțele rezultate pentru orice arhitectură particulară.
6. Cum a afectat calculatorul dezvoltarea tranzistorilor? Ce a înlocuit tranzistorul?
7. Presupunem că o echipă de proiectare realizează îmbunătățirea unui procesor prin adăugarea unui modul hardware MMX (multimedia extension instruction). Atunci când un calcul se realizează în modul MMX pe hardware-ul MMX acesta este de 10 ori mai rapid decât în modul normal de execuție. Se va numi timpul petrecut utilizând modul MMX, procentul de îmbunătățire media.
- Ce procent mediu de îmbunătățire este necesar pentru a obține o îmbunătățire generală a vitezei de 2?
 - Ce procent din timpul de execuție este petrecut în modul MMX dacă se realizează o mărire de două ori a vitezei? (Sfat: trebuie calculat noul timp general).
 - Ce procent de îmbunătățire mediu este necesar pentru a obține o jumătate din maximum de mărire a vitezei (one-half) care se poate obține utilizând modul MMX?

8.

- a. Dacă procesorul A are ceasul mai mare ca procesorul B, și procesorul A are și MIPS-ul mai mare decât procesorul B, explicați când procesorul A va fi mai rapid decât procesorul B. Presupuneți că există două implementări ale aceleiași arhitecturi a setului de instrucțiuni.
 - b. Procesorul A are ciclul de ceas de 20ns și un CPI efectiv de 1,5 pentru unele programe, și procesorul B are ciclul de ceas de 15ns și un CPI efectiv de 1 pentru unele programe. Care procesor este mai rapid pentru acest program și cu cât?
9. Presupunem că un segment de program constă dintr-o parte pur secvențială care ia pentru execuție 25 de cicli, și o buclă iterativă care ia 100 de cicli per iterație. Presupuneți că iterațiile din buclă sunt independente, și nu pot fi paralelizate. Dacă bucla se execută de 100 de ori, care este creșterea maximă de viteză posibilă dacă se utilizează un număr infinit de procesoare (în comparație cu un singur procesor)
10. Procesorul A are un CPI general de 1,3 și poate rula la o frecvență de ceas de 600 MHz. Procesorul B are un CPI de 2,5 și poate rula la o frecvență de ceas de 750MHz. Avem un anumit program pe care dorim să-l executăm. În urma compilării pe procesorul A rezultă 100.000 de instrucțiuni. Câte instrucțiuni trebuie să aibă programul compilat pe procesorul B, pentru ca cele două procesoare să aibă, pentru acest program, exact același timp de execuție?
11. Imaginați-vă că puteți realiza o cursă bazată pe benchmark pentru a compara două calculatoare pe care credeți că vreți să le cumpărați. Precizați o listă cu 5 programe de benchmark sau un scenariu de utilizare pe care vreți să-l folosiți pentru a crea o suită proprie, personalizată de benchmark. Pentru fiecare program l-ați selectat, justificați-l. Pentru suita de benchmark-uri ca întreg, discutați o metodă pentru calculul mediei ponderate a diferitelor programe în execuție.
12. O echipă de proiectanți trebuie să aleagă pentru un procesor simplu cu o singură instrucțiune pe eeciclu, între o implementare pipeline sau fără pipeline. În tabel sunt prezentați câțiva parametri de proiectare pentru cele două opțiuni.

Parameter	Pipelined Version	Non-Pipelined Version
Clock Rate	500MHz	350 MHz
CPI for ALU instructions	1	1
CPI for Control instructions	2	1
CPI for Memory instructions	2.7	1

- a. Pentru un program cu 20% instrucțiuni ALU instrucțiuni, 10% instrucțiuni de control și 70% ezinstrucțiuni de lucru cu memoria, care va fi proiectul mai rapid? Furnizați valoarea medie a CPI pentru fiecare caz.
 - b. Pentru un program cu 80% instrucțiuni ALU instrucțiuni, 10% instrucțiuni de control și 10% instrucțiuni de lucru cu memoria, care va fi proiectul mai rapid? Furnizați valoarea medie a CPI pentru fiecare caz.
13. Un proiectant dorește să îmbunătățească performanțele generale a unui procesor dat cu respectarea unei ținte de tip suită benchmark și se ia în considerare o îmbunătățire X care se aplică la 50% instrucțiunilor originale executate dinamic, și le mărește viteza de execuție fiecareia dintre ele cu un factor 3. Managerul de proiectare are unele îngrijorări despre complexitate și costul efectiv pentru X, și sugerează ca proiectantul să ia în considerare o îmbunătățire Y alternativă. Îmbunătățirea Y, dacă se aplică numai la câteva dintre instrucțiunile executate dinamic le va face doar cu 75% mai rapide. Determinați ce procent dintre toate instrucțiunile executate dinamic trebuie îmbunătățite utilizând optimizarea Y pentru a obține aceeași creștere de viteză propusă de îmbunătățirea X.

14. Descrieți pașii care transformă un program scris într-un limbaj de nivel înalt, precum C, într-o reprezentare care este direct executată de procesor.
15. Presupuneți că un afișor color utilizează 8 biți pentru fiecare dintre culorile de bază (roșu, verde și albastru) per pixel și un cadru având mărimea de 1280 x 1024.
- Care este mărimea minimă în octeți al buffer-ului pentru cadru care memorează cadrul?
 - Când de mult durează trimiterea unui cadru pe o rețea cu viteza de 100 Mbit/s?
16. Se consideră trei procesoare P1, P2 și P3 care execută același set de instrucțiuni. P1 are frecvența de ceas de 3 GHz și un CPI de 1,5, P2 are 2,5 GHz și un CPI de 1 și P3 are 4GHz și un CPI de 2,2.
- Care procesor are performanțele cele mai mari exprimate în MIPS?
 - Dacă fiecare procesor execută un program în 10 secunde, calculați numărul de cicli și numărul de instrucțiuni.
 - S-a încercat reducerea timpului cu 30% dar aceasta a condus la o creștere cu 20% a CPI-ului. Care frecvență de ceas trebuie pentru a avea această reducere?
17. Se consideră două implementări diferite ale aceleiași arhitecturi a setului de instrucțiuni. Instrucțiunile pot fi divizate în 4 clase în funcție de CPI (clasele A, B, C și D). P1 are frecvența de ceas de 2,5 GHz și CPI-urile 1,2,3 și 3, iar P2 are un ceas de 3 GHz și CPI-urile 2,2,2 și 2. Dându-se un program cu conținutul dinamic de instrucțiuni de 1.0E6 instrucțiuni divizate în clase după cum urmează: clasa A 10%, clasa B 20%, clasa C 50%, și clasa D 20%, care implementare este mai rapidă?
- Care este CPI-ul global pentru fiecare implementare?
 - Găsiți ciclurile de ceas necesare în ambele cazuri.
18. Compilatoarele pot avea un impact profund asupra performanțelor unei aplicații. Presupunem că pentru un program, compilatorul A generează un conținut dinamic pentru instrucțiuni de 1.0E9 și are un timp de execuție de 1,1s, în timp ce compilatorul B generează un conținut dinamic pentru instrucțiuni de 1.2E9 și are un timp de execuție de 1,5s.
- Găsiți CPI-ul mediu pentru fiecare program dat fiind faptul că durata unui ciclu de ceas al procesorului este de 1ns.
 - Presupuneți că programele compilate se execută pe două procesoare. Dacă timpul de execuție pe cele două procesoare este același, cât de rapid este ceasul procesorului care rulează codul compilat de compilatorul A în raport cu ceasul procesorului care rulează codul compilat de compilatorul B?
 - Se dezvoltă un nou compilator care utilizează doar 6.0E8 instrucțiuni și care are un CPI mediu de 1,1. Care este creșterea de viteză utilizând acest compilator versus compilatoarele A și B de pe procesorul original?
19. Procesorul Pentium IV Prescott, realizat în 2004, avea un ceas de 3,6 GHz, la o tensiune de lucru de 1,25V. Presupuneți că, pe medie, acesta consuma 10W de energie statică și 90W de energie dinamică. Nucleul Ivy Bridge, realizat în 2012, are un ceas de 3,4 GHz la o tensiune de 0,9V. Presupuneți că în medie consumă 30W de putere statică și 40W de putere dinamică.
- Pentru fiecare procesor găsiți încărcare capacitivă medie.
 - Găsiți procentul din puterea totală disipată ce corespunde puterii statice și raportul dintre puterea statică și cea dinamică pentru fiecare tehnologie.
 - Dacă puterea totală disipată este redusă cu 10%, cât de mult trebuie redusă tensiunea pentru a menține același curent de scurgere? Notă: puterea este definită ca produsul între tensiune și curent.
20. Presupuneți că un procesor are pentru instrucțiunile aritmetice, load/store și de salt un CPI de 1,2, respectiv 5. Presupuneți de asemenea că un program executat pe un singur procesor necesită execuția a 2,56E9 instrucțiuni aritmetice, 1,28E9 instrucțiuni load/store și 256 de milioane de instrucțiuni de salt. Presupuneți că fiecare procesor are un ceas de 2 GHz. Presupuneți că, dacă programul este paralelizat pentru a se executa pe nuclee multiple, numărul de instrucțiuni aritmetice și load/store per procesor este divizat cu 0,7 x p (unde p este numărul de procesoare), dar numărul de salturi per procesor rămâne același.

- a. Calculați timpul total de execuție pentru acest program pe 1, 2, 4 și 8 procesoare, și prezentați creșterea relativă de viteză a celor 2, 4 și 8 procesoare la un singur procesor.
 - b. Dacă CPI-ul instrucțiunilor aritmetice se dublează care este impactul execuției programului pe 1, 2, 4 și 8 procesoare?
 - c. La ce valoare trebuie redus CPI-ul instrucțiunilor load/store pentru un singur procesor pentru ca acesta să aibă aceleași performanțe cu patru procesoare ce utilizează valorile originale ale CPI?
21. Presupunem un wafer de 15 cm diametru care are un cost de 12, conține 84 de die, și are 0,020 defecte/cm². Presupuneți ca un wafer de 20 de cm diametru are un cost de 15, conține 100 de die, și are 0,031 defecte/ cm².
 - a. Calculați yield pentru ambele wafer-e.
 - b. Calculați costul per die pentru ambele wafer-e.
 - c. Dacă numărul de die per wafer este crescut cu 10% și defectele pe unitatea de arie cresc cu 15%, calculația aria die și yield.
 - d. Presupuneți că un proces de fabricație îmbunătățește yield de la 0,92 la 0,95. Calculați defectele pe unitatea de arie pentru fiecare versiune a tehnologiilor, fiind dată aria die ca fiind 200 mm².
22. Rezultatele pentru benchmark-ul SPEC CPU2006 bzip2 ce rulează pe un AMD Barcelona are un contor de instrucțiuni de 2,389E12, un timp de execuție de 750s și un timp de referință de 9650s.
 - a. Calculați CPI-ul dacă perioada ceasului este de 0,333ns.
 - b. Calculați SPECratio.
 - c. Calculați creșterea timpului CPU dacă numărul de instrucțiuni ale banchmark-ului crește cu 10% fără a afecta CPI-ul.
 - d. Calculați creșterea timpului CPU dacă numărul de instrucțiuni ale banchmark-ului crește cu 10% iar CPI-ul crește cu 5%.
 - e. Calculați modificările SPECratio pentru această schimbare.
 - f. Presupuneți că dezvoltați o nouă versiune de procesor AMD Barcelona cu un ceas de 4 GHz. Ați adăugat câteva instrucțiuni adiționale setului de instrucțiuni astfel încât numărul de instrucțiuni a fost redus cu 15%. Timpul de execuție s-a redus la 700s și noul SPECratio este 13,7. Găsiți noul CPI.
 - g. Această valoare este mai mare decât aceea obținută la 1.11.1 deși ceasul a crescut de la 3 la 4 GHz. Determinați situația în care creșterea CPI este similară creșterii frecvenței ceasului. Dacă nu sunt similare explicați de ce.
 - h. Cu cât s-a redus timpul CPU?
 - i. Pentru un al doilea benchmark, libquantum, presupuneți că timpul de execuție este 960ns, CPI de 1,61, și ceasul este de 3 GHz. Dacă timpul de execuție este redus cu 10% fără a afecta CPI-ul și cu un ceas de 4 GHz, determinați numărul de instrucțiuni.
 - j. Determinați ceasul necesar pentru o nouă reducere cu 10% a timpului CPU menținând numărul de instrucțiuni și CPI-ul.
 - k. Determinați ceasul dacă CPI se reduce cu 15% și timpul CPU cu 20%, în timp ce numărul de instrucțiuni rămâne neschimbat.
23. Secțiunea 1.10 citează ca și capcană utilizarea unui subset al ecuațiilor pentru performanță ca metrică pentru performanță. Pentru a ilustra aceasta să considerăm următoarele două procesoare. P1 are un ceas de 4GHz, un CPI mediu de 0,9 și necesită execuția a 5,0E9 instrucțiuni. P2 are un ceas de 3 GHz, un CPI mediu de 0,75, și necesită execuția a 1,0E9 instrucțiuni.
 - a. O interpretare greșită uzuală este să considerăm că procesorul cu ceasul mai mare are performanțe mai bune. Verificați dacă este adevărat pentru P1 și P2.
 - b. O altă interpretare greșită este aceea de a considera că procesorul care execută mai multe instrucțiuni va avea un timp CPU mai mare. Considerând că procesorul P1 execută o secvență de 1,0E9 și că CPI-ul celor două procesoare nu se schimbă, determinați numărul de instrucțiuni pe care le poate executa P2 în același timp necesar ca P1 să execute cele 1,0E9 instrucțiuni.
 - c. O altă interpretare greșită comună este aceea a utilizării MIPS pentru a compara performanțele a două procesoare, și se consideră că procesorul cu MIPS mai mare este mai performant. Verificați dacă este valabil pentru P1 și P2.
 - d. O altă interpretare greșită este utilizarea MFLOPS. Presupuneți că 40% din instrucțiunile executate atât de P1 cât și de P2 sunt instrucțiuni în virgulă flotantă. Găsiți interpretarea MFLOPS greșită pentru programe.
24. O altă capcană citată în secțiunea 1.10 este aceea că ne așteptăm să îmbunătățim performanțele generale ale unui calculator prin îmbunătățirea doar a unui aspect. Să considerăm un calculator

care rulează un program care necesită 250s, cu 70s petrecute executând instrucțiuni FP, 85 de secunde executând instrucțiuni L/S, și 40s executând instrucțiuni de salt.

- a. Cât de mult se reduce timpul total dacă timpul pentru operațiile FP se reduce cu 20%?
- b. Cât de mult se reduce timpul pentru operațiile INT dacă timpul total se reduce cu 20%?
- c. Poate fi redus timpul total cu 20% reducând numai operațiile de salt?

25. Presupuneți că un program necesită execuția a 50 x 10⁶ instrucțiuni FP, 110 x 10⁶ instrucțiuni INT, 80 x 10⁶ instrucțiuni L/S, și 16 x 10⁶ instrucțiuni de salt. CPI pentru fiecare tip de instrucțiune este de respectiv 1, 1, 4, și 2. Presupuneți că procesorul are ceasul de 2 GHz.

- a. Cât de mult trebuie îmbunătățit CPI-ul instrucțiunilor FP dacă dorim ca programul să ruleze de două ori mai repede?
- b. Cât de mult trebuie îmbunătățit CPI-ul instrucțiunilor L/S dacă dorim ca programul să ruleze de două ori mai repede?
- c. De câte ori este îmbunătățit timpul de execuție al programului dacă CPI pentru instrucțiunile INT și FP este redus cu 40% și CPI pentru instrucțiunile L/S și de salt este redus cu 30%.

26. Atunci când un program este adaptat să se execute pe procesoare multiple într-un sistem multiprocesor, timpul de execuție pe fiecare procesor este alcătuit din timpul de calcul și timpul dat de supracontrolul necesar pentru a bloca secțiunile critice și/sau trimiterea datelor de la un procesor la altul. Presupuneți un program care necesită 100s de execuție pe un singur procesor. Atunci când se rulează pe p procesoare necesită t/ps precum și 4s adiționale pentru supracontrol, indiferent de numărul de procesoare. Calculați timpul de execuție per procesor pentru 1, 2, 4, 8, 16, 32, 64, și 128 de procesoare. Pentru fiecare caz, listați creșterea relativă de viteză la un singur procesor și raportul dintre creșterea de viteză actuală și creșterea de viteză ideală (fără supracontrol).

Probleme propuse pentru examenul de la
disciplina Structura și Organizarea
Calculatoarelor

== MOODLE TYPE ==

Prof. dr. ing. Vasile Gheorghiuță GĂITAN
Ș. I. dr. ing. Ionel ZAGAN

Cap. 2 Instrucțiunile: limbajul calculatorului

1. Înainte de anii 1980, calculatoarele erau construite cu seturi instrucțiuni din ce în ce mai complexe. MIPS este un calculator RISC. De ce a existat o mutare de la calculatoarele cu set complex de instrucțiuni la calculatoarele RISC.
2. Scrieți următoarea secvență în cod assembler MIPS.
 $x = x + y + z - q$.
Presupuneți că x,y,z,q sunt memorate în \$s1 - \$s4.
3. În assemblerul MIPS, scrieți o versiune în limbaj de asamblare a următorului segment de cod C:

```
int A[100], B[100];
for (i=1; i<100, i++) {
    A[i] = A[i-1] + B[i];
}
```


La începutul acestui segment de cod, singurele valori din regiștrii sunt adresele de bază ale ariilor A și B în regiștrii \$a0 și \$a1. Evitați utilizarea instrucțiunilor de înmulțire.
4. Unele calculatoare un registru special pentru indicatori care conțin indicatori de stare. Acești biți sunt includ adesea biții *carry* și *overflow*. Descrieți diferența între funcționalitățile acestor doi biți și dați un exemplu de operație aritmetică care să conducă la starea lor pe valori diferite.
5. Setul de instrucțiuni MIPS include câteva instrucțiuni de deplasare. Acestea sunt deplasare logică la stânga, deplasare logică la dreapta, și deplasare aritmetică la dreapta. Alte arhitecturi includ numai o instrucțiune de deplasare aritmetică la dreapta.
 - a. De ce nu conține MIPS un opcode pentru deplasare aritmetică la stânga?
 - b. Cum se poate implementa în assembler o pseudo-operație de deplasare logică la stânga (LSL) pentru un procesor care nu are această instrucțiune? Fiți siguri că instrucțiunea voastră LSL poate realiza deplasarea până la W-biți, unde W este mărimea în biți a cuvântului procesor.
6. Considerați următorul cod în assembler pentru partea 1 și partea 2.

```
r1 = 99

Loop:
    r1 = r1 - 1
    branch r1 > 0, Loop
halt
```

 - a. Pe durata execuției codului de mai sus câte instrucțiuni dinamice sunt executate?
 - b. Presupuneți ca aveți un procesor uniciclu (cu un singur ciclu procesor) standard funcționând la 100 KHz, cât de lung va fi codul de mai sus până la completare?
7. Converteți funcția C de mai jos în limbaj de asamblare MIPS. Fiți siguri că, codul vostru în limbaj de asamblare poate fi apelat de un program C standard (altfel spus, fiți sigur că respectați convențiile MIPS de apel).

```
unsigned int sum(unsigned int n)
{
    if (n == 0) return 0;
    else return n + sum(n-1);
}
```

Acest procesor nu are sloturi de întârziere. Stiva crește spre adrese mai mici de memorie. Următorii regiștrii sunt utilizați în convențiile de apel:

Numele registrului	Numărul registrului	Utilizare
\$zero	0	Constanta 0
\$at	1	Rezervat pentru asamblor

\$v0, \$v1	2, 3	Valorile de revenire ale funcției
\$a0 - \$a3	4 – 7	Valorile argumentelor funcției
\$t0 - \$t7	8 – 15	Temporar (salvate de apelant)
\$s0 - \$s7	16 – 23	Temporar (salvate de apelat)
\$t8, \$t9	24, 25	Temporar (salvate de apelant)
\$k0, \$k1	26, 27	Rezervate pentru nucleul OS
\$gp	28	Pointer la aria globală
\$sp	29	Pointerul de stivă - Stack Pointer
\$fp	30	Pointerul de cadru Frame Pointer
\$ra	31	Adresa de revenire

8. În secvența de cod MIPS ce urmează, de câte ori este accesată instrucțiunea de acces la memorie? De câte ori este accesată memoria de date? (Numărați doar accesele la memorie nu și la regiștrii).

```
lw $v1, 0($a0)
addi $v0, $v0, 1
sw $v1, 0($a1)
addi $a0, $a0, 1
```

9. Utilizați valorile din regiștrii și memorie din tabelul următor pentru următoarea întrebare. Presupuneți un procesor pe 32 de biți. Presupuneți că fiecare din următoarele întrebări pornește de la valorile din tabel; nu utilizați valorile modificate de o întrebare ca propagate în viitoarea parte a întrebării.

Registru	Valoare	Locație de memorie	Valoare
R1	12	12	16
R2	16	16	20
R3	20	20	24
R4	24	24	28

- Dați valorile pentru R1, R2 și R3 după execuția instrucțiunii add R3, R2, R1.
 - Ce valori vor fi în R1 și R3 după execuția instrucțiunii load R3, 12(r1).
 - Ce valori vor fi în regiștrii după execuția instrucțiunii addi R2,R3,#16.
10. Loop unrolling and Fibonacci: fie următorul pseudo C-cod pentru a calcula numărul 5 Fibonacci (F(5)).

```
1 int a,b,i,t;
2 a=b=1;          /* Set a and b to F(2) and F(1) respectively */
3 for(i=0;i<2;i++)
4 {
5 t=a;            /* save F(n-1) to a temporary location */
6 a+=b;          /* F(n) = F(n-1) + F(n-2) */
7 b=t;           /* set b to F(n-1) */
8 }
```

Codul mai poate fi scris și astfel:

```
1 int a,b,t;
2 a=b=1;
3 t=a;
4 a+=b;
5 b=t;
6 t=a;
7 a+=b;
8 b=t;
```

- Converteți pseudo C-codul pentru ambele secvențe de mai sus într-un cod MIPS rezonabil de eficient. Reprezentați fiecare variabilă din pseudo C-cod cu un registru. Încercați să urmăriți pseudocodul cât mai aproape posibil (prima secvență trebuie să conțină o buclă, în timp ce a doua nu).

- b. Presupuneți acum că în loc de numărul 5 Fibonacci decideți să îl calculați pe al 20-lea. Cât de multe instrucțiuni statice vor fi în prima versiune și cât de multe vor fi în versiunea unrolled? Ce puteți spune despre instrucțiunile dinamice? Aici nu trebuie scris codul în asamblare.
11. Scrieți în limbajul de asamblare MIPS următorul segment de cod C:
- ```
for (i = 0; i < 98; i++) {
 C[i] = A[i + 1] - A[i] * B[i + 2]
}
```
- Ariile A, B și C sunt plasate la adresele de memorie 0xA000, 0xB000 și respectiv 0xC000. Încercați să deduceți numărul total de instrucțiuni și numărul de instrucțiuni scumpe precum înmulțirea.
12. Să presupunem că se implementează o nouă instrucțiune MIPS, denumită bcp, pentru a copia un bloc de cuvinte de la o adresă la alta. Presupunem că această instrucțiune cere ca adresa de start a blocului sursă să fie în registrul \$t1, iar adresa blocului destinație în registrul \$t2. Instrucțiunea necesită și numărul de cuvinte ce urmează a fi copiate și pentru aceasta se folosește registrul \$t3 (>0). Mai mult presupunem că valorile acestor regiștrii precum și a registrului \$t4 sunt distruse de această instrucțiune. Realizați următoarele: Scrieți cod în limbajul de asamblare MIPS pentru a implementa copierea unui bloc fără această instrucțiune. Scrieți cod în limbajul de asamblare MIPS pentru a implementa copierea unui bloc cu această instrucțiune. Estimați numărul total de cicli necesari pentru fiecare realizare pentru a copia 100 de cuvinte pe un procesor multiciclu.
13. Care este codul în limbajul de asamblare MIPS pentru următorul program în C? Presupuneți că variabilele f, g, h, și i sunt date și pot fi considerate întregi pe 32 de biți așa cum se declară în programul în C. Utilizați un număr minim de instrucțiuni în assembler MIPS.
- ```
f = g + (h - 5).
```
14. Fiind dat următorul cod în limbajul de asamblare MIPS, care este codul corespunzător în C?
- ```
add f, g, h
add f, i, f
```
15. Care este codul în limbajul de asamblare MIPS pentru următorul program în C? Presupuneți că variabilele f, g, h, i, și j sunt asignate regiștrilor \$s0, \$s1, \$s2, \$s3, și respectiv \$s4. Presupuneți că adresele de bază ale ariilor A și B sunt în regiștrii \$s6 și respectiv \$s7.
- ```
B[8] = A[i-j];
```
16. Care sunt liniile C corespunzătoare codului în limbaj de asamblare MIPS de mai jos? Presupuneți că variabilele f, g, h, i, și j sunt asignate regiștrilor \$s0, \$s1, \$s2, \$s3, și respectiv \$s4. Presupuneți că adresa de bază a ariilor A și B sunt în regiștrii \$s6 și respectiv \$s7.

```
sll $t0, $s0, 2    # $t0 = f * 4
add $t0, $s6, $t0  # $t0 = &A[f]
sll $t1, $s1, 2    # $t1 = g * 4
add $t1, $s7, $t1  # $t1 = &B[g]
lw  $s0, 0($t0)    # f = A[f]
addi $t2, $t0, 4
lw  $t0, 0($t2)
add $t0, $t0, $s0
sw  $t0, 0($t1)
```

17. Pentru codul în limbaj de asamblare din exercițiul 16, rescrieți codul pentru a minimiza numărul de instrucțiuni MIPS (dacă este posibil) necesare pentru a realiza aceeași funcție.
18. Tabelul următor prezintă valorile pe 32 de biți a unei arii stocate în memorie.

Adresa	Data
24	2
38	4
32	3
36	6
40	1

- a. Pentru locațiile de memorie din tabelul de mai sus scrieți cod în C pentru a sorta data de la valoarea cea mai mică spre valoarea cea mai mare, plasând valoarea cea mai mică la adresa de memorie cea mai mică. Presupuneți că datele prezentate în tabel reprezintă o variabilă C denumită Array, care este o arie de tip int, și că primul număr din arie este primul element din aceasta. Presupuneți că procesorul este cu adresare la nivel de octet și că un cuvânt constă din patru octeți.
 - b. Pentru locațiile de memorie din tabelul de mai sus scrieți cod în limbajul de asamblare MIPS pentru a sorta data de la valoarea cea mai mică spre valoarea cea mai mare, plasând valoarea cea mai mică la adresa de memorie cea mai mică. Utilizați un număr minim de instrucțiuni MIPS. Presupuneți că adresa de bază a Array este memorată în registrul \$s6.
19. Prezentați modul în care valoarea 0xabcdef12 va fi aranjată în memorie pe un procesor little-endian respectiv big-endian. Presupuneți că data este aranjată începând cu adresa 0.
 20. Translați în zecimal valoarea 0xabcdef12.
 21. Translați următorul cod în limbaj de asamblare MIPS. Presupuneți că variabilele f, g, h, i, și j sunt asignate regiștrilor \$s0, \$s1, \$s2, \$s3, și respectiv \$s4. Presupuneți că adresele de bază a ariilor A și B sunt în regiștrii \$s6 și respectiv \$s7. Presupuneți că elementele ariilor A și B sunt cuvinte pe 4 octeți: $B[8] = A[i] + A[j]$.
 22. Translați următorul cod scris în limbaj de asamblare MIPS în cod C. Presupuneți că variabilele f, g, h, i, și j sunt asignate regiștrilor \$s0, \$s1, \$s2, \$s3, și respectiv \$s4. Presupuneți că adresele de bază pentru ariile A și B sunt în regiștrii \$s6 și respectiv \$s7.

```

addi $t0, $s6, 4
add $t1, $s6, $0
sw $t1, 0($t0)
lw $t0, 0($t0)
add $s0, $t1, $t0

```
 23. Pentru fiecare instrucțiune MIPS precizați valoarea opcode-ului (OP), și câmpurile pentru registrul sursă (RS), și registrul destinație (RT). Pentru instrucțiunile I-type precizați valoarea câmpului imediat, și pentru instrucțiunile R-type precizați valoarea câmpului pentru registrul destinație (RD).
 24. Presupuneți că regiștrii \$s0 și \$s1 memorează valorile 0x80000000 și respectiv 0xD0000000.
- a. Care este valoarea din registrul \$t0 după execuția următorului cod în limbaj de asamblare MIPS?
 - a. add \$t0, \$s0, \$s1
 - b. Rezultatul din \$t0 este rezultatul dorit sau a apărut un overflow?
 - c. Pentru valorile din regiștrii \$s0 și \$s1 specificate anterior care este valoarea din registrul \$t0 după execuția următorului cod în limbaj de asamblare MIPS?
 - a. sub \$t0, \$s0, \$s1
 - d. Rezultatul din \$t0 este rezultatul dorit sau a apărut un overflow?
 - e. Pentru valorile din regiștrii \$s0 și \$s1 specificate anterior care este valoarea din registrul \$t0 după execuția următorului cod în limbaj de asamblare MIPS?
 - a. add \$t0, \$s0, \$s1
 - b. add \$t0, \$t0, \$s0
 - f. Rezultatul din \$t0 este rezultatul dorit sau a apărut un overflow?
25. Presupunem că \$s0 conține valoarea holds the value 128_{10} .
 - a. Pentru instrucțiunea \$t0, \$s0, \$s1, care este gama (-ele) valorilor din \$s1 pentru care va rezulta overflow?
 - b. Pentru instrucțiunea sub \$t0, \$s0, \$s1, care este gama (-ele) valorilor din \$s1 pentru care va rezulta overflow?
 - c. Pentru instrucțiunea sub \$t0, \$s1, \$s0, care este gama (-ele) valorilor din \$s1 pentru care va rezulta overflow?
 26. Precizați tipul și instrucțiunea în limbaj de asamblare MIPS pentru următoarea valoare binară 0000 0010 0001 0000 1000 0000 0010 0000₂
 27. Precizați tipul și reprezentarea hexazecimală a următoarei instrucțiuni: sw \$t1, 32(\$t2).

28. Precizați tipul, instrucțiunea în limbaj de asamblare, și reprezentarea binară a instrucțiunii descrisă de următoarele câmpuri MIPS: $op=0$, $rs=3$, $rt=2$, $rd=3$, $shamt=0$, $funct=34$.
29. Precizați tipul, instrucțiunea în limbaj de asamblare, și reprezentarea binară a instrucțiunii descrisă de următoarele câmpuri MIPS: $op=0x23$, $rs=1$, $rt=2$, $const=0x4$.
30. Presupuneți că am dori să expandăm fișierul de regiștrii MIPS la 128 de regiștrii, precum și setul de instrucțiuni pentru a se permite de patru ori mai multe instrucțiuni.
- Cum va afecta aceasta mărimea fiecărui câmp de biți pentru instrucțiunile de tip R?
 - Cum va afecta aceasta mărimea fiecărui câmp de biți pentru instrucțiunile de tip I?
 - Cum ar putea fiecare din cele două schimbări propuse să reducă mărimea unui program scris în limbaj de asamblare MIPS?
31. Presupuneți următorul conținut pentru regiștrii: $\$t0 = 0xAAAAAAAA$, $\$t1 = 0x12345678$
- Pentru valorile de mai sus din regiștrii care este valoarea din registrul $\$t2$ pentru următoarea secvență de instrucțiuni?
`sll $t2, $t0, 44`
`or $t2, $t2, $t1`
 - Pentru valorile de mai sus din regiștrii care este valoarea din registrul $\$t2$ pentru următoarea secvență de instrucțiuni?
`sll $t2, $t0, 4`
`andi $t2, $t2, -1`
 - Pentru valorile de mai sus din regiștrii care este valoarea din registrul $\$t2$ pentru următoarea secvență de instrucțiuni?
`srl $t2, $t0, 3`
`andi $t2, $t2, 0xFFEF`
32. Găsiți cea mai scurtă secvență de instrucțiuni MIPS care extrage biții de la 16 la 11 din registrul $\$t0$ și utilizează valorile acestor biți pentru a înlocui biții de la 31 la 26 din registrul $\$t1$ fără a schimba ceilalți 26 de biți din registrul $\$t1$.
33. Furnizați un set minimal de instrucțiuni MIPS care poate fi utilizat pentru a implementa următoarea pseudoinstrucțiune: `not $t1, $t2 // bit-wise invert`.
34. Pentru următoarea secvență de linii scrise în C, scrieți o secvență minimală de instrucțiuni MIPS în limbaj de asamblare care realizează operații identice. Presupuneți că $\$t1 = A$, $\$t2 = B$, și $\$s1$ este adresa de bază pentru C. $A = C[0] \ll 4$.
35. Presupuneți că $\$t0$ memorează valoarea $0x00101000$. Care este valoarea lui $\$t2$ după execuția următoarelor instrucțiuni?
`sll $t2, $0, $t0`
`bne $t2, $0, ELSE`
`j DONE`
`ELSE: addi $t2, $t2, 2`
`DONE:`
36. Presupuneți că program counter (PC) este setat cu valoarea $0x2000\ 0000$. Este posibil să utilizați instrucțiunea MIPS în limbaj de asamblare `jump (j)` pentru a seta PC-ul cu adresa $0x4000\ 0000$? Este posibil ca să utilizați instrucțiunea MIPS în limbaj de asamblare `branch-on-equal (beq)` pentru a seta PC-ul cu aceeași adresă?
37. Următoarea instrucțiune nu este inclusă în setul de instrucțiuni MIPS:
`rpt $t2, loop # if(R[rs]>0) R[rs]=R[rs]-1, PC=PC+4+BranchAddr`
- Dacă această instrucțiune ar trebui implementată în setul de instrucțiuni MIPS care ar fi formatul cel mai corespunzător?
 - Care este cea mai scurtă secvență de instrucțiuni MIPS care ar realiza aceleași operații?
38. Fie următoarea buclă MIPS:

```

LOOP: slt $t2, $0, $t1
      beq $t2, $0, DONE
      subi $t1, $t1, 1
      addi $s2, $s2, 2
      j LOOP
DONE:

```

- Presupuneți că registrul \$t1 este inițializat cu valoarea 10. Care este valoarea din \$s2 presupunând că \$s2 are valoarea inițială zero?
- Pentru fiecare buclă de mai sus, scrieți codul echivalent al unei rutine C. Presupuneți că regiștrii \$s1, \$s2, \$t1, și \$t2 conțin întregii A, B, i, și respectiv temp.
- Presupuneți, pentru buclele de mai sus scrise în limbaj de asamblare MIPS2, că registrul \$t1 este inițializat cu valoarea N. Câte instrucțiuni MIPS sunt executate?

- Translatați următorul cod C în cod în limbaj de asamblare MIPS. Utilizați un număr minim de instrucțiuni. Presupuneți că valorile a, b, i, și j sun în regiștrii \$s0, \$s1, \$t0, li respectiv \$t1. De asemenea presupuneți că registrul \$s2 păstrează adresa de bază pentru aria D.

```

for(i=0; i<a; i++)
  for(j=0; j<b; j++)
    D[4*j] = i + j;

```

- Câte instrucțiuni MIPS sunt necesare pentru a implementa codul C din exercițiul anterior (39)? Dacă variabilele a și b sunt inițializate cu 10 și 1 și toate elementele din D sunt inițial 0, care este numărul total de instrucțiuni MIPS care se execută pentru a completa bucla?
- Translatați următoarea buclă în C. Presupuneți că întregul i de la nivelul C este păstrat în registrul \$t1, \$s2 păstrează întregul de la nivelul C denumit result, iar \$s0 păstrează adresa de bază al întregului.

```

      addi $t1, $0, 0
LOOP: lw $s1, 0($s0)
      add $s2, $s2, $s1
      addi $s0, $s0, 4
      addi $t1, $t1, 1
      slti $t2, $t1, 100
      bne $t2, $s0, LOOP

```

- Rescrieți bucle din exercițiul anterior reducând numărul de instrucțiuni MIPS executate.
- Implementați următorul cod scris în C în cod scris în limbaj de asamblare MIPS. Care este numărul total de instrucțiuni MIPS necesare pentru a executa funcția?

```

int fib(int n){
    if (n==0)
        return 0;
    else if (n == 1)
        return 1;
    else
        return fib(n-1) + fib(n-2);
}

```

- Funcțiile pot fi adesea scrise de compilatoare "in-line." O funcție in-line este atunci când corpul funcției este copiat în spațiul programului, permițând eliminarea supracontrolului dat de apelul funcției. Implementați o versiune "in-line" version of the the **C code in the table in MIPS assembly**. Care este reducerea instrucțiunilor în asamblare MIPS necesare pentru a realiza funcția? Presupuneți că variabila C n este inițializată cu 5.
- Pentru fiecare apel de funcție, prezentați conținutul stivei după ce este realizat apelul de funcție. Presupuneți că poiterul de stivă are adresa originală 0x7ffffffc, și urmează convenția specificată în figura 2.11.

Preserved	Not preserved
Saved registers: \$s0-\$s7	Temporary registers: \$t0-\$t9
Stack pointer register: \$sp	Argument registers: \$a0-\$a3
Return address register: \$ra	Return value registers: \$v0-\$v1
Stack above the stack pointer	Stack below the stack pointer

46. Translați funcția f în cod scris în limbaj de asamblare MIPS. Dacă utilizați regiștrii de la \$t0 la \$t7, utilizați prima dată regiștrii cu numerele mici. Presupuneți că declarația funcției pentru func este "int f(int a, int b);". Codul funcției f: `int f(int a, int b, int c, int d){ return func(func(a,b),c+d);}`
47. Se poate utiliza optimizarea tail-call pentru această funcție? Dacă nu, explicați de ce. Dacă da, care este diferența privind numărul de instrucțiuni executate în f cu și fără optimizare?
48. Chiar înainte ca funcția f (două exerciții mai sus -46) să revină, ce știm despre regiștrii \$t5, \$s3, \$ra, și \$sp? Amintiți-vă faptul că se știe cum arată toată funcția f , dar pentru func știm doar declarația sa.
49. Scrieți un program în limbajul de asamblare MIPS pentru a converti un șir ASCII ce reprezintă întregi zecimali pozitivi sau negativi la un întreg. Programul trebuie să se aștepte că registrul \$a0 conține adresa unui șir terminat cu 0 și care conține combinații de digiți de la 0 la 9. Programul vostru trebuie să calculeze valoarea întreagă echivalentă aceluia șir de digiți, după care să plaseze numărul în registrul \$v0. Dacă apare un caracter non-digit oriunde în șir programul se va opri și va returna valoarea -1 în registrul \$v0. De exemplu dacă registrul \$a0 pointează o secvență de trei octeți 50_{10} , 52_{10} , 0_{10} (șirul terminat cu zero este "24"), Când programul se va opri, registrul \$v0 trebuie să conțină valoarea 24_{10} .
50. Pentru următorul cod:
`lbu $t0, 0($t1)`
`sw $t0, 0($t2)`
 Presupuneți că registrul \$t1 conține adresa 0x1000 0000, iar registrul \$t2 conține adresa 0x1000 0010. Amintim că arhitectura MIPS utilizează adresarea big-endian. Presupuneți că data (în hexazecimal) la adresa 0x1000 0000 este: 0x11223344. Care este valoarea memorată la adresa pointată de registrul \$t2?
51. Scrieți cod în limbaj de asamblare MIPS care creează constanta 0010 0000 0000 0001 0100 1001 0010 0100₂ pe 32 de biți și memorați această valoare în registrul \$t1.
52. Dacă valoarea curentă a PC-ului este 0x00000000, se poate utiliza o singură instrucțiune de salt pentru a prelua în PC adresa din exercițiul precedent? Exercise 2.39 (51)?
53. Dacă valoarea curentă a PC-ului este 0x00000600, se poate utiliza o singură instrucțiune branch pentru a prelua în PC adresa din exercițiul - Exercise 2.39 (51)?
54. Dacă valoarea curentă a PC-ului este 0x1FFFf000, se poate utiliza o singură instrucțiune branch pentru a prelua în PC adresa din exercițiul - Exercise 2.39 (51)?
55. Scrieți cod în limbaj de asamblare MIPS pentru a implementa următorul cod C:
`lock(lk);`
`shvar=max(shvar,x);`
`unlock(lk);`
 Presupuneți că adresa variabilei lk este în \$a0, adresa variabilei shvar este în \$a1, și valoarea variabilei x este în \$a2. Secțiunea voastră critică nu poate conține nici un apel de funcție Utilizați instrucțiunile ll/sc pentru a implementa operația lock(), iar operația unlock() este simplu o instrucțiune ordinară de memorare.
56. Repetă exercițiul anterior (Exercise 2.43-55), dar de data această utilizează ll/sc pentru a realiza o actualizare atomică direct a variabilei shvar, fără a utiliza lock() și unlock(). Amintim că în această problemă nu există variabila lk

57. Utilizând codul din exercițiul (Exercise 2.43 – 55) ca un exemplu, explicați ce se întâmplă când două procesoare încep să execute această secțiune critică în același timp, presupunând că fiecare procesor execută exact o instrucțiune pe ciclu.
58. Presupuneți că pentru un procesor dat CPI pentru instrucțiunile aritmetice este 1, CPI pentru instrucțiunile load/store este 10, și CPI pentru instrucțiunile branch este 3. Presupuneți că programul are următoarele instrucțiuni executate: 500 milioane de instrucțiuni aritmetice, 300 de milioane de instrucțiuni load/store, 100 milioane de instrucțiuni branch.
- Presupuneți că s-au adăugat la setul de instrucțiuni, instrucțiuni aritmetice noi mult mai puternice. Pe medie, prin utilizarea acestor instrucțiuni aritmetice mai puternice, putem reduce cu 25% numărul de instrucțiuni aritmetice executate de program, iar costul creșterii ciclului de ceas este de numai 10%. Este aceasta o bună alegere de proiectare? De ce?
 - Presupuneți că am găsit o cale pentru a dubla performanțele instrucțiunilor aritmetice. Care este creșterea generală de viteză la nivelul procesorului? Care este creșterea generală de viteză la nivelul procesorului dacă găsim o cale de a mări performanțele instrucțiunilor aritmetice de 10 ori?
59. Presupuneți că pentru un program dat 70% din instrucțiunile executate sunt aritmetice, 10% sunt load/store, și 20% sunt branch.
- Fiind dat amestecul de instrucțiuni și presupunând că instrucțiunile aritmetice necesită 2 cicli, o instrucțiune load/store necesită 6 cicli, și o instrucțiune branch necesită 3 cicli, găsiți CPI-ul mediu.
 - Pentru o îmbunătățire în performanță de 25%, cât de mulți cicli, în medie, poate necesita o instrucțiune aritmetică dacă instrucțiunile load/store și branch nu sunt îmbunătățite?
 - Pentru o îmbunătățire în performanță de 50%, cât de mulți cicli, în medie, poate necesita o instrucțiune aritmetică dacă instrucțiunile load/store și branch nu sunt îmbunătățite?

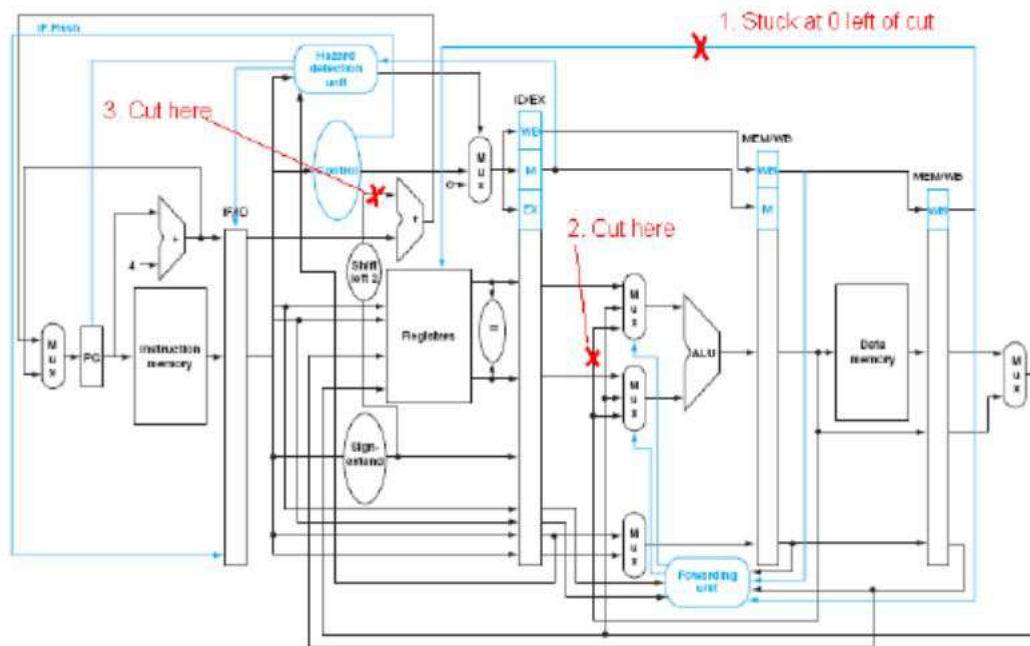
Probleme propuse pentru examenul de la
disciplina Structura și Organizarea
Calculatoarelor

== MOODLE TYPE ==

Prof. dr. ing. Vasile Gheorghiuță GĂITAN
Ș. I. dr. ing. Ionel ZAGAN

Cap. 3 Procesorul

1. Pentru calea de date MIPS din figură, unele linii sunt marcate cu X. pentru fiecare:
 - a. Descrieți în cuvinte consecințele negative ale tăierii acestora relativ la procesorul care lucrează fără modificări.
 - b. Furnizați un fragment de cod care va eșua.
 - c. Furnizați un fragment de cod care totuși va funcționa.



2. Considerați următorul cod în limbajul de asamblare:

```

I0: ADD R4 = R1 + R0;
I1: SUB R9 = R3 - R4;
I2: ADD R4 = R5 + R6;
I3: LDW R2 = MEM[R3 + 100];
I4: LDW R2 = MEM[R2 + 0];
I5: STW MEM[R4 + 100] = R2;
I6: AND R2 = R2 & R1;
I7: BEQ R9 == R1, Target;
I8: AND R9 = R9 & R1;
  
```

Luăți în considerare o pipeline cu avansare (forwarding), detecția hazardului și un slot de întârziere pentru ramificații. Pipeline-ul este unul tipic MIPS cu 5 etaje IF, ID, EX, MEM, WB. Pentru codul anterior, completați diagrama care urmează (instrucțiunile în stânga, ciclul în partea de sus). Înserați caracterele IF, ID, EX, MEM, WB pentru fiecare instrucțiune. Presupuneți că există două niveluri de bypassing, conform cărora a doua jumătate a etajului de decodificare realizează o citire a regiștrilor sursă, iar prima jumătate a etajului write-back scrie în fișierul de regiștrii. Etichetați în tabel cu X toate încetinirile (stalls). Etichetați toate avansările (forwarding) pe care unitatea de avansare le detectează (săgeți între etajele care produc datele și etajele care recepționează datele). Care este timpul final de execuție a codului?

	T0	T1	T2	T3	T4	T5	T6	T7	T8	T9	T10	T11	T12	T13	T14	T15	T16	T17
I0																		
I1																		
I2																		
I3																		
I4																		
I5																		
I6																		
I7																		
I8																		

3. Hazardul structural, de date și de control în mod tipic necesită ca pipeline-ul procesorului să încetinească (stall). În continuare sunt listate o serie de tehnici de optimizare implementate în compilator sau într-un pipeline al unui procesor proiectat pentru a reduce sau a elimina încetinirile (stalls) datorate acestor hazarduri. Pentru fiecare dintre următoarele tehnici de optimizare, precizați care tip de hazard este adresat și cum este adresat. Unele tehnici de optimizare pot adresa mai mult decât un singur tip de hazard, ca urmare fiți sigur ca ați inclus explicații pentru toate hazardurile adresate.

- Predicția salturilor (Branch prediction)
- Planificarea instrucțiunilor (Instruction scheduling)
- Sloturi de întârziere (delay slots)
- Creșterea disponibilității unor unități funcționale (increasing availability of functional units (ALUs, adders, etc))
- Caches

4. Predicția salturilor. Considerați următoarea secvență a rezultatelor pentru o singură ramificație statică. T înseamnă că ramificația este efectuată. N înseamnă că ramificația nu este efectuată. Pentru această întrebare, presupuneți că este singura ramificație din program.

TTTNTNTTTNTNTTTNTN

- Presupuneți că încercăm să prezicem această secvență cu un BHT utilizând un numărător cu un singur bit. Numărătorul din BHT este

inițializat cu starea N. Care dintre salturile din această secvență va prezice greșit?

Predictor state before prediction	Branch outcome	Misprediction?
N	T	
	T	
	T	
	N	
	T	
	N	
	T	
	T	
	T	
	N	
	T	
	N	
	T	
	T	
	T	
	N	
	T	
	N	

- b) Acum presupuneți un predictor cu două niveluri care utilizează un bit pentru istoria ramificației – un BHR cu un bit. Deși există o singură ramificație în program, nu contează cum este BHR concatenat cu PC-ul ramificației ca index în BHT. Presupuneți că BHT utilizează un numărător de un bit și ca urmare, din nou, toate intrările sunt inițializate cu N. Care dintre salturile din secvență vor prezice greșit?

Predictor state before prediction	BHR	Branch outcome	Misprediction?
N	N	T	
		T	
		T	
		N	
		T	
		N	
		T	
		T	
		T	
		N	
		T	
		N	
		T	
		T	
		T	
		N	
		T	
		N	

- c) Care este adresa de revenire de pe stivă? Când este actualizată adresa de pe stivă?

5. Pipeline-ul clasic are 5 etaje IF, ID, EX, MEM, WB. Acest pipeline este proiectat în mod specific pentru a executa setul de instrucțiuni MIPS. MIPS este o arhitectură load-store care realizează o operație cu memoria per instrucțiune, ca urmare este suficient un singur etaj MEM în pipeline. De asemenea, modul de adresare cel mai comun este modul de adresare prin registru cu deplasament. Etajul EX este plasat înainte de MEM pentru a permite ca EX să fie utilizat pentru calculul adresei. În această întrebare se va lua în considerare o variație în setul de instrucțiuni MIPS și interacțiunea acestei variații în structura pipeline.

Variația care se ia în calcul implică comutarea etajelor MEM și EX, creându-se o nouă pipeline IF, ID, MEM, EX, WB. Această schimbare are două efecte asupra setului de instrucțiuni. Primul, ne previne asupra utilizării adresării pe bază de registru cu deplasament (nu există un EX înainte de MEM care să realizeze adunarea cu semn). Totuși, putem utiliza instrucțiuni cu un operand de intrare în memorie, cum ar fi instrucțiuni cu adresare prin registru. De exemplu `multf_m f0,f2,(r2)` înmulțește conținutul registrului f2 cu valoarea locației de memorie pointată de r2 și pune rezultatul în f0.

- Eliminarea modului de adresare prin registru cu deplasament este o mare pierdere, pentru că este modul cel mai utilizat în MIPS. De ce este așa de frecvent? Enumerați două construcții populare care utilizează modul de adresare prin registru cu deplasament (utilizați adresarea cu deplasament cu deplasament diferit de zero).
- Care este diferența între dependențe și un hazard?
- În această întrebare se lucrează cu bucla SAXPY.

do $I = 0, N$

$Z[I] = A * X[I] + Y[I]$

Aici este noul cod în asamblare.

```
0: slli      r2,r1,#3 // I is in r1
1: addi      r3,r2,#X
2: multf_m   f2,f0,(r3) // A is in f0
3: addi      r4,r2,#Y
4: addf_m    f4,f2,(r4)
5: addi      r4,r2,#Z
6: sf        f4,(r4)
7: addi      r1,r1,#1
8: slei      r6,r1,r5 // N is in r5
9: bnez      r6,#0
```

Utilizând numerele pentru instrucțiuni etichetați datele și dependențele de control

- Umpleți diagrama pipeline pentru codul nou al buclei SAXPY. Etichetați încetinirea (stall) ca d^* pentru hazardul de date și s^* pentru hazardul structural. Care este întârzierea pentru o singură iterație? (numărul de cicli între completarea a două `#0` instrucțiuni succesive). Pentru această întrebare, presupuneți că adunarea FP ia 2 cicli, înmulțirea FP ia trei cicli și că toate celelalte operații iau un singur ciclu. Unitățile funcționale nu sunt pipeline. Sumatorul FP, multiplicatorul FP și ALU pentru întregi sunt unități funcționale separate, astfel încât nu există hazard structural între ele. Fișierul

de regiștrii este scris de etajul WB în prima jumătate a ciclului de ceas și este citit de etajul ID în a doua jumătate a ciclului de ceas. În plus, procesorul are o unitate de avansare (forwarding) completă. Procesorul încetinește (stalls) ramificațiile până când rezultatul este disponibil la sfârșitul etajului EX. Procesorul nu are resurse pentru a menține o stare precisă.

Instruction	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24
0: sll r2,r1,#3	F	D	M	X	W																			
1: addi r3,r2,#X		F	D	M	X	W																		
2: mul f_m f2,f0,(r3)			F	D		M	X	X	X	W														
3: addi r4,r2,#Y				F		D	M	X	W															
4: addf_m f4,f2,(r4)						F	D		M	X	X	W												
5: addi r4,r2,#Z							F		D	M	X		W											
6: sf f4,(r4)									F	D		M	X	W										
7: addi r1,r1,#1										F		D	M	X	W									
8: slei r6,r1,r5												F	D	M	X	W								
9: bnez r6,#0													F	D	M	X	W							
0: sll r2,r1,#3																	F	D	M	X	W			

- În pipeline-ul pentru MIPS menționat în text, care este motivul pentru a forța ca operațiile non-memorie să treacă prin etajul MEM în loc să meargă direct în etajul WB?
- Pe lângă pierderea adresării prin registru cu deplasament și a instrucțiunilor necesare pentru a suplini acest neajuns, care sunt alte două dezavantaje ale acestui tip de pipeline?
- Reduceți încetinirea prin planificarea pipeline-ului cu o singură iterație a buclei. Prezentați codul rezultat și umpleți diagrama pipeline. Pentru un răspuns corect nu trebuie să prezentați o planificare optimală.

Instruction	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24
0: sll r2,r1,#3	F	D	M	X	W																			
1: addi r3,r2,#X		F	D	M	X	W																		
2: mul f_m f2,f0,(r3)			F	D		M	X	X	X	W														
3: addi r4,r2,#Y				F		D	M	X	W															
4: addf_m f4,f2,(r4)						F	D		M	X	X	W												
5: addi r4,r2,#Z							F		D	M	X		W											
6: sf f4,(r4)									F	D		M	X	W										
7: addi r1,r1,#1										F		D	M	X	W									
8: slei r6,r1,r5												F	D	M	X	W								
9: bnez r6,#0													F	D	M	X	W							
0: sll r2,r1,#3																	F	D	M	X	W			

6. O întrebare cu două părți

a) Detecția dependențelor

Această întrebare se referă la modul în care ați înțeles dependențele între instrucțiuni. Utilizând codul care urmează, listați toate dependențele de tipul RAW, WAR, WAW. Listați dependențele în tabelul următor (exemplu INST-X la INST-Y) prin scrierea numărului instrucțiunii implicate în dependență.

I0: $A = B + C$;
 I1: $C = A - B$;
 I2: $D = A + C$;
 I3: $A = B * C * D$;
 I4: $C = F / D$;
 I5: $F = A \wedge G$;
 I6: $G = F + D$;

RAW Dependence		WAR Dependence		WAW Dependence	
From Instr	To Instr	From Instr	To Instr	From Instr	To Instr

b) Analiza dependențelor

Fiind date patru instrucțiuni, câte comparații unice (între registrul sursă și destinație) sunt necesare pentru a găsi toate dependențele de tip RAW, WAR și WAW. Răspundeți pentru cazul cu patru instrucțiuni, și apoi derivați o ecuație generală pentru N instrucțiuni. Presupuneți că toate instrucțiunile au un registru destinație și doi regiștrii sursă.

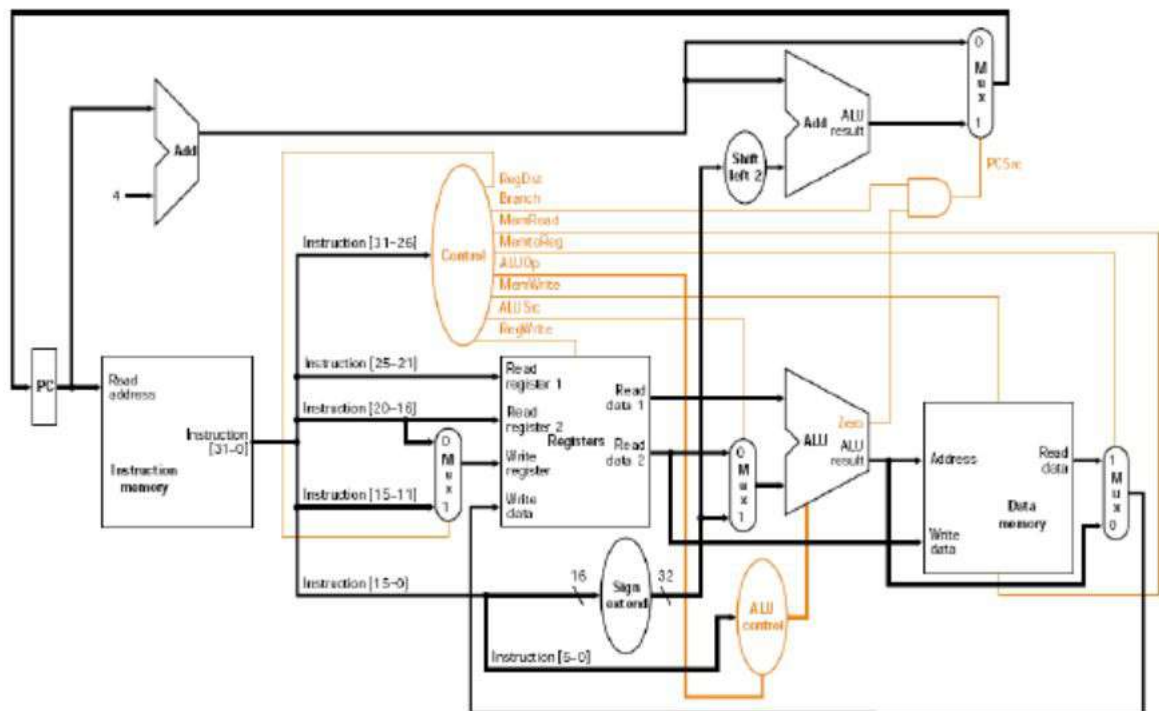
7. Se consideră o implementare cu un singur ciclu a ISA MIPS. Timpul de operare pentru componentele funcționale majore sunt prezentate în tabelul următor:

Component	Latency
ALU	10 ns
Adder	8 ns
ALU Control Unit	2 ns
Shifter	3 ns
Control Unit/ROM	4 ns
Sign/zero extender	3 ns
2-1 MUX	2 ns
Memory (read/write) (instruction or data)	15 ns
PC Register (read action)	1 ns
PC Register (write action)	1 ns
Register file (read action)	7 ns
Register file (write action)	5 ns
Logic (1 or more levels of gates)	1 ns

În figura următoare este o copie a proiectului pentru o cale de date cu un singur ciclu. În această implementare ciclul de ceas este determinat de calea cea mai lungă posibilă. Căile critice pentru diferite tipuri de instrucțiuni care trebuie luate în considerare sunt: Format-R, Load-word, și store-word. Toate instrucțiunile au aceeași pași pentru extragere și decodificare. Transferul de bază între regiștrii instrucțiunilor este:

Fetch/Decode: $\text{Instruction} \leftarrow \text{IMEM}[\text{PC}]$;
 R-type: $R[\text{rd}] \leftarrow R[\text{rs}] \text{ op } R[\text{rt}]$; $\text{PC} \leftarrow \text{PC} + 4$;
 load: $R[\text{rt}] \leftarrow \text{DMEM}[R[\text{rs}] + \text{signext}(\text{offset})]$; $\text{PC} \leftarrow \text{PC} + 4$;

store: $DMEM[R[rs] + \text{signext}(\text{offset})] \leftarrow R[Rt]; PC \leftarrow PC + 4;$



- a) În tabelul care urmează, indicați componentele care determină calea critică pentru respectiva instrucțiune, în cazul în care apare calea critică. Dacă este utilizată o componentă, dar nu este parte a căii critice a instrucțiunii (lucrează în paralel cu altă componentă), nu trebuie pusă în tabel. Fișierul de registre este utilizat pentru citire sau scriere; poate să apară de două ori pentru aceeași instrucțiune. Toate instrucțiunile încep prin a citi registrul PC cu o întârziere de 2ns.

Instruction Type	Hardware Elements Used By Instruction											
R-Format												
Load												
Store												

- b) Plasați întârzierile componentelor pe care ați decis că intră în calea critică a fiecărei instrucțiuni în tabelul care urmează. Calculați suma întârzierilor fiecărei componente pentru fiecare instrucțiune.

Instruction Type	Hardware Latencies For Respective Elements											Total
R-Format	2 ns											
Load	2 ns											
Store	2 ns											

- c) Utilizați coloana cu întârzierea totală pentru a obține următoarele informații privind calea critică:
- ☐ Fiind dată calea de date, care instrucțiune determină calea critică generală?

- ☐ Care va fi ciclul de ceas rezultat pe baza analizei căii critice?
- ☐ La ce frecvență va merge procesorul?

8. Figura următoare ilustrează trei tipuri de predictor,

- ☐ Ultimul luat - predicția este luată pe 1
- ☐ Sus-jos (Up-Down – numărător cu saturație) predicția este luată pentru 11 și 10
- ☐ Automatul A3 – este luat pentru 11 și 10
- ☐ Umpleți tabelul care urmează pentru fiecare predictor de ramificație. Șablonul de execuție pentru ramificația este NTNNTTTN.

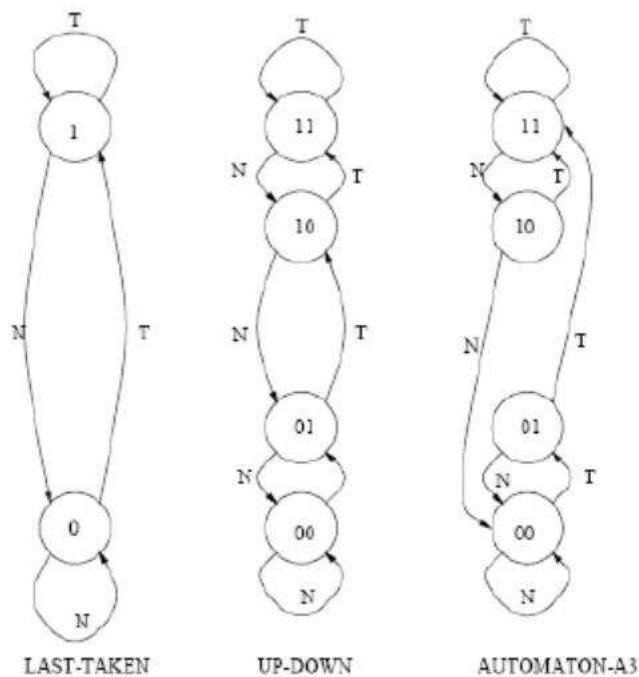


Table 1: Table for last-taken branch predictor

Execution Time	Branch Execution	State Before	Prediction	Correct or Incorrect	State After
0	N	0			
1	T				
2	N				
3	N				
4	T				
5	T				
6	T				
7	N				

Table 2: Table for saturating counter (up-down) branch predictor.

Execution Time	Branch Execution	State Before	Prediction	Correct or Incorrect	State After
0	N	01			
1	T				
2	N				
3	N				
4	T				
5	T				
6	T				
7	N				

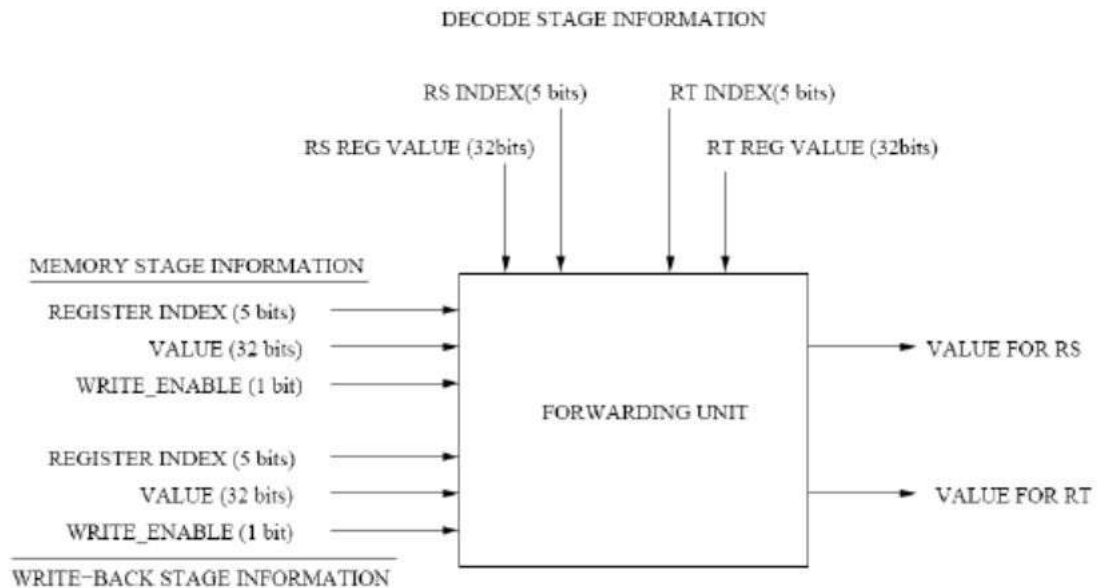
Table 3: Table for Automata-A3 branch predictor.

Execution Time	Branch Execution	State Before	Prediction	Correct or Incorrect	State After
0	N	01			
1	T				
2	N				
3	N				
4	T				
5	T				
6	T				
7	N				

Calculate the prediction rates of the three branch predictors:

Predictor	Prediction Accuracy
Last-taken	
Up-Down	
Automata-A3	

9. Pipeline este utilizat deoarece îmbunătățește fluxul instrucțiunilor. Creșterea nivelului de pipeline descrește cantitatea de muncă la fiecare etaj pipeline, permițând ca în procesor să fie executate mai multe instrucțiuni în același timp, iar o anumită instrucțiune să fie complet executată la o viteză mai mare. Totuși, fluxul nu se va îmbunătăți dacă pipeline-ul crește nedefinit. Dați două motive pentru această observație.
10. Proiectarea logicii de avansare. Pentru această problemă trebuie să proiectați unitatea de avansare pentru un procesor pipeline cu 5 etaje. Unitatea de avansare returnează valoare care trebuie avansată pentru instrucțiunea curentă. Există trei locuri de unde poate să vină valoarea pentru regiștrii RS și RT: etalul de decodificare (fișierul de regiștrii), etajul de memorie și etajul write-back (scrie înapoi).



Informația de la etajele write-back și memorie constă din: `_INDEX` – precizează indexul registrului în lucru care trebuie scris, `_VALUE` – valoarea care trebuie scrisă, `_ENABLE` – dacă este sau nu instrucțiunea din etaj una care scrie.

Etajul de decodificare memorează indexul regiștrilor (pentru RS și RT) și valorile corespunzătoare din fișierul de regiștrii.

În general pot exista trei valori, una din care unitatea de avansare poate alege una dintre cererile pentru regiștrii RS și RT. Etajul de memorie are valoarea MEM, etajul write-back are valoarea WB, și fișierul de regiștrii are valorile RS-REG sau RT-REG. Utilizând tabelul următor care conține informații despre toate etajele instrucțiunilor, indicați care valoare trebuie avansată (forwarding) pentru instrucțiunea curentă: MEM, WB, RS-REG sau RT-REG. Fiecare linie reprezintă o evaluare a unității de avansare; nu există conexiuni între liniile de evaluare din tabel. Nu trebuie să vă îngrijorați asupra detecției hazardului, ci numai de bypass-area valorii.

	Mem Stage		Write-Back Stage		Register Stage		RS Value	RT Value
Evaluation	Index	Write	Index	Write	RS-Index	RT-Index		
0	5	1	23	0	6	7		
1	7	0	16	1	16	8		
2	10	1	10	1	11	10		
3	17	0	12	1	12	12		
4	19	0	19	0	19	25		

11. Considerați un procesor MIPS cu 5 etaje pipeline cu durata unui ciclu de 10ns, Considerați că executați un program unde o fracție, f , din toate instrucțiunile urmează imediat un load de care devin dependente.

- Care este timpul total de execuție pentru N instrucțiuni, din f , cu unitatea de avansare validată?
- Considerați un scenariu pentru care etajul MEM, împreună cu registrul ei de pipeline, necesită 12 ns. Există două opțiuni: se adaugă încă un etaj MEM astfel încât să rezulte etajele MEM1 și MEM2, sau să se mărească durata ciclului mașină la 12 ns astfel încât etajul MEM să se încadreze în timpul noului ciclu și numărul de etaje pipeline să nu fie afectat. Pentru un program care mixează caracteristicile anterior descrise, când este prima opțiune mai bună decât a doua? Răspunsul trebuie să se bazeze pe valoarea lui f .

- c) Procesoarele înglobate (embedded) au două regiuni diferite de memorie – o memorie rapidă apropiată de procesor și o memorie normală mai înceată. Presupuneți că într-un procesor cu 6 etaje (cu MEM1 și MEM2), există o regiune de memorie care este mai rapidă și pentru care valoarea corectă se obține la sfârșitul etajului MEM1 în timp ce restul memoriei are nevoie de ambele etaje MEM1 și MEM2. Din motive de simplitate presupuneți că există două tipuri de instrucțiuni load și anume load.fast și load.slow care indică ce regiune de memorie este utilizată. Cum se schimbă răspunsul la întrebarea precedentă dacă 40% din fracția f menționată anterior își preia valorile din memoria rapidă?

12. Fie un procesor cu 4 etaje pipeline descris în continuare:

- ☐ IF: extragerea instrucțiunii
- ☐ IDE: decodificarea instrucțiunii, extragerea regiștrilor, evaluarea ALU, instrucțiunea de ramificație modifică PC, calculul adresei pentru accesul la memorie,
- ☐ MEM: accesul la memorie pentru instrucțiunile load și store,
- ☐ WB: scrie rezultatul execuției înapoi în fișierul de regiștrii. Scrierea înapoi apare în a doua jumătate a ciclului.

Presupuneți metoda de ramificație întârziată. Pentru programul următor, presupuneți că bucla va itera de 15 ori. Presupuneți că pipeline-ul termină o instrucțiune în fiecare ciclu cu excepția cazului când saltul este luat sau când are loc o interblocare. O interblocare previne instrucțiunile de a fi executate într-o secvență greșită pentru a conserva dependențele originale de date. Presupuneți că regiștrii pot bypass-ați atât de la ieșirea IDE cât și de la ieșirea MEM. Presupuneți de asemenea că $r2$ nu va mai fi folosit după întoarcerea din execuție.

A	ADD	$r3 \leftarrow 0$
B	LOAD	$r1 \leftarrow M(0200)$
C	Loop: LOAD	$r2 \leftarrow M(0208+r1)$
D	ADD	$r3 \leftarrow r2 + r3$
E	SUB	$r1 \leftarrow r1 - 4$
F	BRANCH!	$PC \leftarrow \text{Loop}$ if $r1 \neq 0$
G	NOP	a delay slot to be filled
H	STORE	$M(0204) \leftarrow r3$
I	RETURN	return from function
J	NOP	a delay slot to be filled

- a) Există o interblocare în program? Dacă da, realizați reordonarea codului din program și prezentați noul program fără cicli de interblocare.
- b) Calculați numărul total de cicli necesari pentru a executa toate instrucțiunile înainte și după ce s-a eliminat cicli de interblocare.
- c) Filați sloturile de întârziere. Descrieți reordonarea codului și/sau duplicatele realizate. Prezentați același program după filarea sloturilor de întârziere. Reamintim că RETURN este tot o instrucțiune de ramificare, Utilizați un apostrof (') pentru a marca noua copie a instrucțiunii duplicate. De exemplu dacă ați duplicat D, numele copiei este D'.
- d) Calculați numărul total de cicli necesari pentru a se executa toate instrucțiunile după ce ați filat sloturile de întârziere.

13. Utilizând o optimizare ILP, dublați performanțele următoarei bucle, sau explicați de ce nu este posibil acest lucru. Procesorul poate realiza doar un salt per ciclu, dar altfel are resurse infinite.

r1 = ... ; r1 is head pointer to a linked list

r3 = 0

LOOP:

r2 = M[r1 + 8]

r3 = r3 + r2

r1 = M[r1]

branch r1 != 0, LOOP

... = r3 ; r3 este utilizat când bucla este completă

14. Considerați calea de date de dedesubt. Acest procesor nu suportă cod cu sloturi de întârziere pentru ramificații (prezice că saltul nu se realizează cu o penalizare de 1 ciclu). Pentru fiecare semnal de control listat în tabelul de mai jos, determinați valoarea acestuia în cicli 3 până la 9, inclusiv. De asemenea, arătați instrucțiunile care ocupă fiecare etaj din pipeline în toți cicli (presupuneți că linia de validare a scrierii în IF/ID este setată cu inversul semnalului Stall).

Starea inițială a procesorului este:

PC = 0

Toți regiștrii pipeline conțin 0

Toți regiștrii din fișierul de regiștrii conțin 0

Memoria de date conține 0 în toate locațiile

Memoria de instrucțiuni conține:

00: addiu \$3, \$zero, 4

04: lw \$4, 100(\$3)

08: addu \$2, \$4, \$3

0C: beq \$4, \$zero, 0x14

10: addiu \$3, \$3, 1

14: addu \$2, \$2, \$3

Toate celelalte locații conțin 0.

Utilizați avansarea datelor unde este posibil. Toate intrările multiplexoarelor sunt numerotate vertical de sus în jos, începând cu 0 dacă priviți la calea de date orientată pe orizontală. De asemenea valorile pentru ALUOp sunt:

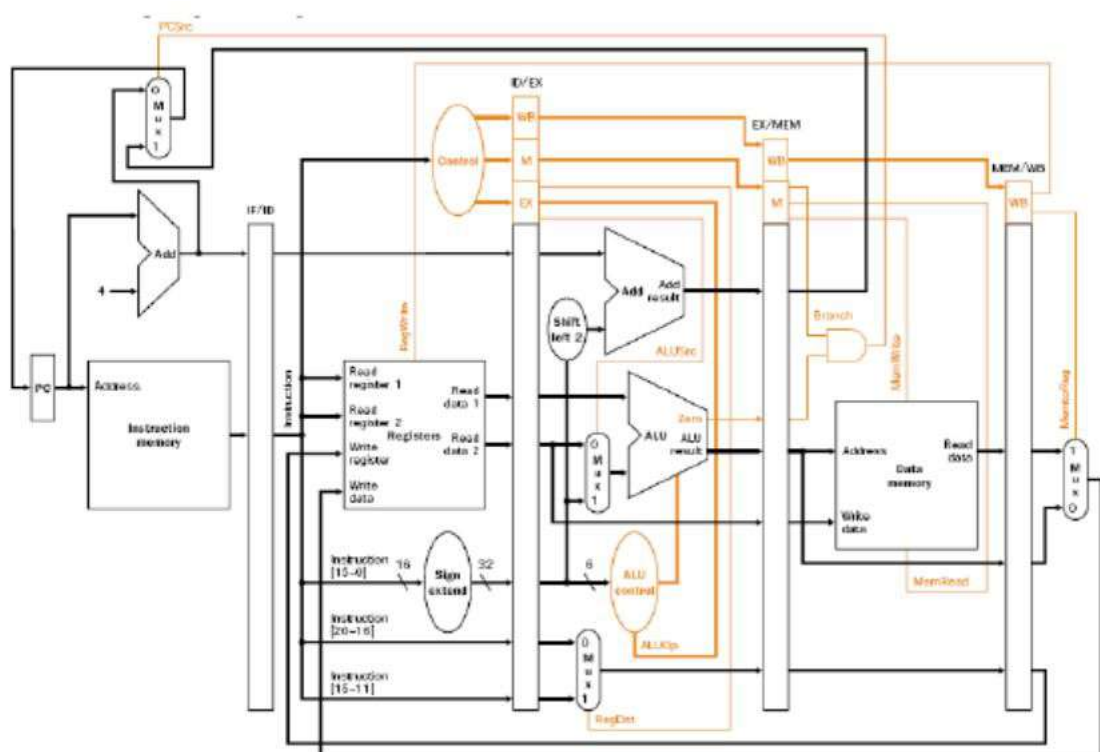
Value	Desired ALU Action
00	Add
01	Subtract
10	Determine by decoding funct field

Instruction formats:

	6 bits	5 bits	5 bits	5 bits	5 bits	6 bits
R:	op	rs	rt	rd	shamt	funct
I:	op	rs	rt	address / immediate		
J:	op	target address				

Time	PCWrite	IF.Flush	Branch	Stall	ALUSrc	ALUOp	RegDst	ForwardA	ForwardB	MemRead	MemWrite	MemtoReg	RegWrite
	IF		ID		EX					M		WB	
0	00: addiu		--		-					--		-	
	1	0	0	0	0	00	0	00	00	0	0	0	0
1	04: lw		00: addiu		-					--		-	
	1	0	0	0	0	00	0	00	00	0	0	0	0
2	08: addu		04: lw		00: addiu					--		-	
	1	0	0	0	1	00	0	00	00	0	0	0	0
3													
4													
5													
6													
7													
8													
9													

15. Considerați următoarea cale de date:



a) Discutați funcționalitatea semnalelor de control RegDst și ALUSrc.

b) Modificați, diagrama pentru a indica schimbările căii de date (și orice multiplexare adițională) necesară pentru a furniza bypass-ul din EX în EX pentru toate hazardurile RAW posibile pentru operațiile aritmetice. Cum se schimbă ALUSrc când se adaugă bypass-ul?

16. Imaginați o instrucțiune a cărei funcție este de a citi din memorie patru cuvinte adiacente pe 32 de biți în patru regiștrii arhitecturali pe 32 de biți specificați. Presupuneți că pipeline-ul cu cinci etaje este filat cu aceste instrucțiuni și NUMAI cu aceste instrucțiuni, care este numărul minim de citiri și scrieri care pot fi cerute la porturile fișierului de regiștrii?

17. Pipelining și bypass. În această întrebare se va explora modul cum bypass-ingul afectează performanțele de execuție a programului. Pentru început se consideră un MIPS clasic cu 5 etaje. Pentru referință luați în considerare figura care urmează. Pentru această întrebare, se va utiliza următorul cod pentru a evalua performanțele pipeline-ului.

```
1 add      $t2, $s1, $sp
2 lw       $t1, $t1, 0
3 addi     $t2, $t1, 7
4 add      $t1, $s2, $sp
5 lw       $t1, $t1, 0
6 addi     $t1, $t1, 9
7 sub      $t1, $t1, $t2
```

Care este întârzierea dată de utilizarea lui load (load-use) pentru MIPS-ul standard cu 5 etaje?

Din nou, utilizând pipeline-ul MIPS standard, identificați unde valoarea pentru fiecare operand din registru vine din bypass sau din fișierul de regiștrii. Pentru claritate, scrieți REG sau BYPASS în fiecare celulă.

Instruction	Src Operand 1	Src Operand 2
1		
2		N/A
3		N/A
4		
5		N/A
6		N/A
7		

Câți cicli va lua programul pentru a se executa pe un pipeline MIPS standard?

Presupuneți, că datorită constrângerilor de circuit, că firele pentru bypass de la etajul de memorie înapoi către etajul de execuție sunt omise din pipeline. Care întârzierea la utilizarea lui load pentru acest pipeline modificat?

Identificați dacă valoarea pentru fiecare operand din registru vine din bypass sau din fișierul de regiștrii pentru pipeline-ul modificat. Pentru claritate scrieți REG sau BYPASS în fiecare celulă.

Instruction	Src Operand 1	Src Operand 2
1		
2		N/A
3		N/A
4		
5		N/A
6		N/A
7		

Cât timp ia execuția programului pe pipeline-ul modificat?

18. Fie pseudo-codul în assembler:

```

1: mov r1 = 0 ;; iteration count i
2: mov r2 = 0 ;; initialize x
loop:
3: sll r3 = r1, 2 ;; r3 = r1*4
4: lw r4 = A(r3) ;; load A[i]
5: add r2 = r2, r4 ;; x += A[i]
6: add r1 = r1, 1 ;; increment counter
7: bne r1, 1024, loop ;; backwards branch

```

Presupuneți că aveți un procesor cu resurse infinite (lărgime infinită, fereastră infinită pentru instrucțiuni, un număr infinit de unități funcționale, resurse infinite de redenumire, un număr infinit de încărcări întârziate/suprapuse), și care este pe deplin pipeline, astfel încât toate instrucțiunile normale iau un ciclu, dar încărcările care lipsesc iau 10 cicli. Presupuneți o predicție perfectă a salturilor.

- a) Prezentați o planificare dinamică a instrucțiunilor pentru primele două iterații ale buclei în tabelul ce urmează, presupunând că lipsa încărcărilor (load) în cache apare iterațiile pare ($r1 = 0, 2, \dots$) și potrivirea încărcărilor (load) în iterațiile impare ($r1 = 1, 3, \dots$). Fiți siguri că ați indicat iterația buclei la care este asociată fiecare instrucțiune. De exemplu, dacă instrucțiunea 5 din iterația 7 poate rula în ciclul 3, atunci se poate scrie "5.7" în una din celulele din rândul ciclului 3. DE notat că primul ciclu a fost filat și celulele din table sunt mai multe decât aveți nevoie.

Cycle	Instructions					
1	1	2				
2						
3						
4						
5						
6						
7						
8						
9						
10						
11						
12						
13						
14						
15						
16						
17						
18						

19. Predicție și predicție.

- a) Considerați următoarea secvență de cod:

```
for (i=0; i < 1000000; i++){  
    a= random(100);  
    if (a >= 50){  
        ....  
    }  
}
```

Presupuneți că `random(N)` returnează un număr aleatoriu uniform distribuit între 0 și $N - 1$ inclusiv. Considerați instrucțiunea de ramificație asociată cu linia `if`. Există acolo un predictor de ramificație care prezice corect? Explicați răspunsul.

- b) Predicarea poate elimina toate ramificațiile condiționale de avansare în program. Salturile înapoi în program sunt, în mod tipic, asociate cu buclele și ca urmare sunt în majoritate luate (se sare la adresa destinație). Ca urmare este posibil de a elimina salturile de avansare și de a prezice static salturile înapoi ca luate și aceasta va elimina complexitatea predictorilor într-un procesor. Dar, procesorul Itanium care are suport hardware pentru predicție totuși păstrează un predictor complex cu două niveluri pentru ramificație. Explicați aceasta în acest caz.

20. Ocazional compilatoarele sunt utile. O buclă comună în programele științifice este așa numita SAXPY:

```
for (i = 0; i < N; i++)  
{  
    Y[i] = A*X[i] + Y[i];  
}
```

- a) Pentru $A = 3$ și $N=300$ converțiți codul C în asamblare MIPS. Puteți presupune că X și Y sunt arii cu întregi pe 32 de biți. Încercați să evitați pseudo-codul datorită părții a doua a acestei întrebări. De asemenea, evitați să introduceți dependențe false. Utilizați vă rog `$s0` pentru a păstra valoarea lui i . Puteți, de asemenea, presupune că `$s2` și `$s3` păstrează inițial adresa de start a lui X , respectiv Y . Sugestie: deoarece $A=3=2+1$ puteți realiza ușor înmulțirea fără a utiliza instrucțiunea mult.
- b) Cât de multe instrucțiuni statice (numărul total de instrucțiuni în memorie) are fragmentul de program? Câte instrucțiuni dinamice (numărul total de instrucțiuni executate) are fragmentul de program?
- c) Presupuneți că dorim să realizăm mai mult de o buclă SAXPY:

```
for (i = 0; i < N; i++)  
{  
    Y[i] = A*X[i] + Y[i];  
}  
for (i = 0; i < N; i++)  
{  
    Z[i] = A*X[i] + Z[i];  
}
```

Utilizând calculele de la întrebările a și b, câte instrucțiuni dinamice și câte statice are acest fragment de program?

- d) Există două optimizări care pot fi realizate cu acest cod. Prima buclă este fuzionată, astfel cele două bucle sunt combinate într-o singură buclă:

```
for (i = 0; i < N; i++)  
{  
  Y[i] = A*X[i] + Y[i];  
  Z[i] = A*X[i] + Z[i];  
}
```

Altă optimizare este eliminarea subexpresiilor comune. Astfel putem scoate în afară calculul repetat:

```
for (i = 0; i < N; i++)  
{  
  T = A*X[i];  
  Y[i] = T + Y[i];  
  Z[i] = T + Z[i];  
}
```

Converteți fiecare din aceste fragmente de program în asamblare MIPS. \$s4 va conține adresa de start pentru Z, Utilizați \$s1 pentru a memora valoarea lui T. Câte instrucțiuni dinamice și câte statice are fiecare fragment? Cât de mult s-a salvat cu fuziunea buclilor. Cât de mult s-a salvat cu eliminarea sub-expresiei?

- e) Adnotați codul în asamblare de la punctul d) cu arcuri pentru dependențele de date (flux, anti, și ieșire – fără control

21. Considerați următoarea instrucțiune:

Instrucțiune: and Rd, Rs, Rt

Interpretare: $\text{Reg}[\text{Rd}] = \text{Reg}[\text{Rs}] \text{ AND } \text{Reg}[\text{Rt}]$

- a) Care sunt valorile semnalelor de control generate de controlul din figura 4.2 pentru instrucțiunea următoare `sw rt, offset(rs)`?

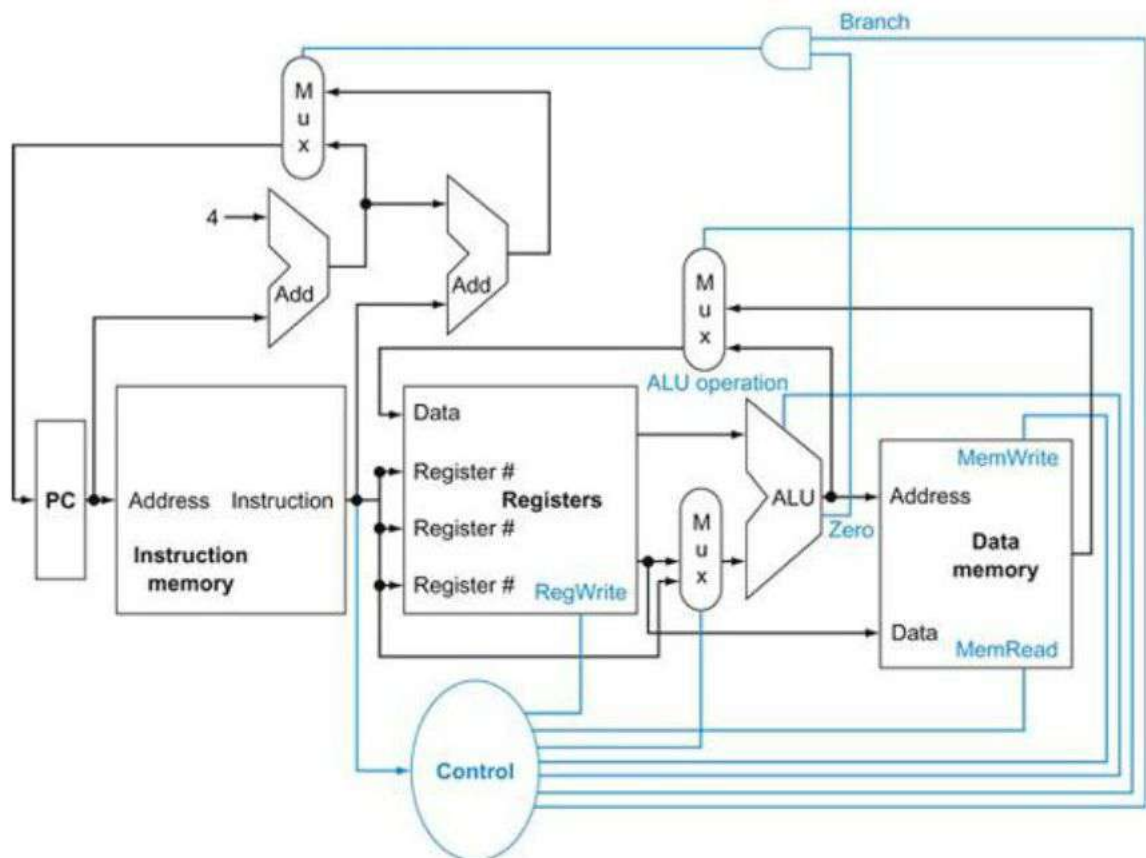


Figura 4.2

- b) Ce resurse (blocuri) realizează o funcție utilă pentru această instrucțiune?
- c) Care resurse (blocuri) produc ieșiri dar ieșirile lor nu sunt folosite pentru această instrucțiune. Care resurse nu produc ieșiri pentru această instrucțiune?

22. Implementarea MIPS de bază din figura 4.2 poate să implementeze doar câteva instrucțiuni. Se pot adăuga noi instrucțiuni la o arhitectură a setului de instrucțiuni (ISA), dar decizia dacă se face sau nu această adăugare depinde, printre alte lucruri, de costul și complexitatea pe care adăugarea propusă o introduce în calea de date și control. Fie instrucțiunea: LWI Rt, Rd(Rs). Interpretare: $\text{Reg}[Rt] = \text{Mem}[\text{Reg}[Rd] + \text{Reg}[Rs]]$.

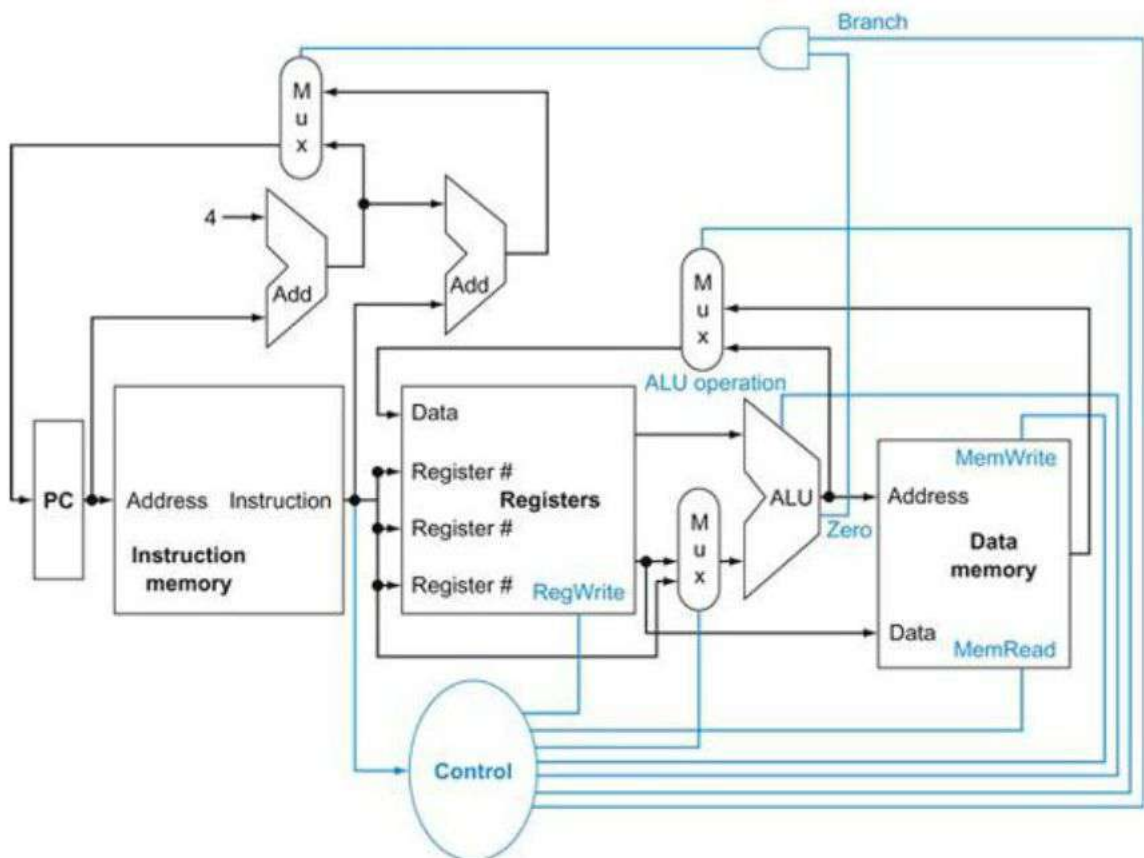


Figura 4.2

- Care dintre blocurile existente (dacă există unul) poate fi utilizat pentru această instrucțiune?
- Ce blocuri funcționale noi (dacă există unul) sunt necesare pentru această instrucțiune?
- Ce semnale noi (dacă există unul) sunt necesare pentru unitatea de control pentru a suporta această instrucțiune?

23. Atunci când un proiectant de procesoare consideră o posibilă îmbunătățire a căii de date, decizia depinde în mod uzual de raportul cost/performanță. Presupunem că se începe cu calea de date din figura 4.2 unde I-Mem, Add, Mux, ALU, Regs, D-Mem și blocurile de control introduc întârzieri de 400 ps, 100ps, 30ps, 120ps, 200ps, 350ps și respectiv 100ps și un cost de 1000, 30, 10, 100, 200, 2000 și respectiv 500. Considerați adăugarea unui multiplicator la ALU. Această adăugare va aduce o întârziere de 300 ps la ALU și un cost de 600 tot la ALU. Rezultatul va fi 5% mai puține instrucțiuni executate deoarece nu se va mai emula instrucțiune MUL de înmulțire.

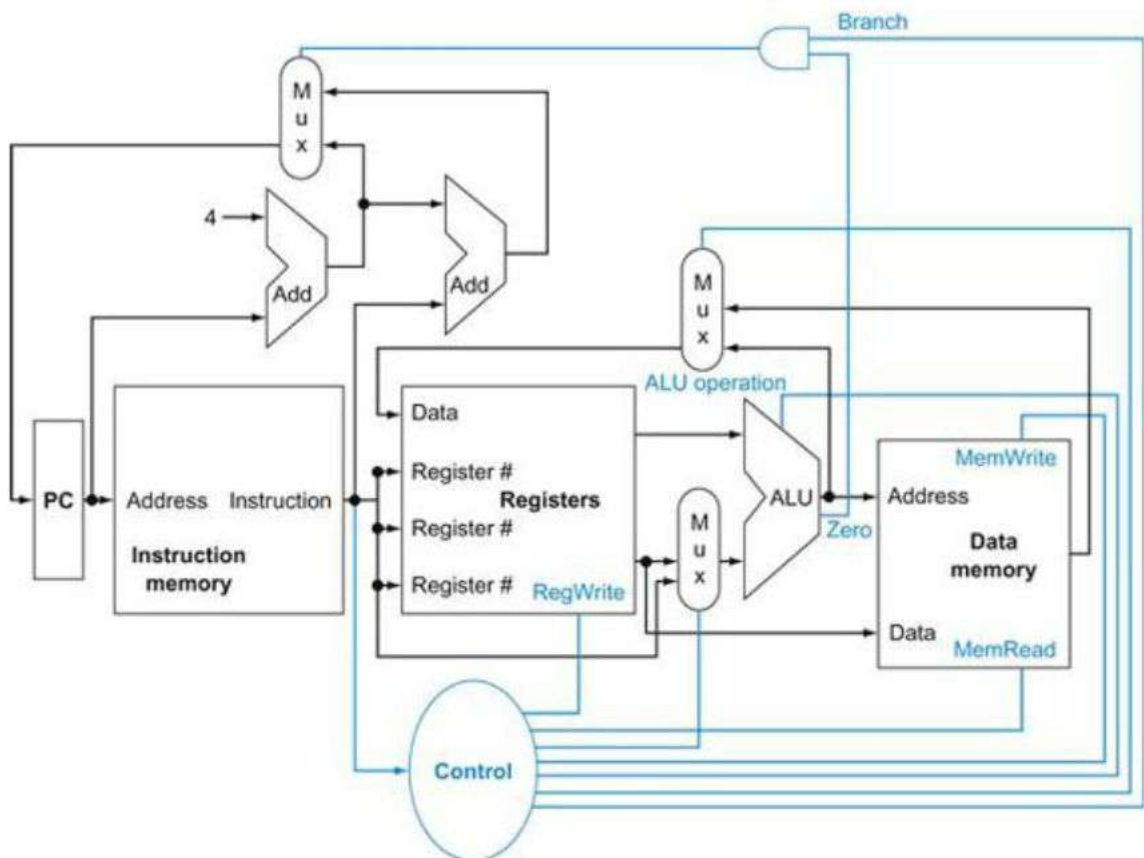


Figura 4.2

- Care este durata unui ciclu de ceas cu și fără această îmbunătățire?
- Care este creșterea de viteză prin adăugarea acestei îmbunătățiri?
- Comparați raportul cost/performanță cu și fără această îmbunătățire.

24. Să presupunem că blocurile logice necesare pentru implementarea unei căi de date au următoarele întârzieri:

I-Mem	Add	Mux	ALU	Regs	D-Mem	Sign-Extend	Shift-Left-2
200ps	70ps	20ps	90ps	90ps	250ps	15ps	10ps

- Dacă singurul lucru pe care trebuie să-l facem în procesor este extragerea de instrucțiuni consecutive (figura 4.6), care va fi durata unui ciclu?

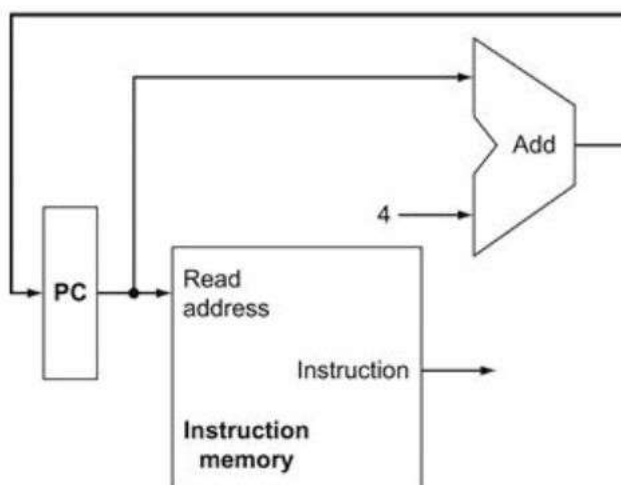


Figura 4.6

- b) Considerați o cale de date similară cu aceea din figura 4.11, dar pentru un procesor care are numai un singur tip de instrucțiune: ramificație necondițională relativă la PC. Care va fi durata ciclului pentru această cale de date?

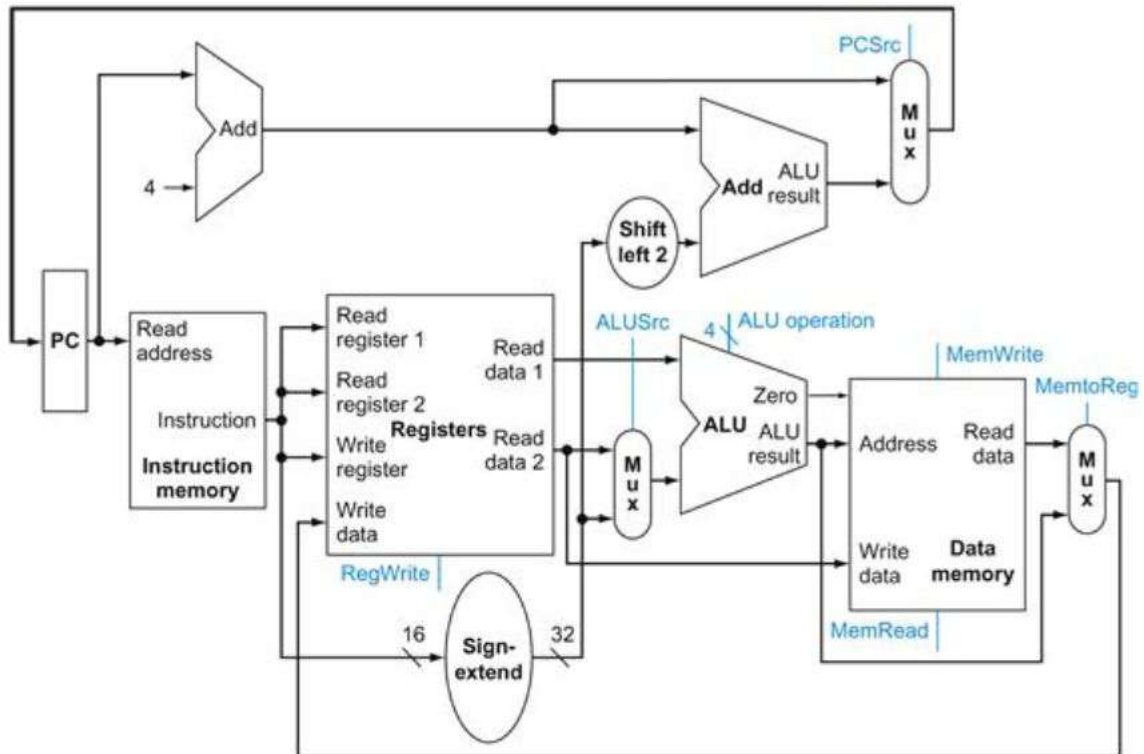


Figura 4.11

- c) Repetați punctul anterior (b) dar de data aceasta trebuie suportată doar ramificație condiționată relativă la PC.
- d) Fie deplasarea la stânga cu 2. Care tipuri de instrucțiuni necesită această resursă?
- e) Fie deplasarea la stânga cu 2. Pentru care tipuri de instrucțiuni (dacă există unul) este această resursă pe calea critică?
- f) Fie deplasarea la stânga cu 2. Presupunând că sunt suportate doar instrucțiunile beq și add, discutați cum schimbările în întârzierile acestei resurse afectează durata unui ciclu al procesorului. Presupuneți că întârzierile altor resurse nu se schimbă.

25. Presupunem că nu există încetiniri (stalls) ale pipeline-ului și că defalcarea instrucțiunilor dată de execuția instrucțiunilor este:

add	addi	not	beq	lw	sw
20%	20%	0%	25%	25%	10%

- a) În câte fracțiuni din toți cicli sunt utilizate datele din memorie?
- b) În câte fracțiuni din toți cicli este utilizat circuitul pentru extensia de semn? Ce face circuitul atunci când intrările sale nu sunt necesare?

26. Atunci când se fabrică cipurile pe siliciu, defectele din materiale (siliciu) și erorile de fabricație pot genera circuite defecte. Un defect foarte des întâlnit este ca un fir să afecteze semnalul din alt fir. Acest defect este denumit convorbire încrucișată (cross-talk). O clasă specială de eroare cross-talk este atunci când un semnal este conectat la un fir care are o valoare logică constantă (de exemplu un fir de alimentare). În acest caz avem o eroare de tip blocat pe 0 sau blocat pe 1, și semnalul afectat are întotdeauna valoarea m0 respectiv 1. Se va lua în considerare în continuare bitul 0 de la intrarea Write Register din fișierul de regiștrii din figura 4.24.

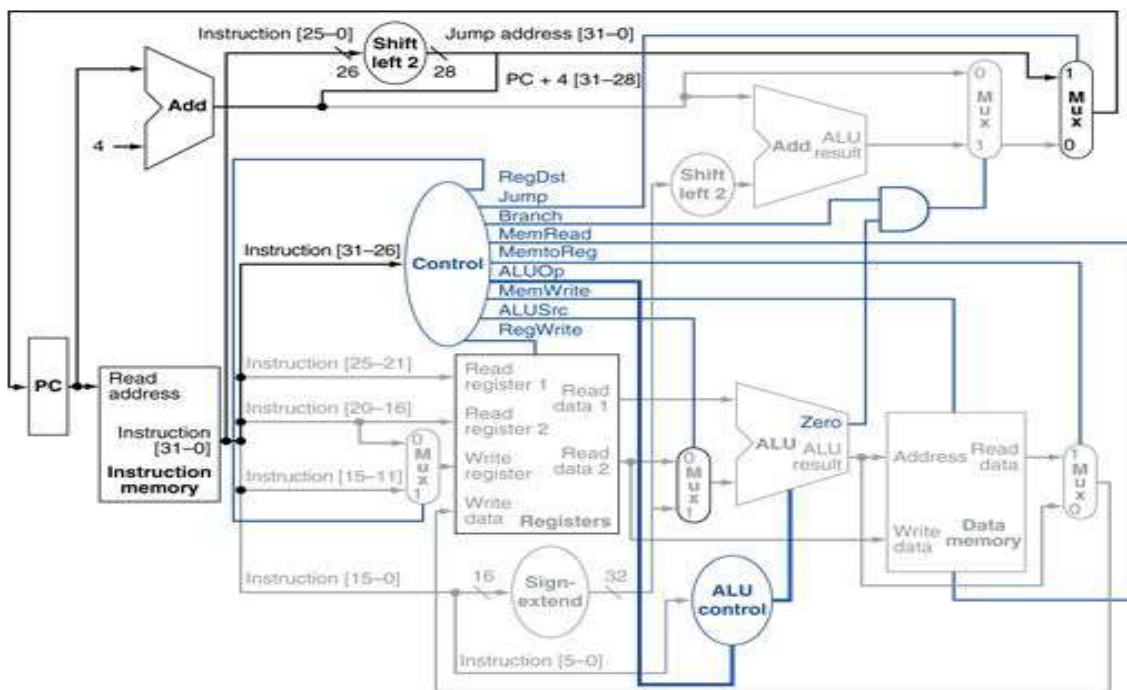


Figura 4.24

- Presupuneți că testarea procesorului este realizată prin umplerea PC-ului, regiștrilor, și a memoriilor de date și instrucțiuni cu unele valori (puteți alege care valori), se permite execuția unei singure instrucțiuni, după care se citesc PC, memoriile și regiștrii. Aceste valori sunt apoi examinate pentru a determina dacă apare o anumită eroare. Puteți proiecta un test (valorile pentru PC, memorii și regiștrii) care va determina o eroare de tip blocare pe 0 pe acest semnal?
- Repetăți punctul anterior pentru eroare de tipul blocat pe 1. Puteți utiliza un singur test atât pentru blocat pe 0 cât și pentru blocat pe 1? Dacă da, explicați cum; dacă nu explicați de ce nu.
- Dacă știm că procesorul are un defect de tipul blocat pe 1 pe acest semnal, mai este procesorul utilizabil? Pentru a fi utilizabil, trebuie să fim capabil să convertim orice program care se execută pe un procesor MIPS normal într-un program care să ruleze pe acest procesor. Puteți presupune că există suficiente memorie de instrucțiuni și date liberă care să vă permită să faceți programul mai lung și să memorați date adiționale. Sugestie: procesorul este utilizabil dacă fiecare instrucțiune afectată de acest defect poate fi înlocuită cu o secvență de instrucțiuni care au același efect.

- d) Repetă punctul a, dar acum eroarea de testat este aceea că semnalul MemRead devine 0 dacă semnalul de control RegDst devine 0. Altfel nu este nici un defect.
- e) Repetă punctul d, dar acum eroarea de testat este aceea că semnalul Jump devine 0 dacă semnalul de control RegDst devine 0. Altfel nu este nici un defect.

27. În acest exercițiu vom examina în detaliu cum se execută o instrucțiune într-o cale de date cu un singur ciclu. Problema se referă la ciclul de ceas în care procesorul extrage următorul cuvânt instrucțiune:

1010110001100010000000000010100

Presupuneți că memoria de date este toată filată cu 0 și că regiștrii procesor au următoarele valori la începutul ciclului când este extrasă instrucțiunea anterioară:

r0	r1	r2	r3	r4	r5	r6	r8	r12	r31
0	-1	2	-3	-4	10	6	8	2	-16

- a) Care sunt ieșirile unităților de extensie de semn și de deplasare la stânga cu 2 pentru jump (vezi figura 4.24 în partea de sus) pentru acest cuvânt instrucțiune?
- b) Care sunt valorile intrărilor unității de control pentru ALU pentru această instrucțiune?
- c) Care este noua adresă din PC după ce este executată această instrucțiune? Accentuați calea pe care este determinată această valoare.
- d) Pentru fiecare Muz arătați valorile ieșirilor de date pe durata execuției acestei instrucțiuni și a acestor valori din regiștrii.
- e) Pentru ALU și cele două unități de adunare, care sunt valorile datelor de intrare?
- f) Care sunt valorile pentru toate intrările unității "Registers"?

28. În acest exercițiu, se va examina cum afectează pipeline-ul durata ciclului de ceas a procesorului. Se presupune că etajele individuale ale căii de date au următoarele întârzieri:

IF	ID	EX	MEM	WB
250ps	350ps	150ps	300ps	200ps

De asemenea, presupuneți că instrucțiunile executate de procesor sunt defalcate după cum urmează:

alu	beq	lw	sw
45%	20%	20%	15%

- a) Care este durata ciclului de ceas într-un procesor pipelined și în unul fără pipeline?
- b) Care este întârzierea totală a unei instrucțiuni LW într-un procesor pipelined și în unul fără pipeline?
- c) Dacă putem împărți un etaj din calea de date pipeline în două noi etaje, fiecare cu jumătate de întârziere față de etajul original, care etaj l-ați înjumătății și care este noua durată a ciclului de ceas a al procesorului?
- d) Presupunând că nu există hazard sau încetiniri (stalls) care este utilizarea memoriei de date?

- e) Presupunând că nu există hazard sau încetiniri (stalls), care este utilizarea portului write-register din unitatea "Registers"?

29. În acest exercițiu, se va examina cum afectează dependențele de date execuția într-un pipeline de bază cu 5 etaje. Întrebările se referă la următoarea secvență de instrucțiuni:

```
or    r1,r2,r3
or    r2,r1,r4
or    r1,r1,r2
```

De asemenea, presupuneți următorii timpi pentru ciclul de ceas pentru fiecare opțiune legată de avansare (forwarding):

Without Forwarding	With Full Forwarding	With ALU-ALU Forwarding Only
250 ps	300 ps	290 ps

- Indicați dependențele și tipul lor.
- Presupuneți că nu există avansare în acest procesor. Indicați hazardul și adăugați instrucțiuni nop necesare pentru a-l elimina.
- Presupuneți că există avansare completă. Indicați hazardul și adăugați instrucțiuni nop necesare pentru a-l elimina.
- Care este timpul total de execuție pentru această secvență de instrucțiuni fără avansare (forwarding) și cu avansare completă? Care este creșterea de viteză obținută prin adăugarea unei avansări complete la o pipeline care nu are avansare?
- Adăugați instrucțiuni nop la acest cod pentru a elimina hazardul dacă există numai o avansare ALU-ALU (nu există avansare de la etajul MEM la etajul EX).
- Care este timpul total de execuție a acestui secvențe de instrucțiuni cu avansare doar de tipul ALU-ALU? Care este creșterea de viteză față de cazul fără avansare?

30. În acest exercițiu, se va examina cum hazardul de resurse, cel de control și proiectarea arhitecturii setului de instrucțiuni (ISA) poate afecta execuția pipeline. Întrebările din această problemă se referă la următoarea secvență de cod MIPS:

```
sw    r16, 12(r6)
lw    r16, 8(r6)
beq   r5, r4, label  # se presupune că r5!=r4
add   r5, r1, r4
slt   r5, r15, r4
```

Presupuneți că etajele pipeline individuale au următoarele întârzieri:

IF	ID	EX	MEM	WB
200ps	120ps	150ps	190ps	100ps

- Pentru această întrebare, presupuneți că toate ramificațiile sunt perfect predictibile (acest fapt elimină hazardul de control) și că nu se utilizează nici un slot de întârziere. Dacă avem o singură memorie (atât pentru instrucțiuni

și date), există un hazard structural de fiecare dată când dorim să extragem o instrucțiune în același ciclu cu accesul la date al altei instrucțiuni. Pentru a garanta progresul avansării, acest hazard trebuie întotdeauna rezolvat în favoarea instrucțiunii care realizează accesul la date. Care este timpul total de execuție a acestei secvențe de instrucțiuni într-un pipeline cu 5 etaje care are o singură memorie? S-a văzut că hazardul de date poate fi eliminat prin adăugarea de cod de tip nop. Se poate realiza la fel și cu hazardul structural? De ce?

- b) Pentru această întrebare, presupuneți că toate ramificațiile sunt perfect predictibile (acest fapt elimină hazardul de control) și că nu se utilizează nici un slot de întârziere. Dacă schimbăm instrucțiunile load/store astfel încât să utilizeze un registru (fără offset) ca adresă, aceste instrucțiuni nu mai au nevoie de ALU. Ca rezultat, etajele MEM și EX pot fi suprapuse și pipeline-ul are numai 4 etaje. Schimbați codul pentru acomodarea la această schimbare a ISA. Presupuneți că această schimbare nu afectează durata ciclului de ceas, care este creșterea de viteză în secvența de instrucțiuni?
- c) Presupuneți că avem încetiniri la ramificații (stall-on-branch) și nu avem sloturi de întârziere, care este creșterea de viteză a acestui cod dacă ieșirile ramificațiilor sunt determinate în etajul ID, relativ la execuția unde ieșirile ramificațiilor sunt determinate în etajul EX?
- d) Fiind date aceste întârzieri în etajele pipeline, repetați creșterea de viteză calculată la punctul b, dar luați în calcul (posibila) modificarea duratei ciclului de ceas. Când etajele MEM și EX sunt executate într-un singur etaj, majoritatea execuției lor poate fi realizată în paralel. Ca rezultat, noul etaj EX/MEM are o întârziere care este mai mare ca cea originală a celor două etaje, plus 20ps necesare pentru execuție care nu poate fi realizată în paralel.
- e) Fiind date aceste întârzieri în etajele pipeline, repetați creșterea de viteză calculată la punctul c, dar luați în calcul (posibila) modificarea duratei ciclului de ceas. Presupuneți că întârzierea etajului ID crește cu 50% și întârzierea etajului EX descrește cu 10ps când ieșirea ramificației este mutată din EX în ID.
- f) Presupuneți că avem încetiniri la ramificații (stall-on-branch) și nu avem sloturi de întârziere, care este noua durată a ciclului de ceas și timpul de execuție a acestei secvențe de instrucțiuni dacă se calculează adresa pentru beq în etajul MEM? Care este creșterea de viteză pentru această modificare? Presupuneți că întârzierea pentru etajul EX se reduce cu 20ps și întârzierea pentru etajul MEM este neschimbată atunci când ieșirea pentru ramificație este mutată din EX în MEM.

31. Presupuneți următoarea buclă:

Loop:

lw r1,0(r1)

and r1,r1,r2

lw r1,0(r1)

lw r1,0(r1)

beq r1,r0,loop

Presupuneți că se utilizează o predicție perfectă a salturilor (nici o încetinire datorită hazardului de control), că nu există sloturi de întârziere, și pipeline-ul are suport complet pentru avansare. Presupuneți că sunt executate multe iterații ale acestei bucle înainte de ieșire.

- a) Prezentați diagrama de execuție pentru a treia iterație a acestei bucle, de la ciclul în care se va extrage prima instrucțiune a acestei iterații până la (fără a include) ciclul în care se va extrage prima instrucțiune a următoarei iterații. Prezentați toate instrucțiunile care sunt în pipeline pe durata acestor cicli (nu doar cei din iterația a treia).
- b) Cât de des (ca procent din toți cicli) există un ciclu în care toate cele cinci etaje pipeline realizează o muncă (execuție) utilă?

32. Acest exercițiu încearcă să vă ajute să înțelegeți costul/ complexitatea/ performanțele avansării (forwarding) într-un procesor pipelined. Exercițiul se referă la calea de date din figura 4.5. Întrebările din problemă presupun că, toate instrucțiunile au un tip particular de dependențe de date de tip RAW (Read After Write – citire după scriere). Tipul de dependență de date RAW este identificat prin etajul care produce rezultatul (EX sau MEM) și instrucțiunea care consumă rezultatul (prima instrucțiune care urmează unei care produce rezultatul, sau a doua instrucțiune, sau ambele). Presupunem că scrierea în registru este realizată în prima parte a ciclului mașină în timp ce citirea se realizează în a doua jumătate a ciclului mașină, astfel încât dependențele de date de tipul EX cu a treia instrucțiune care urmează și MEM cu a treia instrucțiune care urmează nu sunt luate în calcul deoarece nu conduc la hazard de date. De asemenea presupuneți că CPI-ul procesorului este 1 dacă nu există nici un hazard de date.

EX to 1 st Only	MEM to 1 st Only	EX to 2 nd Only	MEM to 2 nd Only	EX to 1 st and MEM to 2 nd	Other RAW Dependences
5%	20%	5%	10%	10%	10%

Presupuneți următoarele întârzieri pentru etajele individuale ale pipeline-ului. Pentru etajul EX, întârzierile sunt date separat pentru un procesor fără avansare și pentru un procesor cu diferite tipuri de avansare.

IF	ID	EX (no FW)	EX (full FW)	EX (FW from EX/MEM only)	EX (FW from MEM/ WB only)	MEM	WB
150 ps	100 ps	120 ps	150 ps	140 ps	130 ps	120 ps	100 ps

add \$14, \$5, \$6	lw \$13, 24(\$1)	add \$12, \$3, \$4	sub \$11, \$2, \$3	lw \$10, 20(\$1)
Instruction fetch	Instruction decode	Execution	Memory	Write-back

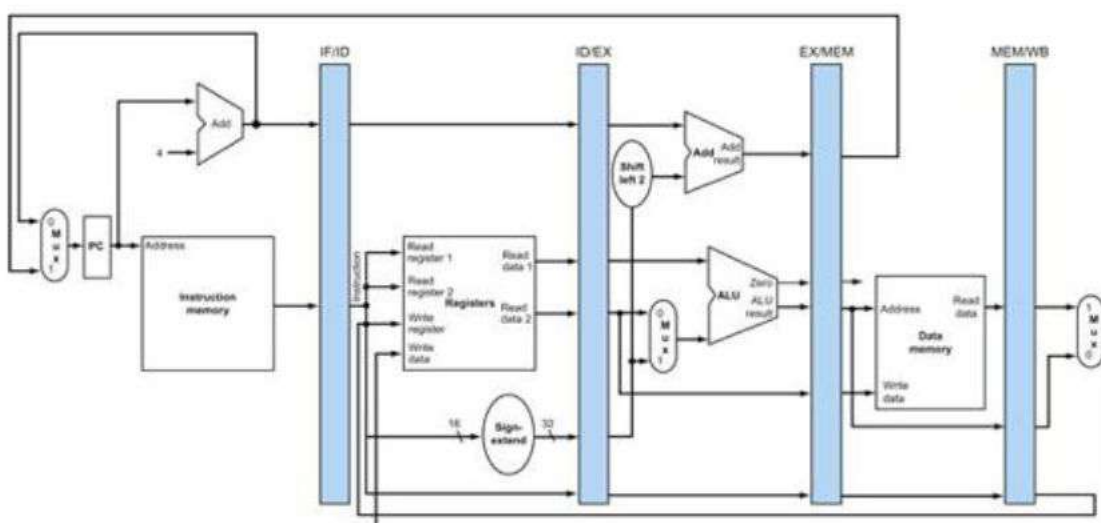


Figura 4.45

- a) Dacă nu se utilizează avansarea, ce fracțiune din ciclul sunt încetiniți datorită hazardului?
- b) Dacă utilizăm avansarea completă (se avansează toate rezultatele care pot fi avansate), ce fracțiune din ciclul sunt încetiniți datorită hazardului?
- c) Să presupunem că nu putem avea multiplexoare cu trei intrări necesare pentru avansarea completă. Trebuie să decidem dacă este mai bine să avansăm numai din registrul pipeline EX/MEM (avansare din ciclul următor), sau numai din registrul pipeline MEM/WB (avansare cu doi cicli). Care dintre aceste două soluții generează mai puține încetiniri pentru ciclul de date?
- d) Pentru o probabilitate dată a hazardului și întârzieri ale etajelor pipeline, care este câștigul în viteză dat de adăugarea la pipeline unei unități de avansare complete față de un pipeline fără avansare?
- e) Care va fi sporul adițional de viteză (relativ la un pipeline cu avansare) dacă se utilizează o avansare prin călătorie în timp (time-trevel forwarding) care elimină toate hazardurile? Presupuneți că acest circuit (neinventat încă) adaugă o întârziere de 100ps la unitatea de avansare completă pentru etajul EX.
- f) Repetați c) dar de data aceasta determinați care dintre cele două opțiuni conduce la un timp mai scurt pe instrucțiune.

33. Acest exercițiu încearcă să vă ajute să înțelegeți relațiile între avansare, detecția hazardului și proiectarea ISA (arhitectura setului de instrucțiuni). Întrebările din această problemă se referă la următoarea secvență de instrucțiuni, și presupuneți că se execută pe o cale de date pipeline cu 5 etaje:

```
add  r5,r2,r1
lw   r3,4(r5)
lw   r2,0(r2)
or   r3,r5,r3
sw   r3,0(r5)
```

- a) Dacă nu există avansare sau detecția hazardului, înserați nops pentru a asigura execuția corectă.
- b) Repetați punctul a) dar acum utilizați nops numai dacă hazardul nu poate fi evitat prin schimbarea sau rearanjarea acestor instrucțiuni. Puteți presupune că, în codul vostru modificat, registrul R7 poate fi utilizat pentru a păstra valori temporare.
- c) Dacă procesorul are avansare, dar am uitat să implementăm unitatea de detecție a hazardului, ce se întâmplă când se execută acest cod?
- d) Dacă există avansare, pentru primii cinci cicli din execuția acestui cod, specificați care semnale sunt activate în fiecare ciclu de detecția hazardului și unitatea de avansare din figura 4.60.
- e) Dacă nu există avansare, ce noi semnale de intrare și ieșire sunt necesare pentru unitatea de detecție a hazardului din figura 4.60? Utilizând această secvență de instrucțiuni ca un exemplu, explicați de ce este necesar fiecare semnal.
- f) Pentru noua unitate de detecție a hazardului de la punctul e) specificați care semnale de ieșire sunt activate în fiecare dintre primii cinci cicli procesor pe durata execuției acestui cod.

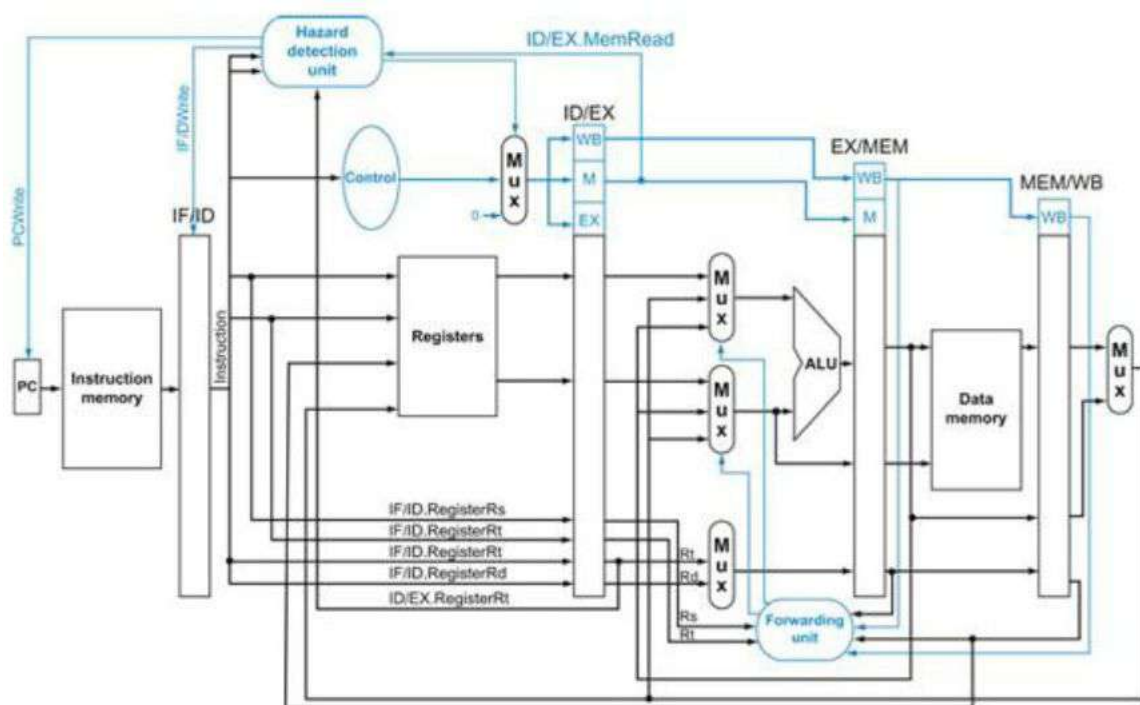


Figura 4.60

34. Acest exercițiu încearcă să vă ajute să înțelegeți relația între sloturile de întârziere, hazardul de control și execuția ramificațiilor într-un procesor cu pipeline. În acest exercițiu, presupunem că se execută următorul cod MIPS pe un procesor cu pipeline cu 5 etaje, avansare completă, și un predictor cu predicția conform căreia ramificația are loc (predict-taken).

```
lw    r2,0(r1)
```

```
label1:
```

```
beq   r2,r0,label2  #nu se ramifică (not taken) o dată, apoi se ramifică (taken)
```

```
lw    r3,0(r2)
```

```
beq   r3,r0,label1  # ramificația se execută (taken)
```

```
add   r1,r3,r1
```

```
label2:
```

```
sw    r1,0(r2)
```

- Desenați diagrama pipeline pentru acest cod, presupunând că nu există sloturi de întârziere și că ramificațiile se execută în etajul EX.
- Repetăți punctul a), dar presupuneți că sunt utilizate sloturile de întârziere. În codul dat, instrucțiunea care urmează ramificației este acum instrucțiunea din slotul de întârziere pentru ramificație.
- O cale de a muta rezoluția ramificației un etaj mai târziu este aceea de a nu avea nevoie de o operație ALU în ramificațiile condiționate. Instrucțiunea de ramificație va fi *bez rd, label* și *bnez rd, label*, și se va ramifica dacă registrul are sau respectiv nu are valoarea zero. Schimbați codul pentru a utiliza aceste instrucțiuni de ramificație în loc de *beq*. Puteți presupune că registrul R8 este disponibil pentru ca să-l utilizați ca registru temporar, și că poate fi utilizată o instrucțiune *seq* (set dacă este egal) de tip R.
- Utilizând prima instrucțiune de ramificație din codul dat ca un exemplu, descrieți logica de detecție a hazardului necesară pentru a suporta execuția ramificației în etajul ID ca în figura 4.62. Care tip de hazard este presupus a fi detectat cu această nouă logică?

- e) Pentru codul dat care este creșterea de viteză datorată mutării execuției ramificației în etajul ID? Explicați răspunsul. În calculul creșterii de viteză, presupuneți că această comparație adițională în etajul ID nu afectează durata ciclului de ceas.
- f) Utilizând prima instrucțiune de ramificație din codul dat ca un exemplu, descrieți suportul pentru avansare care trebuie adăugat pentru a realiza execuția ramificației în etajul ID. Comparați complexitatea acestei noi unități de avansare cu complexitatea unității de avansare existente în figura 4.62.

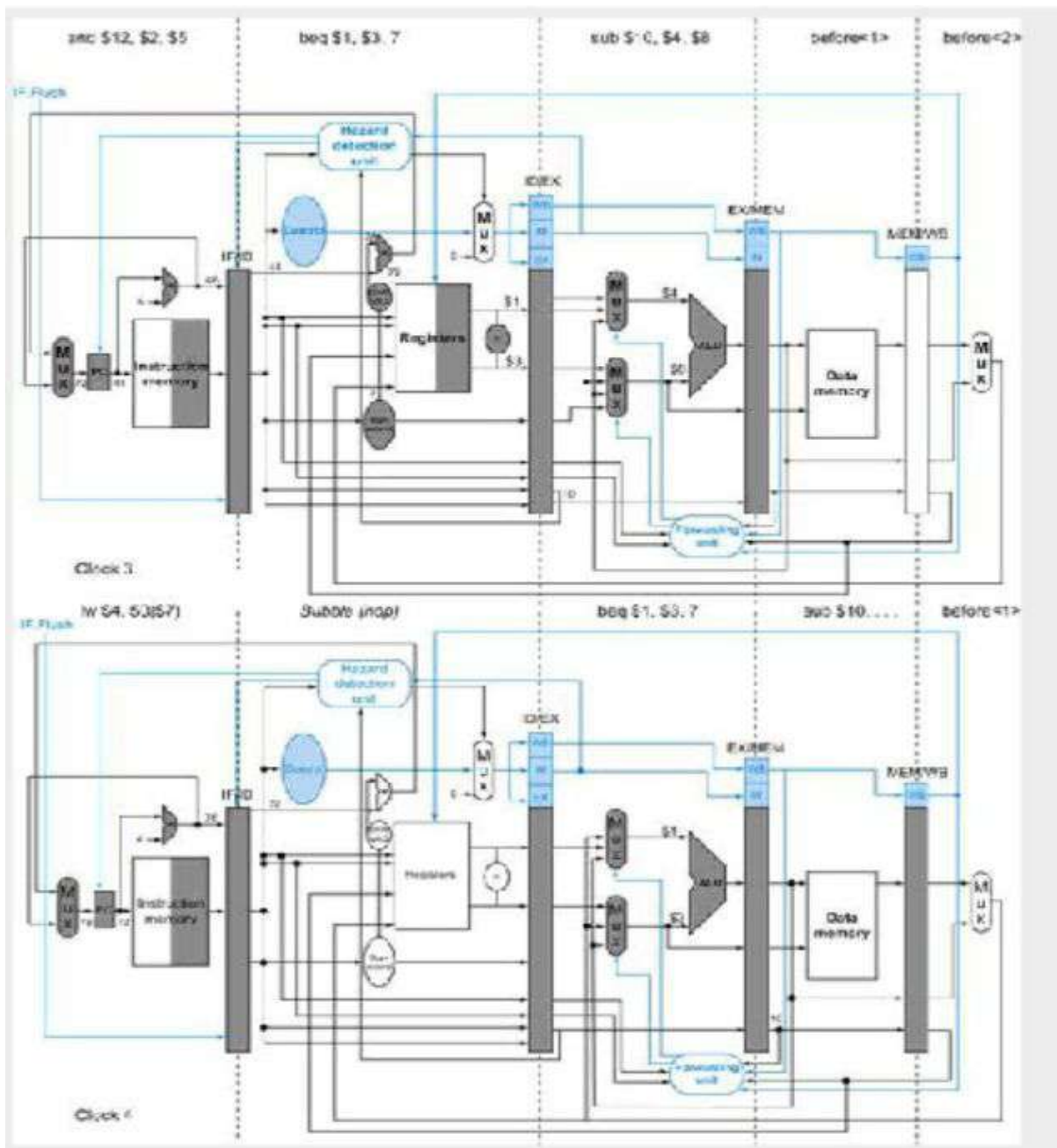


Figura 4.62

35. Importanța de a avea un predictor bun depinde adesea de modul cum sunt executate ramificațiile condiționate. Pe lângă acuratețea predicției ramificațiilor, aceasta va determina cât de mult timp se petrece cu încetirile datorită greșelilor de predicție. În acest exercițiu, presupuneți că eșecurile execuției dinamice a instrucțiunilor pe diferite categorii de instrucțiuni sunt următoarele:

R-Type	BEQ	JMP	LW	SW
40%	25%	5%	25%	5%

Presupuneți de asemenea următoarea acuratețe a predictorului pentru ramificații:

Always-Taken	Always-Not-Taken	2-Bit
45%	55%	85%

- Cicli de încetinire datorită greșelilor de predicție cresc CPI-ul (cicli pe instrucțiune). Care este creșterea CPI dată de greșelile de predicție date de un predictor de tipul ramificația este întotdeauna luată (realizată)? Presupuneți că ieșirea ramificației este determinată în etajul EX, și ca urmare nu apare hazard de date, și nu se utilizează nici un slot de întârziere.
 - Repetăți punctul a) pentru un predictor de tipul “ramificația întotdeauna nu este luată”
 - Repetăți punctul a) pentru un predictor cu 2 biți.
 - Cu un predictor cu 2 biți, care este creșterea în viteză dacă am putea converti o jumătate din instrucțiunea de ramificație într-un mod care să înlocuiască o instrucțiune de ramificație cu o instrucțiune ALU? Presupuneți există aceeași șansă de înlocuire atât pentru instrucțiune prezisă corect cât și pentru aceea prezisă incorect.
 - Cu un predictor cu 2 biți, care este creșterea în viteză dacă am putea converti o jumătate din instrucțiunea de ramificație într-un mod care să înlocuiască fiecare instrucțiune de ramificație cu două instrucțiuni ALU? Presupuneți există aceeași șansă de înlocuire atât pentru instrucțiune prezisă corect cât și pentru aceea prezisă incorect.
 - Unele instrucțiuni de ramificație sunt mult mai predictibile decât altele. Dacă cunoaștem 80% din toate instrucțiunile de ramificație executate, este ușor să prezicem ramificațiile de întoarcere din buclă care sunt întotdeauna prezise corect, care este acuratețea unui predictor pe 2 biți pentru cele 20% de instrucțiuni de ramificație?
36. Acest exercițiu examinează acuratețea diferiților predictor pentru ramificații pentru următorul șablon repetitiv (de exemplu într-o buclă) al ieșirilor ramificației T, NT, T, T, NT.
- Care este acuratețea unor predictor de tipul ramificația este întotdeauna luată sau întotdeauna nu este luată pentru această secvență de ieșire a ramificației?
 - Care este acuratețea unui predictor cu 2 biți pentru primele 4 ramificații din acest șablon, presupunând că predictorul pornește în starea din stânga jos din figura 4.63 (predicția nu este luată)?

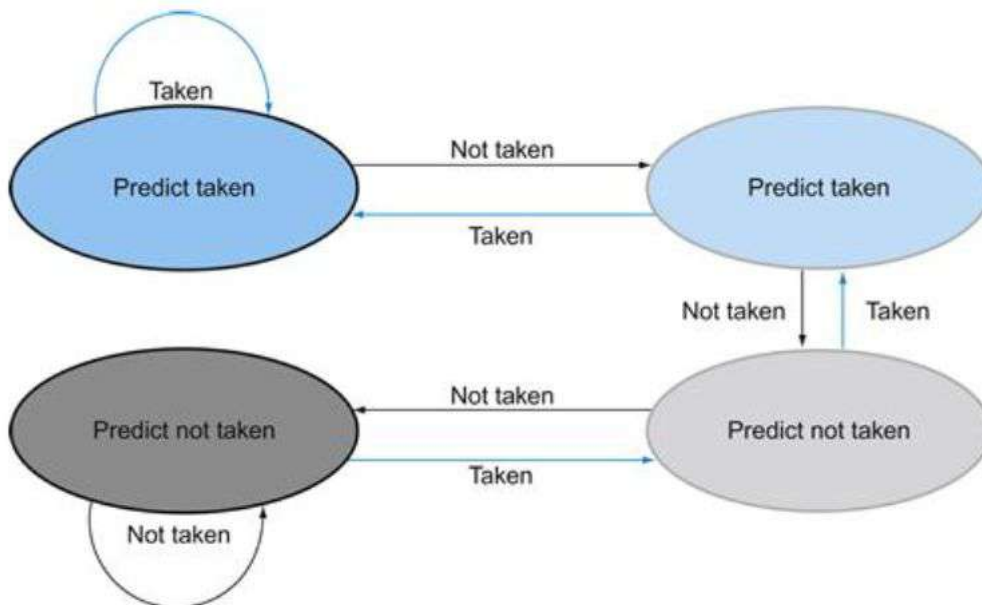


Figura 4.63

- Care este acuratețea unui predictor cu 2 biți dacă acest șablon se repetă la infinit?
- Proiectați un predictor care va atinge o acuratețe perfectă dacă acest șablon se va repeta la infinit. Predictorul vostru trebuie să fie un circuit secvanțial cu o ieșire care furnizează predicția (1 pentru ramificație luată și 0 pentru ramificație care nu este luată).
- Care este acuratețea pentru predictorul de la punctul d) dacă este dat un șablon de repetare exact în opoziție cu cel din problemă?
- Repetăți punctul d), dar acum predictorul va trebui să fie capabil ca eventual (după o perioadă de încălzire pe durata căreia poate face predicții greșite) să înceapă o predicție perfectă atât pentru șablonul din problemă cât și pentru opusul acestuia. Predictorul vostru trebuie o intrarea care să-i spună care a fost ieșirea reală. Sugestie: această intrare permite predictorului vostru să determine care din cele două șabloane este activ.

37. Acest exercițiu explorează modul în care excepțiile afectează proiectarea pipeline-ului. Primele trei întrebări se referă la următoarele două instrucțiuni:

Instruction 1	Instruction 2
BNE R1, R2, Label	LW R1, 0(R1)

- Care excepție poate fi generată de fiecare dintre aceste instrucțiuni? Pentru fiecare dintre aceste excepții, specificați etajul pipeline în care este detectată.
- Dacă există o adresă separată pentru handler-ul fiecărei excepții, prezentați cum trebuie schimbată organizarea pipeline-ului pentru a fi capabilă să gestioneze această execuție. Puteți presupune că adresele acestor handler-e sunt cunoscute atunci când este proiectat procesorul.
- Dacă a doua instrucțiune este extrasă imediat după prima instrucțiune, descrieți ce se întâmplă în pipeline dacă prima instrucțiune cauzează prima excepție pe care ați listat-o la punctul a). Prezentați diagrama de execuție pipeline de la

momentul în care este extrasă prima instrucțiune până când este completată prima instrucțiune din handler-ul de tratare a excepției.

- d) Atunci când se utilizează gestiunea vectorizată a excepțiilor, tabelul cu adresele handler-elor de tratare a excepțiilor este în memorie la o adresă (fixă) cunoscută. Modificați pipeline-ul pentru a implementa acest mecanism de gestiune a excepțiilor. Repetați punctul c) utilizând acest pipeline modificat și gestiunea vectorizată a excepțiilor.
- e) Dorim să emulăm gestiunea vectorizată a excepțiilor pe un procesor care are doar o adresă fixă pentru un handler. Scrieți codul care trebuie să fie la această adresă.

38. În acest exercițiu se compară performanțele unui procesor care execută o instrucțiune pe ciclu și unul care execută două instrucțiuni pe ciclu, luând în calcul transformările din program care pot fi realizate pentru a optimiza execuția cu două căi (două instrucțiuni pe ciclu). Întrebările din acest exercițiu se referă la următoarea buclă (programul este scris în C):

```
for (i=0; i!=j; i+=2)
```

```
  b[i] = a[i] - a[i+1]
```

Când scrieți codul MIPS presupuneți că variabilele sunt păstrate în regiștrii după cum urmează, și că toți regiștrii cu excepția celor indicați ca fiind liberi (Free) sunt utilizați pentru a păstra diferite variabile, astfel încât nu trebuie folosiți în oricare altă parte.

i	j	a	b	c	Free
R5	R6	R1	R2	R3	R10, R11, R12

- a) Translați codul C în instrucțiuni MIPS. Translatarea trebuie să fie directă, fără rearanjarea instrucțiunilor pentru a obține performanțe mai bune.
- b) Dacă există bucle după execuția a doar două iterații, desenați o diagramă pipeline pentru codul MIPS al vostru ca parte a execuției pe procesorul cu două căi prezentat în figura 4.69. Presupuneți că procesorul are un predictor perfect pentru predicția ramificațiilor și poate extrage oricare două instrucțiuni (nu doar instrucțiuni consecutive) în același ciclu.

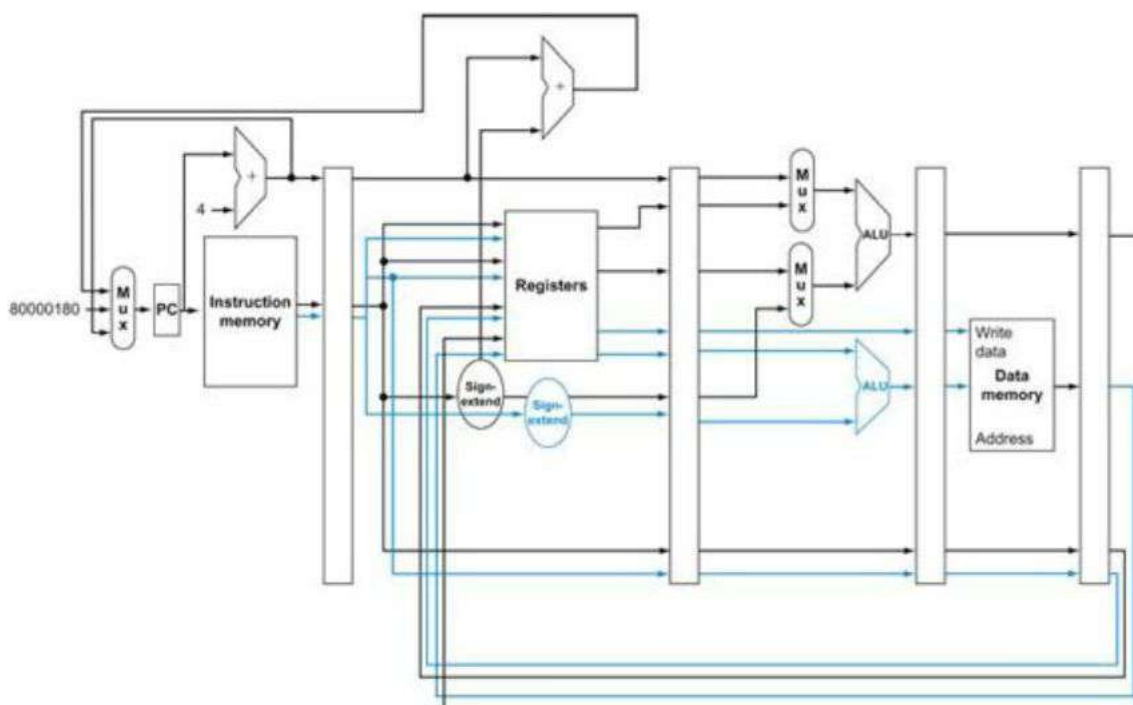


Figura 4.69

- c) Rearanjați codul de la punctul a) pentru a atinge performanțe mai bune pe procesorul cu planificare statică cu două căi din figura 4.69.
 - d) Repetați punctul b) dar de data aceasta utilizați codul MIPS de la punctul c).
 - e) Care este creșterea de viteză obținută prin trecerea de la un procesor cu o cale la procesorul cu două căi din figura 4.69? Utilizați codul de la punctul a) pentru ambele tipuri de procesoare și presupuneți că sunt executate 1000000 de iterații ale buclei. Ca la punctul b) presupuneți că procesorul are un predictor perfect pentru ramificații, și că procesorul cu 2 căi poate extrage oricare două instrucțiuni în același ciclu.
 - f) Repetați punctul e) dar presupuneți de data aceasta că pentru procesorul cu două căi una din instrucțiunile care se pot într-un ciclu poate fi de orice fel dar a doua trebuie să fie o instrucțiune care nu lucrează cu memoria.
39. Acest exercițiu explorează eficiența energetică și relația acesteia cu performanțele. Întrebările din acest exercițiu presupun următoarele consumuri energetice pentru activitățile din memoria de instrucțiuni, regiștrii și memoria de date. Puteți presupune că celelalte componente din calea de date consumă o cantitate neglijabilă d energie.

I-Mem	1 Register Read	Register Write	D-Mem Read	D-Mem Write
140pJ	70pJ	60pJ	140pJ	120pJ

Presupuneți că, componentele din calea de date au următoarele întârzieri. Puteți presupune că celelalte componente din calea de date au întârzieri neglijabile.

I-Mem	Control	Register Read or Write	ALU	D-Mem Read or Write
200ps	150ps	90ps	90ps	250ps

- a) Cât de multă energie este consumată pentru a executa o instrucțiune ADD pentru o proiectare cu un singur ciclu și pentru o proiectare pipeline cu 5 etaje?
- b) Care este cazul cel mai defavorabil de instrucțiune MIPS privind consumul de energie, și care este consumul de energie necesar pentru a o executa?
- c) Dacă reducerea consumului de energie este criteriul de bază al proiectării, cum se va schimba aceasta pentru pipeline? Care este procentul de reducere al consumului energetic pentru o instrucțiune LW după această schimbare?
- d) Care este impactul schimbării pe care ați realizat-o la punctul c) privind performanțele?
- e) Putem elimina semnalul de control MemRead și putem citi datele din memoria de date în fiecare ciclu, adică putem avea semnalul MemRead permanent pe 1. Explicați de ce procesorul totuși mai lucrează corect după această modificare. Care este efectul acestei modificări asupra frecvenței de ceas și a consumului energetic?
- f) Dacă o unitate în așteptare consumă 10% din puterea pe care o consumă dacă este activă, care este energia consumată de memoria de instrucțiuni în fiecare ciclu? Ce procent din energia totală consumată de memoria de instrucțiuni reprezintă această energie când blocul este în așteptare?

Probleme propuse pentru examenul de la
disciplina Structura și Organizarea
Calculatoarelor

== MOODLE TYPE ==

Prof. dr. ing. Vasile Gheorghiuță GĂITAN
Ș. I. dr. ing. Ionel ZAGAN

Cap. 4

Studiul de caz 1: optimizarea performanțelor memoriei cache via tehnici avansate.

Conceptele ilustrate:

- ☐ Cache fără blocare,
- ☐ Optimizările oferite de compilator pentru cache,
- ☐ Pre-extragere software și hardware,
- ☐ Calculul impactului performanțelor cache-ului în procesoare din ce în ce mai complexe.

1 Transpusa unei matrice schimbă între ele rândurile și coloanele așa cu se prezintă în continuare:

```
A11 A12 A13 A14 A21 A22 A23 A24 A31 A32 A33 A34 A41 A42 A43 A44      =>
A11 A21 A31 A41 A12 A22 A32 A42 A13 A23 A33 A43 A14 A24 A34 A44
```

Fie următorul cod C simplu care realizează transpusa:

```
for (i = 1; i < 5; i++) {
    for (j = 1; j < 5; j++) {
        ouput [j][i] = input [i][j];
    }
}
```

Presupuneți că atât matricea de intrare cât și cea de ieșire sunt memorate în ordinea bazată pe rândul major (row major order – ceea ce înseamnă că indicele rândului se modifică cel mai repede). Presupuneți că executați o transpunere dublă precizie 256 x 256 pe un procesor cu memorie cache L1 de 16 KB, pe deplin asociativă cu algoritmul LRU pentru politica de înlocuire a paginilor, organizată pe blocuri de 64 de octeți. Presupuneți că lipsurile sau pre-extragerile (misses) la nivelul L1 necesită 16 cicli de ceas și există întotdeauna potrivire la nivelul memorie cache L2, și că L2 poate procesa o cerere la fiecare doi cicli procesor. Presupuneți că fiecare iterație din bucla interioară necesită patru cicli dacă data este prezentă în L1. Presupuneți că cache-ul are pentru lipsurile la scriere o politică de tipul write-allocate (alocă line cache) fetch-on-write (extrage blocul din nivelul imediat inferior). Deși nerealist, presupuneți că scrierea înapoi a paginilor modificate necesită 0 cicli. Pentru implementarea simplă dată anterior, ordinea de execuție nu va fi ideală pentru matricea de intrare; totuși aplicând o optimizare prin interschimbarea buclilor se va crea o ordine neideală pentru matricea de ieșire. Deoarece interschimbarea buclilor nu este suficientă pentru îmbunătățirea performanțelor, în loc acestea trebuie blocată execuția.

- a. Care trebuie să fie mărimea minimă a cache-ului pentru a profita de avantajul blocării execuției?
- b. Care va fi numărul relativ de lipsuri pentru versiunile cu blocare sau fără blocare în comparație cu mărimea minimă aleasă anterior?
- c. Scrieți cod pentru a realiza transpusa cu mărimea blocului B care utilizează B x B blocuri?
- d. Care este minimul de asociativitate necesară pentru cache-ul L1 pentru a obține performanțe consistente independent de poziția ambelor matrice în memorie?

- e. Încercați transpoziția cu blocare sau fără blocare a unei matrice de 256 x 256 pe un calculator. Cât de aproape se potrivesc rezultatele cu așteptările voastre bazate pe ceea ce cunoașteți despre sistemul de memorie a calculatoarelor? Explicați, dacă este posibil, orice discrepanță.

- 2 Presupuneți că proiectați un prefetcher (un dispozitiv de pre extragere) pentru codul de transpunere fără blocare a unei matrice din problema 1. Cel mai simplu tip de prefetcher hardware pre-extrage numai secvențial blocuri cache la apariția unei miss (lipse). Prefetcher-ele hardware mai complicate de tip “non-unit stride” (pas non-unitate) ¹ pot analiza șirul de referințe care creează lipsuri (miss) și să detecteze și să pre-extragă cu pași non-unitate. În contrast, pre extragerea software poate determina pași non-unitate la fel de ușor cum poate determina și pașii unitate. Presupuneți că pre extragerea scrie direct în cache și că nu există “poluare” (suprascrierea datelor care trebuie utilizate înaintea ca datele să fie pre extrase). Pentru a obține cele mai bune performanțe utilizând un prefetcher non-unitate, câte **pre extrageri trebuie pot fi în așteptare (outstanding la un moment dat,** dacă bucla internă este în starea de echilibru?

3

Studiul de caz 2: Să le adunăm pe toate la un loc: sisteme de memorie cu un grad ridicat de paralelism.

Conceptele ilustrate:

- ☐ Facilități transversale,
- ☐ Proiectarea ierarhiei de memorie .

- 4 Programul ce urmează poate fi utilizat pentru a evalua comportarea unui sistem de memorie. Cheia este aceea de a avea o evaluarea precisă privind timpul și pasul programului prin memorie invocând diferite niveluri ale ierarhiei. Programul este scris în C. Prima parte este o procedură care utilizează un utilitar standard pentru a prelua o măsură cu acuratețe mare privind timpul CPU al utilizatorului; această procedură poate fi schimbată pentru a lucra pe diferite sisteme. Partea a doua, este o buclă cuibărite pentru a citi și a scrie memoria cu diferiți pași (strides) și diferite mărimi ale memorie cache. Pentru a prelua cu acuratețe timpul, acest cod este repetat de multe ori. Partea a treia măsoară, doar pentru a fi extras din timpul total măsurat, timpul pentru supracontrolul din buclele cuibărite pentru a vedea cât de lung este accesul. Rezultatele sunt memorate într-un fișier de ieșire de format .CSV pentru a facilita importul în programe de calcul tabelar. Se poate modifica CACHE_MAX în funcție de întrebarea la care trebuie să răspundeți și mărimea memoriei din sistemul pe care l-ați măsurat. Execuția programului în modul singur utilizator sau cel puțin fără a activa alte

¹ “Stride”(pas) reprezintă numărul de octeți pe care cineva trebuie să-i mute pentru a trece de la un element la altul (la următorul sau la precedentul) . Pasul nu poate fi mai mic decât mărimea elementului, dar poate fi mai mare, indicând un extra spațiu între elemente. O arie cu pasul de aceeași mărime cu mărimea fiecărui element este continuă în memorie. O astfel de arie se spune că are **pasul unitate** (unit stride). Non-unit stride arrays (ariile cu pasul **non-unitate**) sunt uneori mai eficiente în particular pentru matrici (2D) sau pentru arii multidimensionale, depinzând de efectul asupra memoriei cache și de șablonul de acces ales).

programe va da rezultate mult mai consistente. Presupuneti că mărimea tuturor componentelor din ierarhia de memorie sunt puteri ale lui 2. Presupuneti că mărimea unei pagini este mult mai mare ca mărimea unui bloc din al doilea nivel de memorie cache, și că mărimea unui bloc din al doilea nivel de cache este mai mare sau egal cu mărimea unui bloc de pe primul nivel de memorie cache. În figura 2.30 este trasat un exemplu de program; cheile indică mărimea ariilor testate.

```
#include "stdafx.h"
#include <stdio.h>
#include <time.h>
#define ARRAY_MIN (1024) /* 1/4 smallest cache */
#define ARRAY_MAX (4096*4096) /* 1/4 largest cache */
int x[ARRAY_MAX]; /* array going to stride through */

double get_seconds() { /* routine to read time in seconds */
    __time64_t ltime;
    _time64( &ltime );
    return (double) ltime;
}

int label(int i) { /* generate text labels */
    if (i<1e3) printf("%1dB",i);
    else if (i<1e6) printf("%1dK",i/1024);
    else if (i<1e9) printf("%1dM",i/1048576);
    else printf("%1dG",i/1073741824);
    return 0;
}

int _tmain(int argc, _TCHAR* argv[]) {
    int register nextstep, i, index, stride;
    int csize;
    double steps, tsteps;
    double loadtime, lastsec, sec0, sec1, sec; /* timing variables */

    /* Initialize output */
    printf(" ");
    for (stride=1; stride <= ARRAY_MAX/2; stride=stride*2)
        label(stride*sizeof(int));
    printf("\n");

    /* Main loop for each configuration */
    for (csize=ARRAY_MIN; csize <= ARRAY_MAX; csize=csize*2) {
        label(csize*sizeof(int)); /* print cache size this loop */
        for (stride=1; stride <= csize/2; stride=stride*2) {

            /* Lay out path of memory references in array */
            for (index=0; index < csize; index=index+stride)
                x[index] = index + stride; /* pointer to next */
            x[index-stride] = 0; /* loop back to beginning */

            /* Wait for timer to roll over */
            lastsec = get_seconds();
            do sec0 = get_seconds(); while (sec0 == lastsec);
```

```

/* Walk through path in array for twenty seconds */
/* This gives 5% accuracy with second resolution */
steps = 0.0; /* number of steps taken */
nextstep = 0; /* start at beginning of path */
sec0 = get_seconds(); /* start timer */
do { /* repeat until collect 20 seconds */
    for (i=stride;i!=0;i=i-1) { /* keep samples same */
        nextstep = 0;
        do nextstep = x[nextstep]; /* dependency */
        while (nextstep != 0);
    }
    steps = steps + 1.0; /* count loop iterations */
    sec1 = get_seconds(); /* end timer */
} while ((sec1 - sec0) < 20.0); /* collect 20 seconds */
sec = sec1 - sec0;

/* Repeat empty loop to loop subtract overhead */
tsteps = 0.0; /* used to match no. while iterations */
sec0 = get_seconds(); /* start timer */
do { /* repeat until same no. iterations as above */
    for (i=stride;i!=0;i=i-1) { /* keep samples same */
        index = 0;
        do index = index + stride;
        while (index < csize);
    }
    tsteps = tsteps + 1.0;
    sec1 = get_seconds(); /* - overhead */
} while (tsteps<steps); /* until = no. iterations */
sec = sec - (sec1 - sec0);
loadtime = (sec*1e9)/(steps*csize);
/* write out results in .csv format for Excel */
printf("%4.1f,", (loadtime<0.1) ? 0.1 : loadtime);
}; /* end of inner for loop */
printf("\n");
}; /* end of outer for loop */
return 0;
}

```

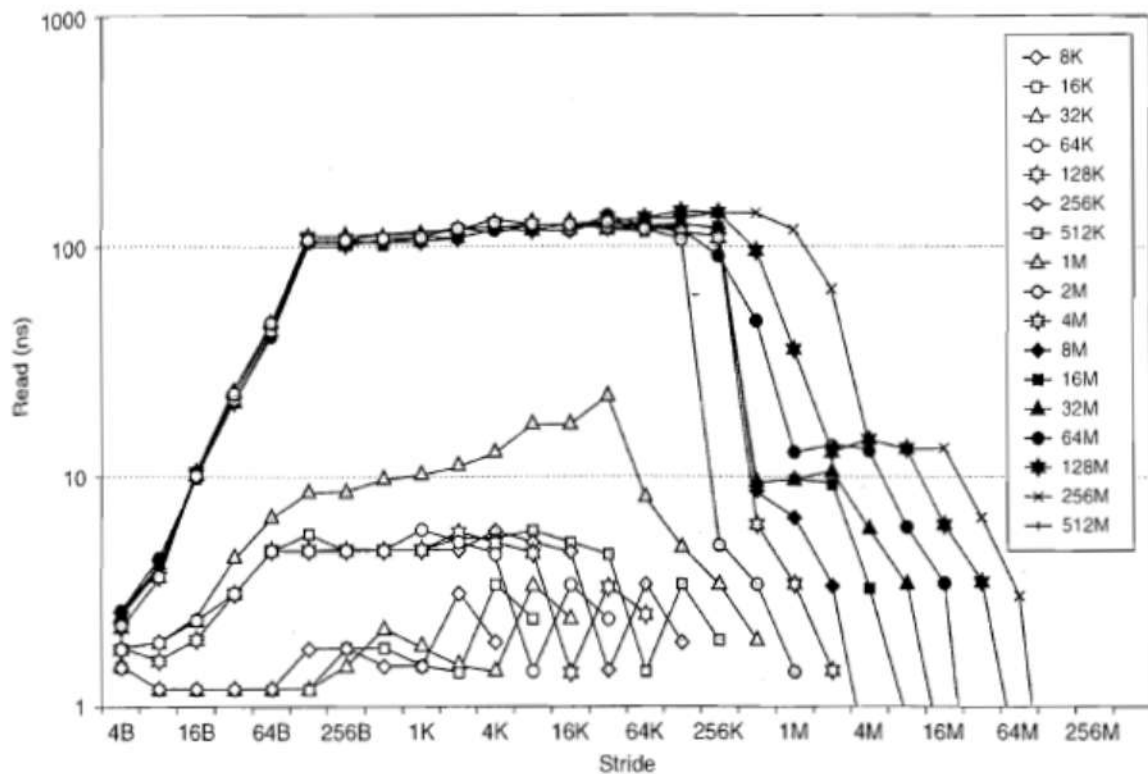


Figure 2.30 Sample results from program in Figure 2.29.

Utilizând exemplul de program prezentat anterior:

- Care este mărimea totală și mărimea unui bloc pentru memoria cache de pe nivelul 2?
- Care este penalitatea dată de lipsuri (miss) de pe nivelul 2 al memoriei cache?
- Care este asociativitatea pe nivelul 2 al memoriei cache?
- Care este mărimea memoriei principale?
- Care este timpul unei pagini dacă mărimea unei pagini este de 4KB?

5 Dacă este necesar modificați codul de la problema 4 pentru a măsura următoarele caracteristici sistem. Trasați rezultatele experimentale cu timpul consumat pe axa y și pasul de memorie pe axa x. Utilizați scalarea logaritmică pentru ambele axe, și desenați o linie pentru fiecare mărime a cache-ului.

- Ce este mărimea unei pagini sistem?
- Cât de multe intrări are buffer-ul de traducere cu anticipare (translation lookaside buffer - TLB)?
- Ce este penalitate pentru lipsuri pentru TLB?
- Ce este asociativitatea TLB-ului ?

6 În sistemele multi – microprocesor, nivelurile inferioare ale ierarhiei memoriei pot să nu fie saturate de un singur procesor dar trebuie să fie capabile să fie saturate de procesoare multiple care lucrează împreună. Modificați codul de la problema 4, și rulați copii multiple în același moment. Puteți determina:

- Câte procesoare actuale sunt în calculatorul vostru și câte procesoare sistem sunt doar contexte adiționale pentru fire multiple de execuție (multithreading)?
- Câte controlere de memorie nare sistemul vostru?

- 7 Puteți să vă gândiți la o cale pentru a testa unele caracteristici a unui cache pentru instrucțiuni utilizând un program? Indicație: compilatorul poate genera un număr larg de instrucțiuni care nu sunt evidente pentru un anumit cod. Încercați să utilizați instrucțiuni aritmetice simple de lungime cunoscută din arhitectura setului vostru de instrucțiuni (Instruction Set Architecture -ISA).
- 8 Această problemă investighează impactul unei memorii cache mici și simple utilizând CACTI și luând în calcul o tehnologie de 65 nm (0,065 μm). (CACTI este disponibil online la <http://quid.hpl.hp.com:9081/cacti/>.)
- Comparați timpul de acces a unui cache de 64KB cu blocuri de 64 de octeți și cu un singur banc. Care este timpul relativ de acces al memoriilor cache asociative cu două căi și cu patru căi în raport cu o organizare cu mapare directă?
 - Comparați timpul de acces al unei memorii cache asociative cu patru căi cu blocuri de 64 de octeți și un singur banc. Care este timpul relativ de acces a unor meorii cache de 32 și respectiv 64 de KB în comparație cu un cache de 16KB?
 - Pentru un cache de 64 KB, găsiți asociativitatea între 1 și 8 cu cel mai mic timp de acces mediu la memorie, cunoscând faptul că lipsurile per instrucțiune pentru o anumită suită de încărcări este de 0,00664 pentru maparea directă, 0,00366 pentru asociativitatea cu două căi, 0,000987 pentru asociativitatea cu patru căi, și 0,000266 pentru asociativitatea cu 8 căi. Per total sunt 0,3 referințe per instrucțiune. Presupuneți că lipsa cache ia 10ns. Pentru a calcula potrivirea în cicli, preluați ieșirea CACTI care calculează timpul unui ciclu, care corespunde frecvenței maxime la care poate opera un cache fără bubble (goluri) în pipeline .
- 9 În această problemă se investighează posibilele beneficii ale unei căi de a prezice (way-predicted) pentru memoria cache L1. Presupuneți că într-un sistem memoria cache L1 având mărimea de 64 KB organizată asociativ cu 4 căi este cea care limitează timpul unui ciclu mașină. Ca o alternativă la organizarea memoriei cache, trebuie să considerați o memorie cache modelată cu o cale de predicție, având 64KB, organizată cu mapare directă și cu o acuratețe a predicției de 80%. Până când nu se precizează altceva, presupuneți că accesul la calea de predicție care asigură potrivirile (hits) în cache presupune un ciclu mașină în plus. Presupuneți că rata pentru penalitățile de lipsuri (miss) și rata de potrivire (hit) sunt cele din problema 5, întrebarea c.
- Care este media timpului de acces la memorie pentru cache-ul curent (în cicli) versus cache-ul care are o cale de predicție?
 - Dacă toate celelalte componente pot opera cu memoria cache cu o cale de predicție care este mai rapidă (incluzând memoria principală), care va fi impactul asupra performanțelor atunci când se utilizează memoria cu o cale de predicție?
 - Memoria cache cu cale de predicție se utilizează în mod uzual pentru memoria cache pentru instrucțiuni care utilizează o coadă sau un buffer pentru instrucțiuni. Imaginați-vă că doriți să încercați calea de predicție pentru o memorie cache pentru date. Presupuneți că aveți o acuratețe de 80% și că operația care urmează (de exemplu operații dependente de accesul la cache al altor instrucțiuni) este lansată presupunând o cale de predicție corectă. Astfel, o cale de predicție greșită presupune golirea benzii și o capcană pentru răspuns (replay trap), care necesită 15 cicli. Este pozitivă sau negativă schimbarea timpului mediu de acces la memorie per instrucțiunea de încărcare (load instruction) cu memorie cache de date cu cale de predicție, și cât de mare este schimbarea?
 - Ca o alternativă la calea de predicție, multe memorii cache nivel L2 cu asociativitate mare serializează accesul la etichetă (tag) și dată, astfel încât trebuie activată doar

aria cerută care conține setul de date. Acest fapt micșorează consumul de putere dar crește timpul de acces. Utilizați interfața WEB CACTI pentru: un proces de 0,065 μm , capacitatea memorie cache de 1MB, organizarea memoriei cu mapare asociativă pe seturi cu 4 căi, mărimea unui bloc de 64B, citirea la ieșire a 144biți, un banc, un singur port de citire/scriere, etichete pe 30 de biți, și tehnologia ITRS-HP cu fire globale. Care este raportul între timpul de acces pentru serializarea etichetelor și a datelor, și accesul paralel?

- 10 Ați fost rugat să investigați pentru un procesor nou performanțele relative a unei memorii cache de date L1 organizată pe bancuri sau sub formă de bandă de asamblare (pipelined). Presupuneți o memorie cache cu mărimea de 64KB, cu mapare asociativă pe seturi cu două căi având mărimea unui bloc de 64 de octeți (64O-octeți sau 64B - byte). Banda de asamblare pentru cache constă din trei etaje, similar în capacitate cu memorie cache de date de la Alpha 21264. O implementare pe bancuri va consta din două bancuri de 32KB cu mapare asociativă pe seturi cu două căi. Pentru a răspunde la următoarele întrebări utilizați CACTI și presupuneți utilizarea unei tehnologii pe 0,065 μm . Durata unui ciclu în versiunea WEB arată la ce frecvență poate lucra o memorie cache fără bule (bubble) în banda de asamblare.
- Care este durata unui ciclu pentru memoria cache în comparație cu timpul său de acces, și câte etaje pipeline va avea cache-ul (până la două zecimale) ?
 - Comparați aria și energia dinamică totală pentru citire per proiectarea de tip pipelined versus proiectarea utilizând bancuri. Precizați care va necesita o arie mai mică și care va necesita putere mai multă, și explicați de ce.
- 11 Considerați pentru lipsurile memoriei cache de nivel 2 (L2) mecanismele denumite cuvântul critic primul (critica word first)² și repornirea timpurie (early restart)³. Presupuneți o memorie cache L2 de 1MB, cu blocuri de 64 de octeți și o cale de reumplere (refill) de 16 octeți. Presupuneți că L2 poate fi scrisă cu 16 octeți la fiecare 4 cicli procesor, că timpul de recepție a primului bloc de 16 octeți din controlerul de memorie este de 120 de cicli, că fiecare bloc adițional de 16 octeți din memorie principală necesită 16 cicli, și că data poate fi direct trimisă (direct bypass) la portul de citire al memoriei cache L2. Ignorați orice ciclu de transfer a unei cereri de tip lipsă (miss) la L2 și datele cerute pentru memoria cache L1.
- Câți cicli vor fi necesari pentru o lipsă L2 cu și fără mecanismele denumite cuvântul critic primul (critica word first) și repornirea timpurie (early restart)?
 - Credeți că mecanismele denumite cuvântul critic primul (critica word first) și repornirea timpurie (early restart) vor fi mult mai importante pentru memorie cache L1 sau pentru L2, și care sunt factorii care contribui la importanța lor relativă?
- 12 Proiectați un buffer de scriere între memoria cache L1 cu scrierea de tip scriere imediată (write through) și memoria cache L2 cu scrierea de tip scriere întârziată (write-back). Magistrala de date pentru L2 este de 16B și poate realiza o scriere la o adresă cache independentă în 4 cicli procesor.
- Câți octeți trebuie să aibă fiecare intrare în buffer-ul de scriere?

² Cuvântul critic primul (critica word first) – cere prima dată de la memorie cuvântul lipsă și trimite-l către UCP imediat ce sosește; lasă UCP să-și continue execuția în timp ce se umple restul de cuvinte din bloc. Se mai numește extragere adaptată (wrapped fetch)

³ Repornire timpurie (early restart) – cât de curând (imediat ce) sosește cuvântul cerut din bloc, acesta este trimis de urgență la UCP (procesor) lăsându-l astfel pe acesta să-și continue execuția.

- b. La ce creștere de viteză ne putem aștepta în starea de echilibru utilizând un buffer de scriere cu fuzionare (merging write buffer)⁴ în loc de un buffer de scriere fără fuzionare atunci când se filează cu 0 memorie prin execuția unei memorări pe 64 de biți dacă toate celelalte instrucțiuni pot fi lansate în paralel cu memorarea și blocurile sunt prezente în cache-ul L2.
- c. Care poate fi efectul unei posibile lipse (miss) la nivelul memoriei cache L1 relativ la numărul de intrări cerute de la buffer-ul de scriere, pentru sisteme cu memorie cache cu blocare și fără blocare?
- 13 Considerați un sistem de tip desktop cu un procesor conectat la un RAM de 2GB cu cod de corectare a erorilor (error-correcting code - ECC). Presupuneți că există un singur canal de memorie cu 72 de biți, dintre care 64 sunt pentru date și 6 pentru ECC.
- Câte chipuri DRAM are DIMM dacă sunt utilizate chipuri DRAM de 1 GB, și câte operații I/O trebuie executate pentru fiecare DRAM dacă doar un DRAM este conectat la fiecare pin de date al DIMM-ului?
 - Ce lungime de salvă (burst) este necesară pentru a suporta blocuri de memorie cache L2 de mărime 32B.
 - Calculați vârful de lărgime de bandă pentru DIMM-urile DDR2-667 și DDR2-533 pentru citiri dintr-o pagină activă excluzând supracontrolul dat de ECC.
- 14 În figura 2/31 este prezentat un exemplu de diagramă de timp pentru DDR2 SDRAM.

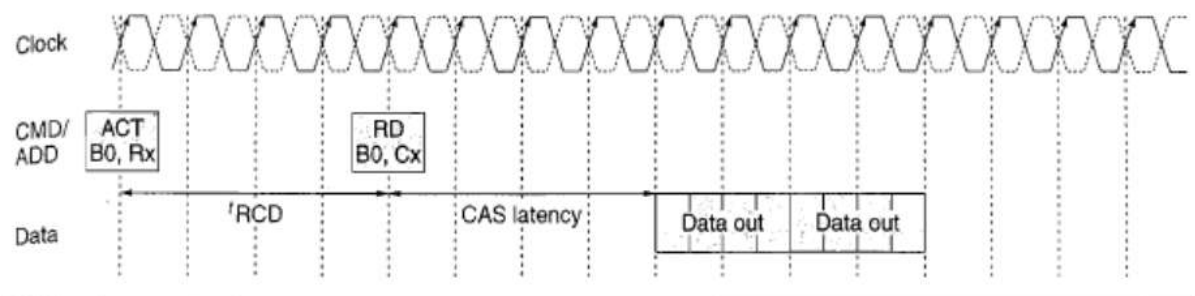


Figure 2.31 DDR2 SDRAM timing diagram.

t_{RCD} este timpul necesar pentru a activa un rând într-un banc, și latența (CL) pentru strobul de adresă pentru coloană (CAS) este numărul de cicli necesari pentru a citi la ieșire o coloană din rând. Presupuneți că RAM-ul este pe DIMM-ul DDR2 cu ECC standard, având 72 de linii de date. Presupuneți de asemenea că lungimea unei salve de 8 citește la ieșire 8 biți cea ce înseamnă un total de 64B de la DIMM. Presupuneți $t_{RCD} = CAS \text{ (or CL)} * \text{clock_frecuency}$, și că $\text{clock_frecuency} = \text{transfer_per_seconds}/2$. Latența pe chip la o lipsă cache de la nivelul 1 la nivelul 2 și înapoi, nu include accesul la DRAM, în cele 20ns.

- Cât timp este necesar de la prezentarea unei comenzi active până la ultimul bit de date cerut de la tranziția DRAM de la valid la invalid pentru DIMM-ul DDR2-667 1GB cu CL = 5. Presupuneți că pentru fiecare cerere se va preextrage automat o altă linie cache adiacentă din aceeași pagină.
- Care este latența relativă, atunci când se utilizează DIM-ul DDR2-667, a unei cereri de citire a unui banc activ versus una pentru o pagină deja deschisă, incluzând și timpul necesar pentru a procesa o lipsă în interiorul procesorului?

⁴ Bufferul de scriere cu fuzionare (Merging write buffer) combină scrierile care au ca destinație adrese consecutive într-o singură intrare în buffer. Altfel, acestea vor ocupa intrări separate care vor crește șansa de oprire a benzii de asamblare a procesorului.

- 15 Presupuneți că pentru 130\$ de dolari puteți achiziționa un DIMM de tip DDR2-667 de 2 GB cu CL=5 și la 100\$ un DIMM DDR2-533 de 2 GB cu un CL =4. Presupuneți că cele două DIMM-uri sunt utilizate în sistem, și restul sistemului costă 800\$. Considerați performanțele sistemului utilizând DIMM-urile DDR2-667 și DDR2-533 în condițiile unei execuții cu 3,33 lipsuri la 1K de instrucțiuni, și presupuneți că 80% din toate citirile din DRAM necesită o activare. Presupunând că este restantă la un moment o lipsă în L2 și că avem un nucleu care lucrează în ordine (in-order) cu un CPI de 1,5 fără a include timpul de acces dat de lipsurile L2, care este cost-permanța întregului sistem când se utilizează DIMM-uri diferite?
- 16 Ați achiziționat un server cu 8 nuclee cu frecvența de 8GHz CMO, care poate executa o încărcare cu un CPI de 2 (presupunând că nu sunt întârziate lipsurile de reumplere la memoria cache L2). Linia cache L2 are 32 de octeți. Presupuneți că sistemul utilizează DIMM-uri DDR2-667, câte canale de memorie independente de memorie trebuie furnizate astfel ca sistemul să nu fie limitat de lărgimea de bandă a memoriei dacă lărgimea de bandă necesară este uneori de două ori mai mare ca media? Încărcarea implică, per media 6,67 lipsuri L2 la 1K instrucțiuni.
- 17 Activarea unei pagini poate necesita o cantitate mare de energie (mai mult de trei ori) pentru DRAM (vezi <http://download.micron.com/pdf/technotes/ddr2/TN4704.pdf> și www.micron.com/systemcalc). Presupuneți că veți construi un sistem cu 2GB de memorie utilizând fie 8 bancuri de tip 2Gb x 8 DDR2 DRAM sau 16 bancuri de 1 Gb x 16 DRAM, ambele cu aceeași viteză. Ambele utilizează pagini de 1KB și ultimul nivel pentru liniile cache au 64B. Presupuneți că DRAM-urile care nu sunt active sunt în așteptare cu consum redus pentru preîncărcare (precharged standby) și disipă o putere neglijabilă. Presupuneți că timpul de tranziție din starea de așteptare cu consum redus în starea activă este nesemnificativ.
- De la care tip de DRAM se va aștepta să furnizeze performanțe sistem mai mari? Explicați de ce?
 - Din punctul de vedere al puterii, cum se compară un DIMM de 2 GB care se construiește cu un DRAM de 1GB x 8 DDR2 cu un DIMM de capacitate similară care se construiește cu 1Gb x 4 DDR2?
- 18 Pentru a accesa date de la un DRAM tipic, trebuie să activăm prima dată rândul corespunzător. Presupuneți ca această activare aduce o întreagă pagină de mărime 8KB în bufferr-ul pentru rânduri. Apoi, selectăm o coloană din buffer-ul pentru rânduri. Dacă accesese imediat următoare la DRAM sunt la aceeași pagină, atunci se poate sări pasul de activare; altfel trebuie să închidem pagina curentă și să pre-încărcăm linia de bit pentru următoarea activare. O altă politică populară pentru DRAM este de a închide pro-activ pagina și de a pre-încărca linia de bit imediat ce apare un acces. Presupuneți că fiecare citire sau scriere din/ în DRAM este de 64B și latența DDR (Data aut în figura 2.31) pentru a trimite 512 biți este T_{ddr}.
- Presupunând că DDR2-667 ia 5 cicli pentru preîncărcare, cinci cili pentru activare, și patru cicli pentru a citi o coloană, pentru ce valoare a fracției de potrivire la bufferul de rând (r) alegeți o politică în loc de cealaltă pentru a obține cel mai bun timp de acces? Presupuneți că fiecare acces la DRAM este separat de timp suficient pentru a termina un nou acces aleatoriu.
 - Dacă 10% din accesul total la DRAM se întâmplă unul după altul sau continuu fără nici un interval de timp pierdut, cum se va schimba decizia?
 - Calculați diferența mediei energiei consumate de DRAM per acces între cele două politici utilizând fracția de potrivire pentru buffer-ul de rând calculată anterior.

Presupuneți că preîncărcarea necesită 2 nJ, activarea 4 nJ și 100pJ sunt necesari pentru a citi sau scrie în buffer-ul de rând.

- 19 Ori de câte ori un calculator este în așteptare (idle), fie îl putem trece în starea cu consum redus (stanby), când DRAM încă este activ, sau îl putem lăsa să hiberneze. Presupuneți că, pentru hibernare, trebuie doar să copiem conținutul DRAM-ului pe un mediu nevolatil precum un Flash. Dacă citire sau scrierea unei linii Flash de 64 B necesită 2,56 μ J, iar DRAM-ul necesită 0,5 nJ, iar dacă consumul în idle al DRAM-ului este de 1,6W (pentru 8 GB), cât timp trebuie un sistem trecut în idle pentru a beneficia de hibernare? Presupuneți că memoria principală are 8GB.
- 20 Mașinile virtuale (VM) au potențial pentru a aduce multe facilități benefice unui sistem de calcul, precum îmbunătățirea costului total al utilizării (TCO) sau a disponibilității. Pot fi VM-urile utilizate pentru a furniza următoarele facilități? Dacă da, cum pot acestea oferi aceste facilități?
- Testarea aplicațiilor în mediile de producție utilizând sisteme de dezvoltare?
 - Repornirea rapidă a aplicațiilor în caz de dezastru sau cădere?
 - Performanțe înalte în aplicații cu operații intensive de I/O?
 - Izolarea defectelor între diferite aplicații, rezultând o disponibilitate crescută pentru servicii?
 - Permite mentenanța software-ului pe sisteme în timp ce aplicațiile rulează fără întreruperi semnificative?
- 21 Mașinile virtuale pot pierde performanța datorită unor evenimente cum ar fi: execuția instrucțiunilor privilegiate, lipsuri în TLB, capcane și operații I/O. Aceste evenimente sunt uzual gestionate de codul sistem. Astfel, o cale de a estima încetinirea atunci când se lucrează pe VM-uri este procentul de timp când aplicația se execută în modul sistem versus modul utilizator. De exemplu, o aplicația care petrece 10% din execuția sa în modul sistem poate încetini cu 60% când se execută pe o VM, Figura 2.32 listează performanțele timpurii pentru diferitele apeluri de sistem sun execuție nativă, virtualizare pură și para virtualizare pentru LMBench utilizând sisteme Xen pe Itanium cu măsurarea timpului în microsecunde.

Benchmark	Native	Pure	Para
Null call	0.04	0.96	0.50
Null I/O	0.27	6.32	2.91
Stat	1.10	10.69	4.14
Open/close	1.99	20.43	7.71
Install sighandler	0.33	7.34	2.89
Handle signal	1.69	19.26	2.36
Fork	56.00	513.00	164.00
Exec	316.00	2084.00	578.00
Fork + exec sh	1451.00	7790.00	2360.00

Figure 2.32 Early performance of various system calls under native execution, pure virtualization, and paravirtualization.

- La ce tipuri de programe vă așteptați să aibă încetiniri mai mici când aceste rulează sub VM?

- b. Dacă încetinirea este liniară ca funcție de timpul sistem, fiind dată încetinirea de la punctul anterior, cât de mult încetinește un program ce petrece 20% din timpul său de execuție în **modul sistem** ?
- c. Care este încetinirea mediană a apelurilor sistem din tabel sub virtualizarea pură și paravirtualizare?
- d. Care funcție din tabel are cea mai mare încetinire? Care credeți că poate fi cauza acesteia?

22 Popek și Goldberg au definit mașina virtuală spunând că nu poate fi distinsă față de o mașină reală decât prin performanțele sale. În această întrebare, se va utiliza acea definiție care să ne permită să știm dacă avem acces la execuția nativă pe un procesor sau execuția este pe o mașină virtuală. Tehnologia Intel VT-x oferă efectiv un al doilea set de niveluri privilegiate pentru a fi utilizate de mașinile virtuale. Ce trebuie să facă o mașină virtuală ce rulează pe o altă mașină virtuală, dacă se utilizează tehnologia VT-x?

- 23 O dată cu adoptarea suportului pentru virtualizare pe arhitecturile x86, mașinile virtuale au evoluat în mod activ și au devenit “de masă”. Comparați și scoateți în evidență diferențele între tehnologiile Intel VT-x și AMD – V. (Informațiile despre AMD-V pot fi găsite la <http://sites.amd.com/us/business/it-solutions/virtualization/Pages/resources.aspx>.)
- a. Care dintre ele poate furniza performanțe superioare pentru aplicații cu acces intensiv la memorie în cazul unor memorii mari?
 - b. Informațiile despre AMD IOMMU pot fi găsite la <http://developer.amd.com/us/documentation/articles/pages/892006101.aspx>.) Ce face tehnologia de virtualizare și unitatea de gestiune a memoriei (IOMMU) pentru a îmbunătăți performanțele de virtualizare a operațiilor I/O?

24 Deoarece paralelismul la nivelul instrucțiunilor poate fi efectiv exploatat atât în cazul procesoarelor superscalare cu execuție în-ordine, cât și pentru procesoarele de tip VLIW (Very large Instruction Word – cu instrucțiuni cu cuvânt foarte mare) cu execuție speculativă, un motiv foarte important pentru construirea unui procesor cu execuție în afara ordinii (Out-of-Order - OOO) este acela de a avea abilitatea de a tolera întârzierea nepredictibilă a memorie cauzată de lipsurile cache. Ca urmare, puteți gândi că hardware-ul care suportă OOO este parte a sistemului de memorie! Priviți amprenta pe siliciu a procesorului Alpha 21264 din figura 2.33 pentru a găsi aria relativă a cozielor și a dispozitivelor de mapare pentru întregi și virgulă flotantă versus memoriile cache. Cozile planifică lansarea instrucțiunilor în execuție, iar dispozitivele de mapare redenumesc specificatorii de regiștrii. Ca urmare acestea sunt elemente adiționale pentru a suporta execuția de tip OOO. Procesorul 21264 are pe chip doar cache de nivel L1 pentru date și instrucțiuni, ambele având mărimea de 64KB cu mapare asociativă pe seturi cu două căi. Utilizați un simulator superscalar precum SimpleScalar (www.cs.wisc.edu/~mscalar/simplescalar.html) cu un test (benchmark) intensiv pe accese la memorie pentru a găsi câtă performanță se pierde dacă aria pentru cozi și dispozitive de mapare este utilizată pentru memoria de date cache L1 într-un procesor superscalar cu execuția în ordine, în loc să fie utilizată pentru execuția OOO așa cum este modelată această în 21264. Asigurați-vă că celelalte aspecte ale procesoarelor sunt cât mai similare posibil pentru a avea o comparație cinstită. Ignorați orice creștere în acces sau timpul unui ciclu dată de memoriile cache mari, și efectul memoriilor cache de date mari asupra repartiției pe zone a chipului (de notat că această comparație nu va fi în totalitate corectă, deoarece codul nu va fi planificat de compilator pentru un procesor cu execuție în ordine).

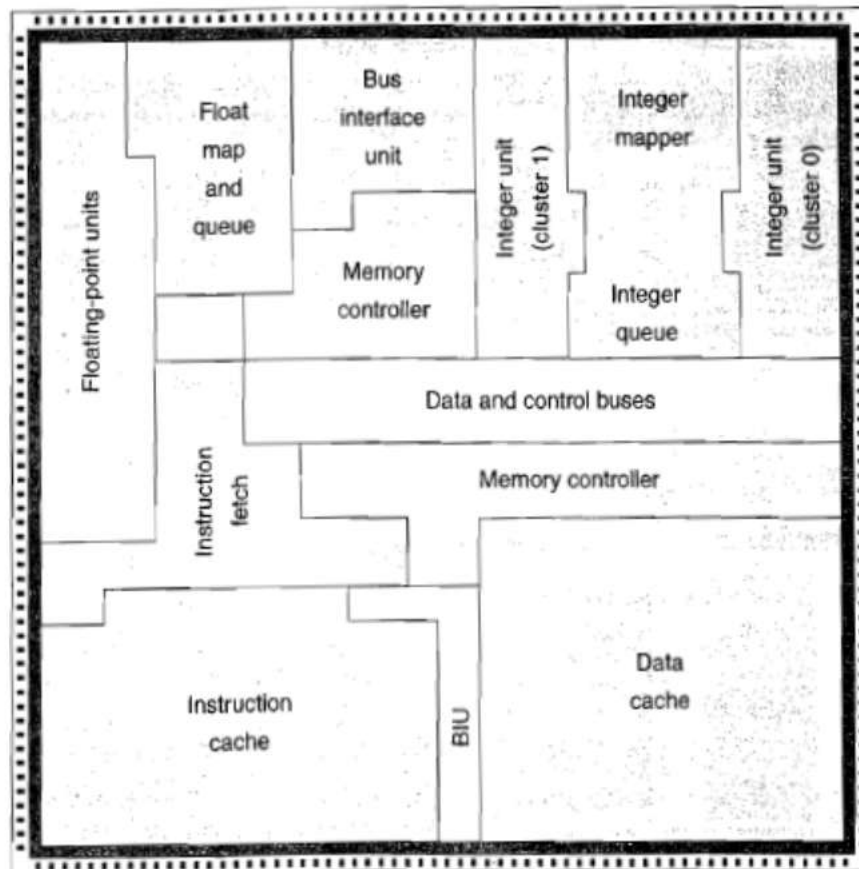


Figure 2.33 Floorplan of the Alpha 21264 [Kessler 1999].

25 Analizorul Intel VTune pentru performanțe poate fi utilizat pentru a realiza multe măsurători privind performanțele memorie cache. O versiune de evaluare liberă atât pentru Windows cât și pentru Linux poate fi descărcată de la <http://software.intel.com/en-us/articles/intel-vtune-amplifier-xe/>.

Programul (aca_ch2_cs2.c) utilizat în studiul 2 de caz a fost modificat astfel încât să poată lucra cu ieșirea din VTune din Microsoft Visual C++. Programul poate fi descărcat de la www.hpl.hp.com/research/cacti/aca_ch2_cs2_vtune.c. Au fost înserate funcții VTune speciale pentru a exclude supracontrolul dat de inițializare și bucle pe durata procesului de analizare a proceselor. Direcțiile de setare detaliate ale VTune sunt date în secțiunea README din program. Programul păstrează bucla pentru 20 de secunde pentru fiecare configurație. În experimentul care urmează puteți găsi efectul mărimii datelor în cache și performanțele generale ale procesorului. Rulați programul VTune pe un procesor Intel cu mărimea datelor de intrare de 8KB, 128KB, 4MB și 32MB, și păstrați un pas de 64 de octeți (pasul unei linii cache pe procesoarele i7). Colectați statisticile privind performanțele generale și performanțele memoriilor cache de date L1, L2 și L3.

- Listați numărul de lipsuri per 1K instrucțiuni în cache-ul de date L1, L2 și L3 pentru fiecare mărime a setului de date pentru viteza și modelul procesorului vostru. Pe baza acestor rezultate, ce puteți spune despre mărimea memoriilor cache L1, L2 și L3 de pe procesorul vostru? Explicați observațiile.
- Listați instrucțiunile per ceas (instructions per clock - IPC) pentru fiecare mărime a setului de date, modelul vostru de procesor și viteza acestuia. Pe baza acestor rezultate, ce puteți spune despre penalitățile referitoare la lipsurile memoriilor cache L1, L2 și L3 de pe procesorul vostru?

- c. Rulați programul VTune cu mărimi ale setului de date de 8KB și 128KB pe un procesor Intel cu OOO. Pentru ambele configurații listați numărul de lipsuri cache L1 și L2 per 1K instrucțiuni și CPI-ul. Ce puteți spune despre eficacitatea tehnicilor de ascundere a latenței memoriei în procesoarele OOO de mare performanță? Sfat: trebuie să găsiți întârzierea dată de latența de la nivelul memorie cache L1 pentru procesorul vostru. Pentru i7 recente aceasta este de 11 cicli.

Probleme propuse pentru examenul de la
disciplina Structura și Organizarea
Calculatoarelor

== MOODLE TYPE ==

Prof. dr. ing. Vasile Gheorghiuță GĂITAN
Ș. I. dr. ing. Ionel ZAGAN

Cap. 5. Instruction-Level Parallelism and Its Exploitation

1. Care vor fi performanțele de referință (în cicli, per o iterație a buclei) a secvenței de cod din figura 3.48 dacă nu se poate iniția nici o instrucțiune nouă până când nu s-a completat execuția vechii instrucțiuni. Ignorați extragerea și decodificarea instrucțiunii. Presupuneți acum că execuția nu este încetinită dacă lipsește următoarea instrucțiune, dar pot fi lansată doar o instrucțiune/ciclu. Presupuneți că saltul este luat, și că există un slot de întârziere pentru ramificații.

			Latencies beyond single cycle	
Loop:	LD	F2,0(RX)	Memory LD	+4
I0:	DIVD	F8,F2,F0	Memory SD	+1
I1:	MULTD	F2,F6,F2	Integer ADD, SUB	+0
I2:	LD	F4,0(Ry)	Branches	+1
I3:	ADD	F4,F0,F4	ADD	+1
I4:	ADD	F10,F8,F2	MULTD	+5
I5:	ADDI	Rx,Rx,#8	DIVD	+12
I6:	ADDI	Ry,Ry,#8		
I7:	SD	F4,0(Ry)		
I8:	SUB	R20,R4,Rx		
I9:	BNZ	R20,Loop		

Figure 3.48 Code and latencies for Exercises 3.1 through 3.6.

2. Gândiți-vă ce înseamnă în realitate numerele date pentru latență – ele indică numărul de cicli pe care o funcție dată îi cere pentru a produce ieșirea și nimic mai mult. Dacă se cunosc la nivelul global al liniei pipeline cicli de latență pentru fiecare unitate funcțională, atunci cel puțin se poate garanta că orice pereche de instrucțiuni una după alta (un producător urmat de un consumator) se va executa corect. Dar, nu toate perechile de instrucțiuni au o relație producător/consumator. Uneori două instrucțiuni adiacente nu au nimic de a face una cu alta. Câți cicli va necesita corpul buclei din secvența de cod din figura 2.48 dacă linia pipeline detectează dependențele reale de date și produce încetinirea pentru acestea, mai degrabă decât să încetinească orbește orice doar pentru că o unitate funcțională este ocupată? Prezentați codul cu <stall> înserat acolo unde este necesar pentru a realiza latențele stabilite.
3. Luați în considerare o proiectare cu lansări multiple în execuție. Presupuneți că aveți două linii pipeline pentru execuție, fiecare fiind capabilă să înceapă execuția unei instrucțiuni per ciclu și, în față, o lărgime de bandă suficientă pentru fetch/decode astfel încât să nu fie încetinită execuția. Presupuneți că rezultatul poate fi imediat avansat de la o unitate de execuție la alta, sau de la ea însăși. Mai departe presupuneți că singurul motiv pentru care pipeline-ul încetinește este acela de a observa o dependență adevărată de date. Ca urmare câți cicli necesită bucla?
4. Pentru proiectarea cu lansări multiple în execuție de la exercițiul 3, se pot recunoaște unele subtilități. Chiar dacă cele două pipeline au exact același repertoriu, ele nu sunt nici identice, nici interschimbabile, deoarece există o ordine implicită între ele care trebuie să reflecte ordinea instrucțiunilor din programul original. Dacă instrucțiunea N+1

Își începe execuția în Pipe 1 în același timp cu instrucțiunea N din Pipe 0, și dacă N+1 necesită o latență în execuție mai mică decât N, atunci N+1 se va completa înainte de N (deși ordinea programului necesită completarea inversă). Descrieți cel puțin două motive de ce această situație poate fi hazard și necesită o considerare specială în micro-arhitectură. Dați un exemplu de două instrucțiuni din codul din figura 3.48 care demonstrează acest hazard.

5. Reordonați instrucțiunile pentru a îmbunătăți performanțele codului din figura 3.48. Presupuneți o mașină cu două pipeline-uri ca aceea din exercițiul 3 și că execuția în afara ordinii din exercițiul 4 a fost distribuită cu succes. Fiți doar îngrijorați pe moment doar de a observa dependențele reale de date și latențele unităților funcționale. Câți cicli ia codul reordonat?
6. Fiecare ciclu care nu inițiază o nouă operație în pipeline este o oportunitate pierdută, în sensul în care hardware-ul nu funcționează la întreg potențialul său.
 - a) Ce fracție din toți cicli codului reorganizat la exercițiul 5, luând în considerare ambele pipeline-uri, a fost pierdută?
 - b) Desfășurarea buclelor este una dintre tehnicile standard ale compilatoarelor pentru a găsi mai mult paralelism la nivelul codului, cu scopul de a minimiza pierderile privind oportunitățile pentru performanță. Desfășurați manual două iterații ale buclei din codul reordonat la exercițiul 5.
 - c) Ce creștere de viteză ați obținut (pentru acest exercițiu doar colorați cu verde instrucțiunea iterației a (N+1) –a, pentru a o distinge de iterația a N-aș. Dacă desfășurați bucla trebuie să reasignați regiștrii pentru a preveni coliziunile între iterații).
7. Calculatoarele petrec o parte însemnată din timpul lor în bucle, astfel buclele cu iterații multiple sunt locurile cele mai adecvate unde se poate găsi execuție speculativă care să țină ocupate resursele CPU. Nimic nu este vreodată ușor; compilatorul emite doar o singură copie a codului buclei, astfel chiar dacă iterațiile multiple gestionează date distincte, ele apar ca utilizând aceeași regiștrii. Pentru a evita coliziunile datorate utilizării regiștrilor de către iterațiile multiple, se vor redenumii regiștrii. Figura 4.49 prezintă un exemplu de cod care pe dorim să o redenumescă hardware-ul. Compilatorul poate să desfășoare simplu bucla și să utilizeze regiștrii diferiți pentru a evita conflictele, dar dacă ne așteptăm ca hardware-ul nostru se desfășoare bucle, și acesta trebuie să facă redenumirea regiștrilor. Cum? Presupuneți că hardware-ul nostru are o arie de regitrii temporari (denumiți registrii T, și presupuneți că 64 dintre ei, T0 la T63) care pot substitui regiștri desemnați de compilator. Acest hardware de redenumire este indexat de denumirea src (sursă) a registrului, și valoarea în tabel este registrul T al ultimei destinații care avea ca țintă acel registru. (Gândiți-vă la aceste valori din tabel ca la producători, iar regiștrii src sunt consumatori; nu contează prea mult unde pun producătorii rezultatelor lor atât timp cât consumatorii le pot găsi). Considerați secvența de cod din figura 3.49. De fiecare dată când vedeți în cod un registru destinație, substituiți-l cu următorul T disponibil, începând cu T9. Apoi actualizați corespunzător toți regiștrii src, astfel încât dependențele reale de date să fie menținute. Prezentați codul rezultat.

Loop:	LD	F4,0(Rx)
I0:	MULTD	F2,F0,F2
I1:	DIVD	F8,F4,F2
I2:	LD	F4,0(Ry)
I3:	ADDD	F6,F0,F4
I4:	SUBD	F8,F8,F6
I5:	SD	F8,0(Ry)

Figure 3.49 Sample code for register renaming practice.

I0:	LD	T9,0(Rx)
I1:	MULTD	T10,F0,T9
...		

Figure 3.50 Hint: Expected output of register renaming.

8. Exercițiul 3.7 a explorat redenumirea simplă a regiștrilor: atunci când un bloc hardware de redenumire a regiștrilor vede un registru sursă, substituie regiștrul destinație T al ultimei instrucțiuni care avea ca țintă acel registru sursă. Atunci când un tabel de redenumire vede un registru destinație, acesta substituie pentru acesta următorul T disponibil, dar proiectarea superscalarilor trebuie să gestioneze instrucțiuni multiple pe un singur ceas la nivelul fiecărui etaj, incluzând cel pentru redenumirea regiștrilor. Prin urmare, un procesor scalar simplu trebuie să caute atât maparea pentru src pentru fiecare instrucțiune și să aloce o nouă mapare destinație per ciclu de ceas. Procesoare superscalare trebuie să fie capabile să facă același lucru, dar trebuie să asigure că orice relație dest-to-src între două instrucțiuni concurente este gestionată corect. Considerați exemplu de cod din figura 3.51. Presupuneți că dorim să redenumim simultan primele două instrucțiuni. Mai mult presupuneți că următorii doi T regiștrii disponibili pentru a fi utilizați sunt cunoscuți la începutul ciclului de ceas în care cele două instrucțiuni vor fi redenumite. Conceptual, ceea ce dorim este ca pentru prima instrucțiune să se facă căutarea în tabelul de redenumire și apoi să realizăm actualizarea tabelului destinație cu registrul T. Apoi a doua instrucțiune va face exact același lucru, și orice dependență între instrucțiuni va fi gestionată corect. Dar nu este suficient timp pentru a scrie registrul T desemnat în tabelul de redenumire și apoi să căutăm din nou pentru a doua instrucțiune totul în același ciclu de ceas. Ca urmare substituția regiștrilor trebuie realizată imediat (în paralel cu actualizarea tabelului cu redenumirea regiștrilor). Figura 3.52 prezintă diagrama circuitului, care utilizează multiplexoare și comparatoare, care vor asigura redenumirea “din zbor”. Sarcina voastră este să prezentați starea ciclu cu ciclu a tabelului de redenumire pentru fiecare instrucțiune din codul din figura 3.51. Presupuneți că tabelul pornește cu fiecare intrare egală cu indexul său (T0=0, T1 = 1, ...)

I0:	SUBD	F1, F2, F3
I1:	ADDD	F4, F1, F2
I2:	MULTD	F6, F4, F1
I3:	DIVD	F0, F2, F6

Figure 3.51 Sample code for superscalar register renaming.

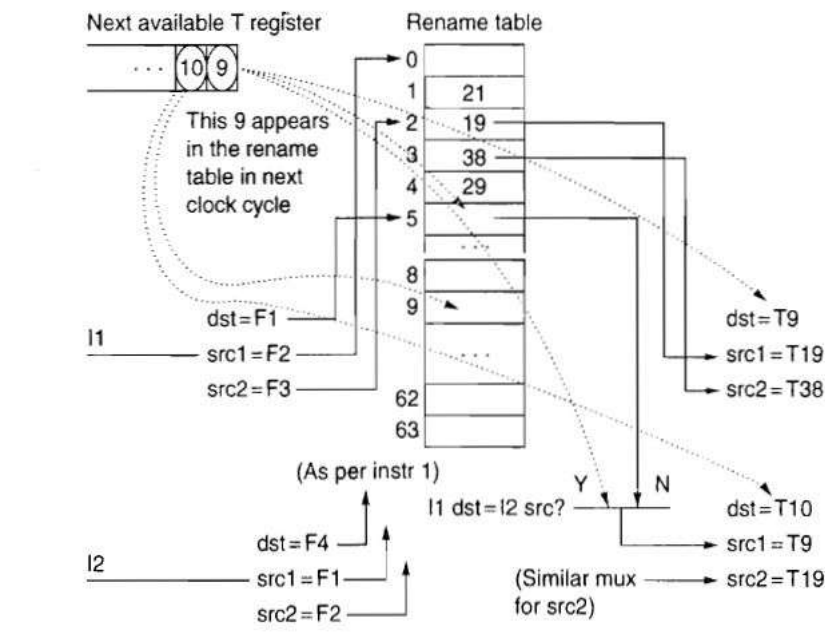


Figure 3.52 Rename table and on-the-fly register substitution logic for superscalar machines. (Note that src is source, and dst is destination.)

9. Dacă nu înțelegeți ce face redenumirea regiștrilor, mergeți înapoi la codul în asamblare pe care l-ați executat, și întrebați-vă ce urmează să se întâmple pentru a obține rezultatul corect. De exemplu, considerați un bloc de redenumire pentru un superscalar cu trei căi care să redenumască concurrent aceste trei instrucțiuni.

```

ADDI R1,R1,R1
ADDI R1,R1,R1
ADDI R1,R1,R1

```

Dacă valoarea din R1 este 5, care va fi valoarea atunci când va fi executată această instrucțiune?

10. Proiectanții pentru Very Long Instruction Word - VLIW au câteva opțiuni de bază pentru a face reguli arhitecturale de utilizare a regiștrilor. Presupuneți că un VLIW este proiectat cu pipeline-uri cu auto-drenare: odată ce este inițiată o operație, rezultatul acesteia va apărea în registrul destinație la cel mult L cicli mai târziu (unde L este latența operației). Nu există niciodată suficienți regiștrii, ca urmare există tentația de a stoarce maximum din regiștrii care există. Considerați figura 3.53. Dacă load are latebța de 1 + 2 cicli, desfășurați această buclă o dată și arătați cum un VLIW capabil de a executa două load și două adunări per ciclu poate utiliza numărul minim de regiștrii, în absența oricărei întreruperi sau încetiniri a pipeline-ului. Dați un exemplu de

un eveniment care, în prezența pipeline-ului cu autodrenare, poate deranja pipeline-ul și conduce la rezultate eronate.

```

Loop: LW      R4,0(R0) ; ADDI   R11,R3,#1
      LW      R5,8(R1) ; ADDI   R20,R0,#1
      <stall>
      ADDI    R10,R4,#1;
      SW      R7,0(R6) ; SW     R9,8(R8)
      ADDI    R2,R2,#8
      SUB     R4,R3,R2
      BNZ     R4,Loop

```

Figure 3.53 Sample VLIW code with two adds, two loads, and two stalls.

11. Presupuneți o micro-arhitectură cu un singur pipeline și cinci etaje (fetch, decode execute, memory, and write-back) și codul din figura 3.54. Toate operațiile sunt de 1 ciclu cu excepția LW și SW care sunt de 1 + 2 cicli, și ramificațiile care sunt de 1 + 1 ciclu. Nu există avansare. Prezentați fazele pentru fiecare instrucțiune per ciclu de ceas pentru o iterație în buclă.

```

Loop: LW      R3,0(R0)
      LW      R1,0(R3)
      ADDI    R1,R1,#1
      SUB     R4,R3,R2
      SW      R1,0(R3)
      BNZ     R4, Loop

```

Figure 3.54 Code loop for Exercise 3.11.

- Câți cicli de ceas per o iterație în buclă sunt pierduți datorită supracontrolului introdus de ramificație?
 - Presupuneți prezența unui predictor static capabil să recunoască salturile înapoi în etajul Decode. Acum, câți cicli de ceas sunt pierduți de supracontrolul dat de ramificații?
 - Presupuneți prezența unui predictor dinamic. Câți cicli sunt pierduți la o predicție corectă?
 - d)
12. Să considerăm ce poate realiza planificarea dinamică. Presupunem că avem arhitectura din figura 3.55. Presupuneți că ALU poate face toate operațiile aritmetice (MULTD, DIVD, ADDD, ADDI, SUB) și ramificațiile, și că stația de rezervare RS poate dispeceriza cel puțin câte o operație la fiecare unitate funcțională per ciclu (o operație la fiecare ALU plus o operație cu memorie la LD/ST).

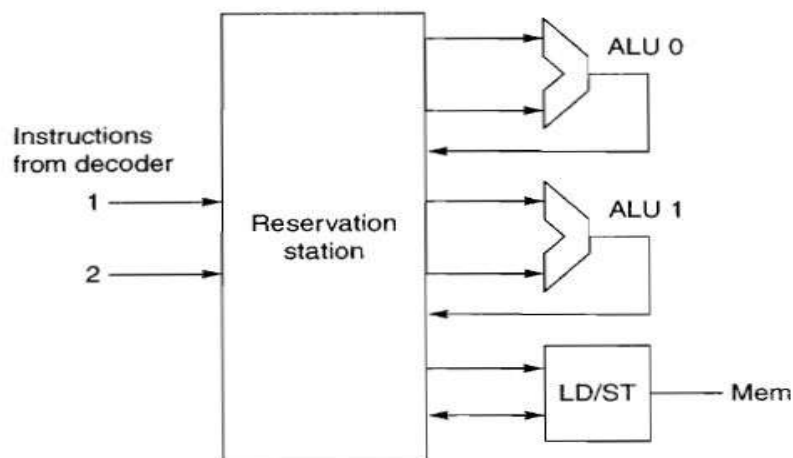


Figure 3.55 An out-of-order microarchitecture.

- a) Presupuneți că toate instrucțiunile din secvența din figura 3.48 sunt prezente în RS, fără execuția redenumirii. Marcați oricare instrucțiune din cod pentru care redenumirea registrului conduce la îmbunătățirea performanțelor. Presupuneți aceleași latențe pentru unitățile funcționale ca în figura 3.48.
- b) Presupuneți că versiunea cu redenumirea regiștrilor de la punctul a este rezidentă în RS în ciclul N, cu latențele date în figura 3.48. Arătați cum trebuie RS să dispecerizeze aceste instrucțiuni în afara ordinii, ceas cu ceas, pentru a obține performanțe optime pentru acest cod (presupuneți aceleași restricții pentru RS ca la punctul a, presupuneți că rezultatele trebuie scrise în RS înainte ca acestea să fie disponibile pentru utilizare – fără bypassing). Câți cicli de ceas ia execuția acestei secvențe de cod?
- c) Punctul b lasă stației de rezervare RS să încerce să planifice optimal aceste instrucțiuni. Dar în realitate, întreaga secvență de interes nu este uzual prezentă în RS. În loc, diferite evenimente șterg RS, și din noua secvență de instrucțiuni venită de la decodificator, RS trebuie să aleagă să dispecerizeze ceea ce are. Presupuneți că RS este goală. În ciclul 0, primele două instrucțiuni cu redenumirea regiștrilor din această secvență apar în RS. Presupuneți că este necesar un singur ciclu de ceas pentru a dispeceriza orice operație, și presupuneți că latențele unităților funcționale sunt cele din exercițiul 3.2. Mai mult presupuneți că în față (decodificator/ blocul de redenumire a regiștrilor) va continua să alimenteze continuu cu două noi instrucțiuni per ciclu de ceas. Câți cicli de ceas necesită acum secvența de cod?
- d) Dacă se dorește îmbunătățirea performanțelor de la punctul (c) ce ne va ajuta mai mult: (1) încă un ALU? (2) Altă unitate LD/ST? (3) un bypassing complet a rezultatelor ALU către operațiile subsecvențiale? Sau (4) Împărțirea în două a celei mai mari latențe? Care este creștere de viteză?
- e) Să considerăm acum speculația, actul de extragere, decodificare și execuție dincolo de una sau mai multe ramificații condiționate. Motivația noastră de a face acest lucru este dublă: Planificarea dispecerizării de la punctul (c) are o mulțime de nop-uri, și cunoaștem calculatoare care își petrec majoritatea din timp executând bucle (care implică salturi înapoi la vârful buclei și care sunt ușor de prezis). Buclele ne spun unde să găsim mai multă treabă de făcut: planificarea împrăștiată a dispecerizării ne sugerează că avem oportunități să facem o parte din această muncă mai timpuriu decât după. La punctul (d) s-a găsit calea critică prin buclă. Imaginați o pliere a unei a doua copie a acelei căi obținută prin planificarea de la punctul (b). Câți cicli suplimentari sunt necesari pentru a executa aceste două bucle

cu valoare în execuție (Presupuneți că toate instrucțiunile sunt rezidente în RS)? (Presupuneți că toate unitățile funcționale sunt complet implementate sub formă de pipeline).

13. În acest exercițiu, veți explora performanțele de compromis între trei procesoare care implică diferite tipuri de multithreading (execuție de fire multiple). Fiecare dintre aceste procesoare este superscalar, utilizează pipeline-uri în ordine, necesită o încetinire fixă de trei cicli ce urmează pentru toate încărcările și ramificațiile, și au cache-ul L1 identic. Instrucțiunile din același fir lansate în execuție în același ciclu sunt citite în ordinea programului și trebuie să nu conțină nicio dependență de date sau control.

- ✓ Procesorul A este o arhitectură SMT (Multithreading simultan) superscalară, capabilă să lanseze două instrucțiuni per ciclu de la două fire.
- ✓ Procesorul B este o arhitectură MT fină, capabilă să lanseze 4 instrucțiuni per ciclu pentru un singur fir și să comute firele atunci când pipeline-ul este încetinit.
- ✓ Procesorul C este o arhitectură MT brută, capabilă să lanseze până la 8 instrucțiuni per ciclu dintr-un singur fir și să comute firele la lipsurile de la cache-ul de nivel 1 L1.

Aplicația avută în vedere este una de căutare în liste, care scanează o regiune de memorie pentru o valoare specifică memorată în R9 în gama de adrese specificată de R16 și R17. Este paralelizată fie prin divizarea în 4 blocuri ocupând spații continue de memorie și asignând câte un fir de căutare pentru fiecare bloc (implicând 4 fire). Majoritatea timpului de execuție a fiecărui fir este petrecut în corpul următoarei bucle desfășurate:

```
loop: LD R1,0(R16)
      LD R2,8(R16)
      LD R3,16(R16)
      LD R4,24(R16)
      LD R5,32(R16)
      LD R6,40(R16)
      LD R7,48(R16)
      LD R8,56(R16)
      BEQAL R9,R1,match0
      BEQAL R9,R2,match1
      BEQAL R9,R3,match2
      BEQAL R9,R4,match3
      BEQAL R9,R5,match4
      BEQAL R9,R6,match5
      BEQAL R9,R7,match6
      BEQAL R9,R8,match7
      DADDIU R16,R16,#64
      BLT R16,R17,loop
```

Presupuneți următoarele:

- ✓ Se utilizează o barieră pentru a se asigura că toate firele încep simultan.
- ✓ Prima lipsă la nivelul cache-ului L1 apare după două iterații ale buclei.
- ✓ Nici este luată (Taken) nici una dintre ramificațiile BEQAL.

- ✓ BLT este întotdeauna luat.
 - ✓ Toate trei procesoarele planifică firele de execuție într-o manieră roud-robin.
- Determinați câți cicli sunt necesari pentru fiecare procesor pentru a completa primele două operații ale buclei.

14. În acest exercițiu, vom examina modul în care o tehnică software poate extrage paralelismul la nivel instrucțiunilor (ILP) într-o buclă vector uzuală. Următoarea buclă este așa numita buclă DAXPY (calculul în dublă precizie aX plus Y) care este operația centrală în eliminarea Gaussian. Următorul cod implementează operația DAXPY, $Y = aX + Y$, pentru un vector cu lungimea 100. Inițial, R1 este setat cu adresa de bază a ariei X și R2 cu adresa de bază a ariei Y:

```

        DADDIU    R4,R1,#800    ; R1 = upper bound for X
foo:    L.D       F2,0(R1)      ; (F2) = X(i)
        MUL.D    F4,F2,F0      ; (F4) = a*X(i)
        L.D       F6,0(R2)      ; (F6) = Y(i)
        ADD.D    F6,F4,F6      ; (F6) = a*X(i) + Y(i)
        S.D       F6,0(R2)      ; Y(i) = a*X(i) + Y(i)
        DADDIU    R1,R1,#8      ; increment X index
        DADDIU    R2,R2,#8      ; increment Y index
        DSLTU     R3,R1,R4      ; test: continue loop?
        BNEZ      R3,foo        ; loop if needed

```

Presupuneți că latențele unităților funcționale sunt prezentate în tabelul următor. Presupuneți că întârzierea de un ciclu a ramificațiilor se rezolvă în etajul ID. Presupuneți că rezultatele sunt complet bypassing.

Instruction producing result	Instruction using result	Latency in clock cycles
FP multiply	FP ALU op	6
FP add	FP ALU op	4
FP multiply	FP store	5
FP add	FP store	4
Integer operations and all loads	Any	2

- Presupuneți un pipeline cu o singură lansare în execuție. Prezentați cum va arăta bucla atât neplanificată de compilator și după ce compilatorul o planifică atât pentru operațiile în virgulă flotantă și întârzierile ramificațiilor, incluzând încetinirile sau cicli de ceas de așteptare. Care este timpul de execuție (în cicli) per element al vectorului rezultat, Y , neplanificat și planificat? Cu cât trebuie să fie mai rapid ceasul pentru hardware-ul procesorului pentru ca acesta singur să atingă îmbunătățirea performanțelor dată de planificarea compilatorului? (neglijați orice efect al creșterii vitezei ceasului asupra performanțelor memoriei).
- Presupuneți un pipeline cu o singură lansare în execuție. Desfășurați bucla de atâtea ori cât este necesar pentru a o planifica fără nici o încetinire, reducând instrucțiunile care dau supracontrol buclei. De câte ori trebuie desfășurată bucla?

Prezentați planificarea instrucțiunilor, Care este timpul de execuție per element al rezultatului?

- c) Presupuneți utilizarea unui procesor VLIW care are instrucțiuni care conțin cinci operații așa cum se arată în figura 3.16. Se vor compara două grade de desfășurare a buclei. Primul, desfășoară bucla de 6 ori pentru a extrage ILP și o planifică fără nici o încetinire (de ex. ciclu de lansare complet gol), reducând instrucțiunile care dau supracontrol buclei, și apoi se repetă procesul dar se desfășoară bucla de 10 ori. Se ignoră slot-urile de întârziere pentru ramificații. Prezentați cele două planificări. Care este timpul de execuție per element al vectorului rezultat pentru fiecare planificare? Ce procentaj pe slot-ul de operare s-a utilizat pentru fiecare planificare? Cât de mult diferă mărimea codului pentru cele două planificări? Care este cererea totală de regiștrii pentru cele două planificări?

Memory reference 1	Memory reference 2	FP operation 1	FP operation 2	Integer operation/branch
L.D F0,0(R1)	L.D F6,-8(R1)			
L.D F10,-16(R1)	L.D F14,-24(R1)			
L.D F18,-32(R1)	L.D F22,-40(R1)	ADD.D F4,F0,F2	ADD.D F8,F6,F2	
L.D F26,-48(R1)		ADD.D F12,F10,F2	ADD.D F16,F14,F2	
		ADD.D F20,F18,F2	ADD.D F24,F22,F2	
S.D F4,0(R1)	S.D F8,-8(R1)	ADD.D F28,F26,F2		
S.D F12,-16(R1)	S.D F16,-24(R1)			DADDUI R1,R1,#-56
S.D F20,24(R1)	S.D F24,16(R1)			
S.D F28,8(R1)				BNE R1,R2,Loop

Figure 3.16 VLIW instructions that occupy the inner loop and replace the unrolled sequence. This code takes 9 cycles assuming no branch delay; normally the branch delay would also need to be scheduled. The issue rate is 23 operations in 9 clock cycles, or 2.5 operations per cycle. The efficiency, the percentage of available slots that contained an operation, is about 60%. To achieve this issue rate requires a larger number of registers than MIPS would normally use in this loop. The VLIW code sequence above requires at least eight FP registers, while the same code sequence for the base MIPS processor can use as few as two FP registers or as many as five when unrolled and scheduled.

15. În acest exercițiu, vom analiza modul în care performează algoritmul lui Tomasulo atunci când execută bucla din exercițiul 3.14.

Status	Wait until	Action or bookkeeping
Issue all instructions		<pre> if (RegisterStat[rs].Busy) /*in-flight instr. writes rs*/ {h ← RegisterStat[rs].Reorder; if (ROB[h].Ready) /* Instr completed already */ {RS[r].Vj ← ROB[h].Value; RS[r].Qj ← 0;} else {RS[r].Qj ← h;} /* wait for instruction */ } else {RS[r].Vj ← Regs[rs]; RS[r].Qj ← 0;} RS[r].Busy ← yes; RS[r].Dest ← b; ROB[b].Instruction ← opcode; ROB[b].Dest ← rd; ROB[b].Ready ← no; </pre>
FP operations and stores	Reservation station (r) and ROB (b) both available	<pre> if (RegisterStat[rt].Busy) /*in-flight instr writes rt*/ {h ← RegisterStat[rt].Reorder; if (ROB[h].Ready) /* Instr completed already */ {RS[r].Vk ← ROB[h].Value; RS[r].Qk ← 0;} else {RS[r].Qk ← h;} /* wait for instruction */ } else {RS[r].Vk ← Regs[rt]; RS[r].Qk ← 0;} </pre>
FP operations		<pre> RegisterStat[rd].Reorder ← b; RegisterStat[rd].Busy ← yes; ROB[b].Dest ← rd; </pre>
Loads		<pre> RS[r].A ← imm; RegisterStat[rt].Reorder ← b; RegisterStat[rt].Busy ← yes; ROB[b].Dest ← rt; </pre>
Stores		<pre> RS[r].A ← imm; </pre>
Execute FP op	(RS[r].Qj == 0) and (RS[r].Qk == 0)	Compute results—operands are in Vj and Vk
Load step 1	(RS[r].Qj == 0) and there are no stores earlier in the queue	<pre> RS[r].A ← RS[r].Vj + RS[r].A; </pre>
Load step 2	Load step 1 done and all stores earlier in ROB have different address	Read from Mem[RS[r].A]
Store	(RS[r].Qj == 0) and store at queue head	<pre> ROB[h].Address ← RS[r].Vj + RS[r].A; </pre>
Write result all but store	Execution done at r and CDB available	<pre> b ← RS[r].Dest; RS[r].Busy ← no; ∀x(if (RS[x].Qj==b) {RS[x].Vj ← result; RS[x].Qj ← 0}); ∀x(if (RS[x].Qk==b) {RS[x].Vk ← result; RS[x].Qk ← 0}); ROB[b].Value ← result; ROB[b].Ready ← yes; </pre>
Store	Execution done at r and (RS[r].Qk == 0)	<pre> ROB[h].Value ← RS[r].Vk; </pre>
Commit	Instruction is at the head of the ROB (entry h) and ROB[h].ready == yes	<pre> d ← ROB[h].Dest; /* register dest, if exists */ if (ROB[h].Instruction==Branch) {if (branch is mispredicted) {clear ROB[h], RegisterStat; fetch branch dest;}} else if (ROB[h].Instruction==Store) {Mem[ROB[h].Destination] ← ROB[h].Value;} else /* put the result in the register destination */ {Regs[d] ← ROB[h].Value;} ROB[h].Busy ← no; /* free up ROB entry */ /* free up dest register if no one else writing it */ if (RegisterStat[d].Reorder==h) {RegisterStat[d].Busy ← no;} </pre>

Figure 3.14 Steps in the algorithm and what is required for each step. For the issuing instruction, rd is the destination, rs and rt are the sources, r is the reservation station allocated, b is the assigned ROB entry, and h is the head entry of the ROB. RS is the reservation station data structure. The value returned by a reservation station is called the result. RegisterStat is the register data structure, Regs represents the actual registers, and ROB is the reorder buffer data structure.

Unitățile funcționale (FU) sunt descrise în tabelul următor.

FU Type	Cycles in EX	Number of FUs	Number of reservation stations
Integer	1	1	5
FP adder	10	1	3
FP multiplier	15	1	2

Presupuneri următoarele:

- ✓ Unitățile funcționale nu sunt de tip pipeline.
- ✓ Nu există forwarding între unitățile funcționale; rezultatele sunt comunicate de către common data bus (CDB).
- ✓ Etajul de execuție EX realizează atât calculul efectiv al adresei cât și accesul la memorie pentru loads și stores. Astfel pipeline-ul este IF/ID/IS/EX/WB.
- ✓ Loads necesită un ciclu de ceas.
- ✓ Etajele de lansare în execuție IS și WB necesită fiecare un ciclu de ceas.
- ✓ Există 5 slot-uri în buffer-ul de încărcare (load) și 5 slot-uri în buffer-ul de memorare (store).
- ✓ Presupunem că instrucțiunea BNEZ necesită un ciclu de ceas.

a) Pentru această problemă utilizați un pipeline MIPS Tomasulo cu o singură lansare în execuție ca aceea din figura 3.6 cu latențele definite în tabelul anterior. Prezentați numărul de cili de încetinire pentru fiecare instrucțiune și în care ciclu de ceas a început execuția fiecărei instrucțiuni (intră în primul ciclu de EX) pentru trei iterații ale buclei. Câți cicli durează fiecare iterație a buclei? Raportați răspunsul sub forma unui tabel cu următoarele capete de tabel:

- ✓ Iterația (numărul iterației buclei)
- ✓ Instrucțiunea
- ✓ Lansarea (ciclu când a fost lansată instrucțiunea)
- ✓ Execuția (cicli când instrucțiunea se execută)
- ✓ Accesul la memorie (cicli când este accesată memoria)
- ✓ Scrierea pe CDB (cicli când rezultatul este scris pe CDB)
- ✓ Comentarii (descrierea oricărui eveniment la care așteaptă instrucțiunea)

Prezentați în tabel trei iterații ale buclei.

b) Repetați punctul (a) dar presupuneți se utilizează algoritmul lui Tomasulo cu două lansări iar FPU este complet pipeline.

16. Algoritmul lui Tomasulo are un dezavantaj: se poate calcula per clock, per CDB un singur rezultat. Utilizați configurația hardware și latențele de la întrebarea precedentă și găsiți o secvență de cod nu mai mare de 10 instrucțiuni unde algoritmul lui Tomasulo trebuie să încetinească datorită disputei pentru CDB. Indicați unde apar acestea în secvența voastră.

17. Un predictor de ramificație cu corelație (m,n) utilizează comportarea celor mai recente m ramificații executate pentru a alege din 2^m predictorii, fiecare dintre aceștia fiind un predictor pe n biți. Un predictor local cu două niveluri lucrează într-o manieră asemănătoare, dar ține doar urma comportării trecute a fiecărei ramificații individuale pentru a prezice comportarea viitoare. Există un compromis în proiectare implicat de un astfel de predictor: predictorii cu corelare necesită mai puțină memorie pentru istorie ceea ce le permite să mențină predictorii pe 2 biți pentru un număr mare de ramificații individuale (reducând probabilitatea ca o instrucțiune de ramificație să reutilizeze un același predictor), în timp ce predictorii locali necesită substanțial mai multă memorie pentru a păstra istoria și sunt astfel limitați la urmărirea unui număr relativ mic de instrucțiuni de ramificație. Pentru acest exercițiu, considerați un predictor de corelație (1,2) care poate urmări patru ramificații (necesită 16 biți) versus un predictor local (1,2) care poate urmări două ramificații utilizând aceeași cantitate de memorie. Pentru fiecare ieșire următoarei ramificații, furnizați fiecare predicție, tabela de intrare utilizată pentru a realiza predicția, orice actualizare a tabelului ca rezultat al predicției, și rată predicției finale greșite pentru fiecare predictor. Presupuneți că toate ramificațiile până la acest punct au fost luate. Inițializați fiecare predictor după cum urmează:

Correlating predictor			
Entry	Branch	Last outcome	Prediction
0	0	T	T with one misprediction
1	0	NT	NT
2	1	T	NT
3	1	NT	T
4	2	T	T
5	2	NT	T
6	3	T	NT with one misprediction
7	3	NT	NT

Local predictor			
Entry	Branch	Last 2 outcomes (right is most recent)	Prediction
0	0	T,T	T with one misprediction
1	0	T,NT	NT
2	0	NT,T	NT
3	0	NT	T
4	1	T,T	T
5	1	T,NT	T with one misprediction
6	1	NT,T	NT
7	1	NT,NT	NT

Branch PC (word address)	Outcome
454	T
543	NT
777	NT
543	NT
777	NT
454	T
777	NT
454	T
543	T

18. Presupuneți că aveți un procesor cu un pipeline adânc, pentru care se implementează un buffer pentru țintele ramificațiilor (branch-target buffer) numai pentru ramificațiile condiționate. Presupuneți că penalitatea pentru predicție greșită este întotdeauna de 4 cicli și penalitatea pentru lipsa din buffer este întotdeauna de 3 cicli. Presupuneți că 90% este rata de potrivire, 90% acuratețea și 15 % este frecvența ramificațiilor. Cât de rapid este un procesor cu un buffer pentru țintele ramificațiilor (BTB) versus un procesor care are o penalitate fixă de doi cicli per ramificație? Presupuneți un ciclu de ceas de bază per instrucțiune (CPI) fără încetinire la ramificație.
19. Considerați un buffer pentru țintele ramificațiilor (BTB) care are penalități de zero, doi, și doi cicli de ceas pentru predicția corectă a ramificațiilor condiționate, predicție incorectă și respectiv o lipsă din buffer. Considerați o proiectare a unui BTB care distinge ramificațiile condiționate de cele necondiționate, și care memorează adresa țintă pentru ramificațiile condiționate și instrucțiunea țintă pentru ramificațiile necondiționate.
- Care este penalitatea în cicli de ceas atunci când o ramificație necondiționată este găsită în buffer?
 - Determinați îmbunătățirea dată de plierea ramificațiilor pentru ramificațiile necondiționate. Presupuneți o rată de potrivire de 90%, o frecvență a salturilor necondiționate de 5%, și o penalitate de 2 cicli pentru o lipsă din buffer. Cât de multă îmbunătățire se câștigă cu această sporire? Cât de mare trebuie să fie rata de potrivire pentru această sporire pentru a furniza câștig de performanță?

Probleme propuse pentru examenul de la
disciplina Structura și Organizarea
Calculatoarelor

== MOODLE TYPE ==

Prof. dr. ing. Vasile Gheorghiuță GĂITAN
Ș. I. dr. ing. Ionel ZAGAN

Cap. 6 Data-Level Parallelism in Vector, SIMD, and GPU Architectures

- 1) Presupuneți constantele prezentate în figura 4.2. Prezentați codul pentru MIPS și VMIPS. Presupuneți că nu putem utiliza încărcările și memorările împrăștiate-adunate. Presupuneți că adresele de start pentru tiPL, tiPR, clL, cșR și clP sunt în RtiPL, RtiPR, RclL, RclR și respectiv RclP. Presupuneți că lungimea registrului VMIPS este programabilă de către utilizator și poate fi asignată prin setarea registrului special VL (de exemplu li VL4). Pentru a facilita reducția pentru adunarea vectorilor, presupuneți că am adăugat următoarele instrucțiuni la VMIPS:

SUMR Fd, Vs Vector Summation Reduction Single Precision:

Această instrucțiune realizează reducția sumării pe un registru vector Vs, scriind suma într-un registru scalar Fd.

Constants	Values
AA,AC,AG,AT	0,1,2,3
CA,CC,CG,CT	4,5,6,7
GA,GC,GG,GT	8,9,10,11
TA,TC,TG,TT	12,13,14,15
A,C,G,T	0,1,2,3

Figure 4.32 Constants and values for the case study.

- 2) Presupuneți seq_lenght == 500, care este contorul pentru execuția dinamică pentru ambele implementări?
- 3) Presupuneți că instrucțiunea de reducere a vectorilor este executată în unitatea funcțională pentru vectori, similar cu instrucțiunea de adunare a vectorilor. Prezentați modul în care secvența de cod stabilită în convoi presupune o singură instanță a fiecărei unități funcționale. Cât de multe “chimes” necesită codul? Cât de mulți cicli per FLOP sunt necesari, ignorând supracontrolul dat de lansarea în execuție a instrucțiunii pe vector?
- 4) Presupuneți acum că putem utiliza încărcări și memorări împrăștiate-adunate LVI și SVI. Presupuneți că tiPL, tiPR, clR, și clP sunt ranjate consecutiv în memorie. De exemplu, dacă seq_lenght == 500, aria tiPR va începe la 500 * 4 octeți în memorie după aria tiPL. Cum va afecta aceasta modul în care puteți scrie codul VMIPS pentru acest nucleu? Presupuneți că puteți inițializa regiștrii vector cu întregi utilizând următoarele tehnici care pot, de exemplu, inițializa registrul vector V1 cu valorile (0, 0, 2000, 2000):

```

LI R2,0
SW R2,vec
SW R2,vec+4
LI R2,2000
SW R2,vec+8
SW R2,vec+12
LV V1,vec

```

Presupuneți că lungimea maximă a vectorului este 64. Există vreun mod în care pot fi îmbunătățite performanțele utilizând încărcări împrăștiate-adunate? Dacă da, cu cât?

- 5) Presupuneți că dorim să implementăm nucleul MrBayes pe un GPU utilizând un singur bloc de tip fir. Rescrieți codul în C al kernel-ului utilizând CUDA. Presupuneți că pointerul la tabelul cu tranzițiile de probabilitate și cu probabilitățile sunt specificați ca parametri pentru kernel. Invocați un fir pentru fiecare interacțiune din buclă. Încărcați și reutilizați valorile în memoria partajată înainte de a realiza operații pe ea.
- 6) Cu CUDA se poate utiliza paralelismul grosier – fir la nivelul bloc pentru a calcula probabilitățile condiționale a mai multor noduri care lucrează în paralel. Presupuneți că vom calcula probabilitățile condiționale de la baza spre vârful unui arbore. Presupuneți că ariile cu probabilitățile condiționale și probabilitățile de tranziție sunt organizate în memorie, așa cum se prezintă în figura 4.33. Schimbați metoda prin care calculați indicele în răspunsul de la Exercițiul 4.5. pentru a include numărul blocului.

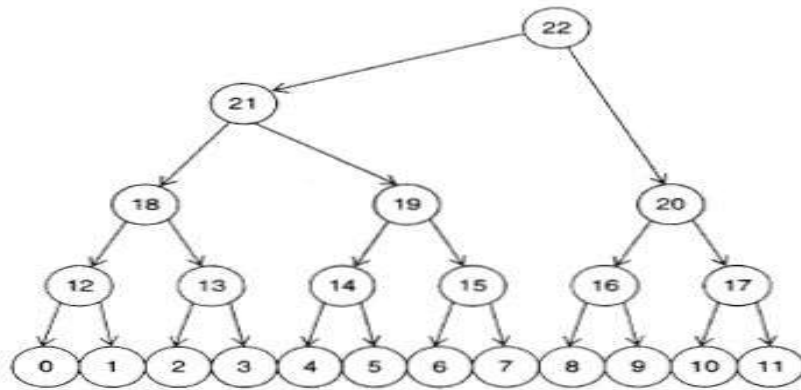


Figure 4.33 Sample tree.

- 7) Converteți codul de la exercițiul 4.6 în cod PTX. Cât de multe instrucțiuni sunt necesare pentru kernel?
- 8) Cât de bine vă așteptați ca acest cod să lucreze pe un GPU? Explicați răspunsul.

- 9) Considerați următorul cod, care multiplică doi vectori care conțin valori complexe în virgulă flotantă simplă precizie:

```
for (i = 0; i < 300; i++) {  
    c_re[i] = a_re[i] * b_re[i] - a_im[i] * b_im[i];  
    c_im[i] = a_re[i] * b_im[i] + a_im[i] * b_re[i];  
}
```

Presupuneți că procesorul rulează la 700 Mhz și are lungimea maximă a vectorului de 64. Unitatea load/store are un supracontrol de start de 15 cicli; unitatea de înmulțire, 8 cicli; și unitatea de adunare/ scădere de 5 cicli.

- Care este intensitatea aritmetică a kernel-ului? Justificați răspunsul.
 - Converțiți această buclă în asamblare VMIPS utilizând mineritul datelor pe benzi (strip-mining)
 - Presupunând o singură memorie pipeline înlănțuită, cât de mulți “chimes” sunt necesari? Cât de mulți cicli de ceas sunt necesari pentru rezultate complexe, incluzând supracontrolul la pornire (start-up)?
 - Dacă secvența de vectori este înlănțuită, cât de mulți cicli de ceas sunt necesari pentru valorile rezultat complexe, incluzând supracontrolul?
 - Presupuneți acum că procesorul are trei pipeline-uri și înlănțuiri. Dacă nu există conflict de banc în accesele din buclă, cât de mulți cicli de ceas sunt necesari per rezultat?
- 10) În această problemă, vom compara performanțele unui procesor vectorial cu un sistem hibrid care conține un procesor scalar și un coprocesor GPU. În sistemul hibrid, procesorul gazdă are performanțe scalare superioare față de GPU, astfel în caest caz tot codul scalar este executat pe procesorul gazdă în timp ce tot codul vectorial este executat pe GPU. Ne vom referi la primul sistem ca la calculatorul vectorial și la al doilea sistem ca la sistemul hibrid. Presupuneți că aplicația țintă conține un kernel vector cu o intensitate aritmetică de 0,5 FLOP per octet DRAM accesat; totuși, aplicația are o componentă scalară care trebuie realizată înainte și după kernel cu scopul de a pregăti vectorii de intrare, respectiv pe cei de ieșire. Pentru un eșantion de date de intrare, porțiunea scalară de cod necesită 400 ms de execuție atât pe procesorul vectorial cât și pe procesorul gazdă în sistemul hibrid. Kernel-ul citește vectorii de intrare constând din 200 Mo de date și are date de ieșire constând din 100 Mo de date. Procesorul vectorial are un flux de vârf al lărgimii de bandă cu lucrul cu memoria de 30 GB/sec și GPU are un flux de vârf al lărgimii de bandă cu lucrul cu memoria de 150 GB/sec. Sistemul hibrid are un supracontrol adițional care necesită ca toți vectorii de intrare să fie transferați între memoria gazdă și memoria locală a GPU înainte și după ce a fost invocat kernel-ul. Sistemul hibrid are o interfață cu acces direct la memorie (DMA) cu lărgimea de bandă de 10 GB/sec și o latență medie de 10 ms. Presupuneți că atât procesorul vectorial cât și GPU sunt mărginiți ca performanță de lărgimera de bandă a memoriei. Calculați timpul de execuție cerut de ambele calculatoare utilizate pentru această aplicație.

- 11) Secțiunea 4.5 discută despre operația de reducere care reduce un vector la un scalar prin aplicarea repetată a unei operații. O reducere este un tip special de recurență a unei bucle. Un exemplu este prezentat în continuare:

```
dot 0.0;  
for (i == 0; i < 64; i++) dot = dot + a[i] * b[i];
```

Un compilator cu vectorizare poate aplica o transformare denumită *scalar expansion*, care expandează dot într-un vector și împarte bucla astfel încât poate fi realizată înmulțirea cu o operație pe vector, lăsând reducția ca o operație scalară separată:

```
for (i == 0; i < 64; i++) dot[i] = a[i] * b[i];
```

```
for (i == 0; i < 64; i++) dot[0] = dot[0] + dot[i];
```

Așa cum s-a menționat în secțiunea 4.5 dacă permitem ca adunarea în virgulă flotantă asociativă, există câteva tehnici disponibile pentru reducția paralelizării.

- Una din tehnici este denumită dublarea recurenței, care adaugă secvențe pentru progresul vectorilor scurți (de exemplu doi vectorii cu 32 de elemente, apoi doi vectori cu 16 elemente, etc). Arătați cum va arăta codul C pentru execuția celeie de a doua bucle în acest caz.
- În unele procesoare vectoriale, elementele individuale di regiștrii vector sunt adresabili. În acest caz, operandii pentru o operație pe vectori pot fi două părți diferite ale aceluiași registru vector. Aceasta permite o altă soluție pentru reducție denumită sumă parțială (*partial sum*). Ideea este de a reduce vectorul la m sume unde m este latența totală prin unitatea funcțională pentru vectori, incluzând timpul de citire și de scriere a operandului. Presupuneți că regiștrii vector ai VMIPS sunt adresabili (de exemplu poți iniția o operație pe vector cu operandul V(16), indicând că operandul de intrare începe cu elementul 16). De asemenea, presupuneți că latența locală pentru adunări, incluzând citirea și scrierea rezultatului este de 8 cicli. Scrieți o secvență de cod VMIPS care reduce conținutul lui V1 la 8 sume parțiale.
- Când se realizează o reducție pe un GPU, un fir este asociat cu fiecare element al vectorului de intrare. Primul pas este, pentru fiecare fir să scrie valoarea corespunzătoare în memorie partajată. În următorul pas, fiecare fir intră într-o buclă în care adună fiecare pereche de valori de intrare. Aceasta reduce numărul de elemente la jumătate după fiecare iterație, ceea ce înseamnă că numărul de fire active de asemenea se vor reduce la jumătate după fiecare iterație. Cu scopul de a minimiza performanțele reducției, numărul de locații deplin populate (warps) trebuie maximizate în timpul rulării buclei. Cu alte cuvinte, firele active trebuie să fie continue. De asemenea, fiecare fir trebuie să indexeze o arie partajată într-un astfel de mod pentru a evita conflictele de banc în memoria partajată. Următoarea buclă violează numai prima opțiune din acest ghid și utilizează de asemenea operatorul modulo care este foarte scump pentru GPU-uri:

```
unsigned int tid = threadIdx.x;
for(unsigned int s=1; s < blockDim.x; s *= 2) {
    if ((tid % (2*s)) == 0) {
        sdata[tid] += sdata[tid + s];
    }
    __syncthreads();
}
```

Rescrieți bucle pentru a îndeplini aceste ghiduri și a elimina utilizarea operatorului modulo. Presupuneți că există 32 de fire per warp și un conflict de banc apare atunci când două sau mai multe fire de la același warp referențiază un index a cărui modulo cu 32 este egal.

- 12) Următorul kernel realizează o porțiune din metoda cu diferențe finite în domeniul timp (FDTD) pentru calculul ecuațiilor lui Maxwell într-un spațiu tridimensional, partea 1 a benchmark-ului SPEC06fp:

```
for (int x=0; x<NX-1; x++) {  
    for (int y=0; y<NY-1; y++) {  
        for (int z=0; z<NZ-1; z++) {  
            int index = x*NY*NZ + y*NZ + z;  
            if (y>0 && x >0) {  
                material = IDx[index];  
                dH1 = (Hz[index] - Hz[index-incrementY])/dy[y];  
                dH2 = (Hy[index] - Hy[index-incrementZ])/dz[z];  
                Ex[index] = Ca[material]*Ex[index]+Cb[material]*(dH2-dH1);  
            }  
        }  
    }  
}
```

Presupuneți că dH1, dH2, Hy, Hz, dy, dz, Ca, Cb, și Ex sunt arii de valori în virgulă în virgulă simplă. Presupuneți că IDx este o arie de întregi fără semn.

- Care este intensitatea aritmetică a acestui kernel?
 - Este kernel-ul maleabil cu execuția SIMD a vectorului?
 - Presupuneți că acest vector se execută pe un procesor care are lărgimea de bandă pentru memorie de 30 GB/sec. Va fi acest kernel mărginit de memorie sau de calcul?
 - Dezvoltați un model de tip roofline pentru acest procesor, presupunând că are un vârf de putere de calcul de 85 FFLOP/sec.
- 13) Presupuneți o arhitectură GPU care conține 10 procesoare SIMD. Fiecare SIMD instrucțiune are o lărgime de 32 și fiecare procesor SIMD conține 8 plane pentru aritmetica în virgulă simplă și instrucțiunile load/store ceea ce înseamnă că o instrucțiune SIMD nedivergentă poate produce 32 de rezultate la fiecare 4 cicli. Presupuneți un kernel care are ramificații divergente care cauzează o medie de 80% de fire ca să fie active. Presupuneți că 70 % din toate instrucțiunile SIMD executate sunt în virgulă flotantă singură precizie și 20 % sunt load/store. Deși nu sunt acoperite toate latențele memoriei, presupuneți că lansarea unei instrucțiuni SIMD se face cu o rată de 0,85. Presupuneți că GPU are un ceas de 1,5 GHz.
- Calculați fluxul, în GFLOP/sec, pentru acest kernel pe acest GPU.
 - Presupuneți că aveți următoarele alegeri:
 - Creșterea numărului de plane pentru simplă precizie la 16
 - Creșterea numărului de procesoare SIMD la 15 (presupuneți că această creștere nu afectează alte metrice privind performanțele și că se scalează codul cu procesoarele suplimentare)
 - Adăugați un cache care efectiv va reduce latența memoriei cu 40%, care va crește viteza de lansare a instrucțiunilor în execuție cu 0,95.Care este creșterea de viteză a fluxului de calcul pentru fiecare dintre aceste îmbunătățiri?

- 14) În acest exercițiu, vom examina câteva bucle și potențialul lor pentru paralelizare.

- Are următoarea buclă o dependență de tipul transport din buclă (loop-carried)?

```
for (i=0; i<100; i++) {
    A[i] = B[2*i+4];
    B[4*i+5] = A[i];
}
```

- b) În bucla următoare, găsiți toate dependențele reale, dependențele de ieșire și antidependențele. Eliminați prin redenumire dependențele de ieșire și anti dependențele.

```
for (i=0; i<100; i++) {
    A[i] = A[i] * B[i]; /* S1 */
    B[i] = A[i] + c; /* S2 */
    A[i] = C[i] * c; /* S3 */
    C[i] = D[i] * A[i]; /* S4 */
}
```

- c) Considerați următoarea buclă:

```
for (i=0; i < 100; i++) {
    A[i] = A[i] + B[i]; /* S1 */
    B[i+1] = C[i] + D[i]; /* S2 */
}
```

Există dependențe între S1 și S2? Este această buclă paralelă? Dacă nu, arătați cum o facem paralelă?

- 15) Listează și descrie cel puțin patru factori care influențează performanțele unui nucleu GPU. Cu alte cuvinte, care este comportarea în timpul execuției care este cauzată de execuția codului kernel care cauzează reducția în utilizarea resurselor pe durata execuției kernel -ului.

- 16) Presupuneți un GPU ipotetic cu următoarele caracteristici:

- ✓ Viteza ceasului de 1,5 GHz
- ✓ Conține 16 SIMD procesoare, fiecare conținând 16 unități în virgulă flotantă simplă precizie.
- ✓ Are 100 GB/sec lărgimea de bandă off-chip.

Fără a considera lărgimea de bandă a memoriei, care este vârful fluxului pentru operațiile în virgulă flotantă simplă precizie pentru acest GPU în GFLOP/sec, presupunând că toată latența memoriei poate fi ascunsă? Este acest flux sustenabil fiind dată limitarea privind lărgimea de bandă a memoriei?