

Laborator 4 - ISA, Pipeline, Hazarde

[1. Arhitectura Setului de Instrucțiuni](#)

[1.1 RISC \(Reduced Instruction Set Computer\)](#)

[1.2 CISC \(Complex Instruction Set Computer\)](#)

[1.3 RISC vs CISC](#)

[2. Paralelism - Pipelining](#)

[3. Hazarde](#)

[3.1 Tipuri](#)

[3.2 Soluții](#)

În cadrul acestui seminar se vor studia elemente de bază din arhitectura procesoarelor. Arhitectura unui procesor este definită prin ISA și prin microarhitectură.

1. Arhitectura Setului de Instrucțiuni

Arhitectura setului de instrucțiuni (ISA - Instruction Set Architecture) este folosită pentru a abstractiza funcționarea internă a unui procesor. ISA definește "personalitatea" unui procesor: cum funcționează procesorul d.p.d.v. al programatorului, ce fel de instrucțiuni execută, care este semantica acestora, în timp ce microarhitectura se referă la cum este implementată ISA, pipeline-ul de instrucțiuni, fluxul de date dintre unitățile de execuție (dacă sunt mai multe), registre și cache. Astfel, o ISA poate fi implementată de diferite microarhitecturi: la ARM, ARMv6 este un exemplu de ISA, și este implementată de 4 procesoare, la Pentium numele ISA este IA-32 și este implementată de diferite procesoare produse de Intel, AMD, Via, Transmeta. Astfel, ISA este cea mai importantă parte a design-ului unui procesor; alte aspecte cum sunt interacțiunea cu cache-ul, pipeline-ul, fluxul de date în procesor putând fi schimbate de la o versiune la alta a procesorului.

La ora actuală există două filozofii de design pentru un procesor: Complex Instruction Set Computer (CISC) și Reduced Instruction Set Computer (RISC). În afară de acestea există și ISA-uri pentru procesoare specializate, cum sunt GPU-urile pentru plăci grafice și DSP-urile pentru procesare de semnal.

1.1 RISC (Reduced Instruction Set Computer)

Caracteristicile care definesc RISC sunt:

- puține instrucțiuni: ~100-200; complexitatea care este eliminată din ISA este de fapt transferată programatorului în limbaj de asamblare și compilatoarelor; codul conține mai multe instrucțiuni, și în multe cazuri (mai ales în trecut (anii '70, '80)) acest fapt constituie o problemă
- instrucțiunile sunt de obicei codificate pe lungime fixă; permite o implementare mai simplă a UC

- execuția instrucțiunilor într-un singur ciclu (folosind pipeline)
- **arhitectură load-store** - numai instrucțiunile de tip LOAD și STORE lucrează cu memoria
- un număr mai mare de registre (din cauza load-store): 32-64; familii de arhitecturi: ARM, AVR, AVR32, MIPS, PowerPC (Freescale, IBM), SPARC (SUN)

Exemplu: încărcarea unui cuvânt de date din memorie într-un registru la arhitectura AVR32:

```
ld.w Rd, Rp++
```

Rd - registrul destinație

Rp - registrul ce conține adresa cuvântului; aceasta este incrementată după copierea conținutului

1.2 CISC (Complex Instruction Set Computer)

Arhitecturile de tip CISC pun accentul pe hardware având instrucțiuni complexe ce reduc dimensiunea codului și simplificând partea de software. Cele mai importante caracteristici sunt:

- instrucțiunile pot accesa memoria, având deja încorporate load-ul și store-ul
- varietate mare de moduri de adresare a memoriei
- instrucțiunile au dimensiuni variabile și necesită mai mulți cicli de ceas pentru a fi executate
- necesită un nivel suplimentar între hardware și ISA, controlul microprogramului - instrucțiunile sunt implementate în microcod.
- familii de arhitecturi: IA-32, x86-64

Exemplu:

La IA-32: `add [ebx], eax` (se aduna operandul de la adresa data de registrul ebx cu cel din registrul eax si rezultatul se stocheaza la adr primului operand)

La o arhitectura RISC:

```
load r1, $10
```

```
load r2, $11
```

```
add r1, r2
```

```
store $10, r1
```

\$10, \$11 sunt locații de memorie (adrese), numele instrucțiunilor sunt generice, nu aparțin unei arhitecturi anume

Pentru o arhitectura RISC am avea întâi de încărcat operanzii din memorie folosind instrucțiuni de tip LOAD și apoi să apelăm instrucțiunea pentru înmulțire, în final stocând rezultatul în memorie folosind o instrucțiune de tip STORE.

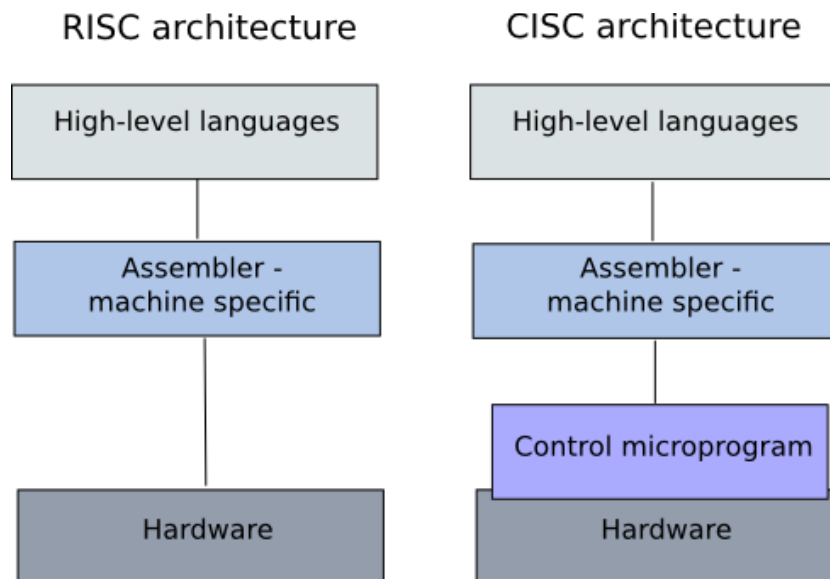


Fig 1.

1.3 RISC vs CISC

Regula 80-20: 80% din timp se folosesc 20% din instrucțiunile unui procesor. Multe instrucțiuni CISC sunt nefolosite de către programatori, iar majoritatea instrucțiunilor complexe pot fi sparte în mai multe instrucțiuni simple.

Performanța:

$$\frac{\text{timp}}{\text{program}} = \frac{\text{timp}}{\text{ciclu}} * \frac{\text{instrucțiuni}}{\text{program}} * \frac{\text{cicli}}{\text{instrucțiune}}$$

Ambele abordări fac compromisuri pentru a minimiza una dintre aceste fracții: în cazul RISC, se preferă un număr mai mare de instrucțiuni pe program pentru a micșora numărul de cicli pe instrucțiune, în timp ce la CISC este invers.

Totuși, lucrurile nu stau atât de simplu. Deși arhitecturile RISC au fost devansate în anii 80 de către CISC datorită lipsei software-ului și al unui segment de piață, precum și datorită prețului ridicat al memoriei, o dată cu progresele tehnologice, procesoarele cu ISA RISC au monopolizat domeniul embedded. În plus, cele două arhitecturi au început să se apropie, majoritatea procesoarelor actuale având o arhitectura hibridă, de exemplu Intel Pentium având un interpretor CISC peste un nucleu RISC.

RISC	CISC
costuri mici de fabricatie, mai puține tranzistoare	hardware mai complex, costuri mai mari
consum redus de energie	consum mai mare
mai multe registre	mai mult hardware, controlul microprogramului
mai multe linii de cod per program	puțin cod

2. Paralelism - Pipelining

Chiar dacă un program este în mod tradițional interpretat secvențial, instrucțiune cu instrucțiune, o parte dintre acestea nu sunt dependente între ele și pot fi executate în paralel. Această metodă de a extrage paralelism la nivelul unui singur flux de instrucțiuni de numește paralelism la nivel de instrucțiune sau Instruction-level parallelism (ILP).

O metodă de a implementa ILP este banda de asamblare (pipeline). Banda de asamblare presupune că atunci când o instrucțiune i este executată, instrucțiunea $i+1$ este decodificată și instrucțiunea $i+2$ este citită, astfel fiecare subsistem din procesor este ocupat la un moment dat.

Modelul de bază pentru un pipeline la o arhitectura RISC este cel din 4 etape (Fig. 2): Instruction Fetch, Instruction Decode, Execute, Register Write Back. În stagiul Fetch se aduce instrucțiunea din memorie, de la adresa dată de contorul program (PC - *program counter*). În stagiul Decode se decodifică instrucțiunea, în Execute se execută efectiv operația de către unitatea aritmetico-logică, iar în Writeback se stochează rezultatele în regiștrii sau memorie. Pentru CISC, avem de obicei un număr mai mare de etape (>12).

Banda de asamblare **superpipeline** este mai lungă decât banda obișnuită (cele 4 stagii F, D, E, W sunt împărțite la rândul lor în mai multe stagii, Fig. 3). Teoretic cu cât banda de asamblare este mai segmentată, mai multe instrucțiuni pot fi executate în paralel și, fiecare stagiu făcând mai puține operații, se poate mări frecvența ceasului. Însă, alți factori caracteristici ai design-ului pot scădea performanțele și nu întotdeauna un procesor cu superpipelining este mai bun decât unul cu bandă de asamblare normală. Pentru o arhitectură și un set de instrucțiuni date există un număr optim de etape ale benzii de asamblare. Peste acest număr se reduce performanța. Pentru procesoarele RISC, unde instrucțiunile se execută într-un singur ceas, această durată de execuție este dată de instrucțiunea cea mai lungă.

Banda de asamblare mărește productivitatea (throughput-ul), adică numărul de instrucțiuni executate în unitatea de timp, însă timpul de execuție (latența) a unei instrucțiuni nu se micșorează, ci se poate chiar mări.

Pentru a paraleliza și mai mult execuția instrucțiunilor, procesoarele **superscalare** introduc benzi de asamblare independente ce funcționează în paralel. Se mărește astfel numărul de resurse și complexitatea procesorului, însă mai multe instrucțiuni pot fi aduse, decodificate și executate în același timp, după cum este ilustrat și în Fig. 4.

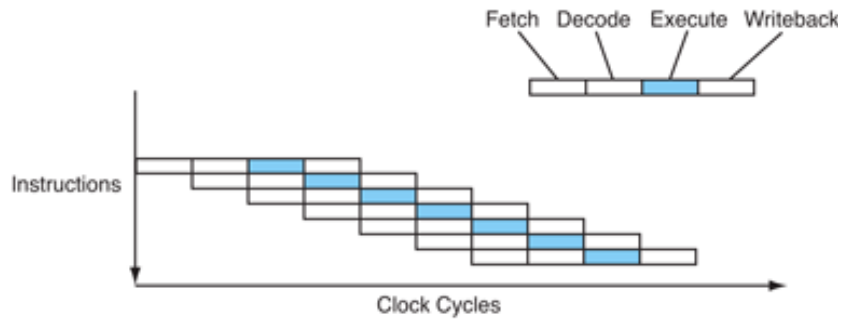


Fig 2. Execuția instrucțiunilor în banda de asamblare.

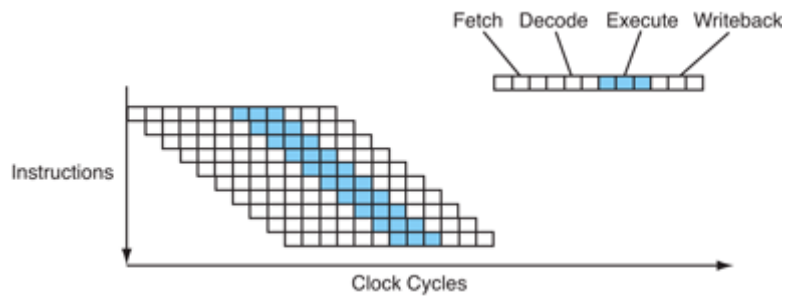


Fig 3. Execuția instrucțiunilor în superpipeline

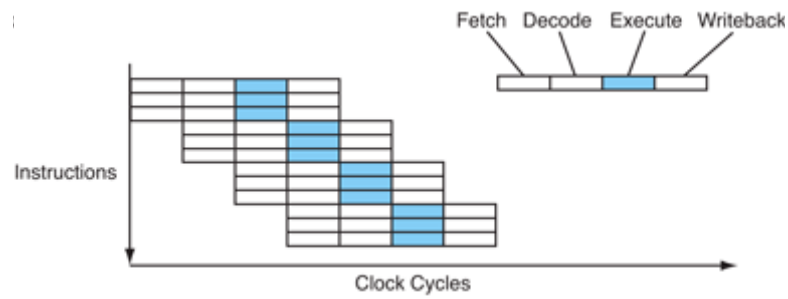


Fig 4. Execuția instrucțiunilor într-un procesor superscalar

Sursa imaginilor este [aici](#).

3. Hazarde

În cazul unui procesor cu bandă de asamblare, pentru a obține performanță maximă, este necesar ca o instrucțiune să fie prezentă în fiecare stadiu al pipeline-ului. Acest lucru nu este posibil însă datorită faptului că instrucțiunile pot avea diferite dependențe între ele ceea ce face ca execuția lor în pipeline să ducă la rezultate nedorite, situație care poartă numele de hazard.

3.1 Tipuri

Există trei tipuri de dependențe care pot conduce la hazarde:

1. **Structurale**: când două instrucțiuni aflate în execuție în pipeline doresc acces la aceeași structură hardware în același timp.

Exemplu:

```
add r1, r2
load r3, [r4 + offset]
```

- ambele instrucțiuni necesită acces la unitatea de adunare: prima pentru a executa o adunare, iar a doua pentru a calcula adresa de la care se face citirea din memorie.

2. **De date**: când două instrucțiuni folosesc același registru sau aceeași locație de memorie și cel puțin una dintre ele este o scriere

• **RAW (Read after Write)**: instrucțiunea I_{i+1} citește o locație de memorie scrisă de instrucțiunea I_i

Exemplu:

```
add r1, r2
add r3, r1
```

- valoarea lui `r1` este disponibilă în a doua instrucțiune abia după ce prima instrucțiune a scris rezultatul

• **WAR (Write after Read)**: instrucțiunea I_{i+1} scrie o locație de memorie care este citită de instrucțiunea I_i

Exemplu:

```
add r2, r1
add r1, r3
```

- valoarea lui `r1` nu poate fi modificată în a doua instrucțiune până când prima instrucțiune nu a terminat citirea lui

• **WAW (Write after Write)**: instrucțiunea I_{i+1} scrie o locație de memorie care este scrisă și de instrucțiunea I_i

Exemplu:

```
mov r1, r2
mov r1, r3
```

- este necesar ca cele două scrieri în `r1` să se execute în ordinea din program pentru a nu modifica valoarea văzută de instrucțiunile următoare

3. **De control:** când o instrucțiune de salt (branch, jump, call, ret) este executată; instrucțiunile de salt pot modifica PC (contorul program), sărind la o altă instrucțiune decât cea imediat următoare; în acest caz adresa următoarei instrucțiuni care trebuie citită nu se cunoaște până când instrucțiunea de salt nu este executată

Exemplu:

```
if (a == b)           (i1)
    d = a+b           (i2)
    c = 1              (i3)
else
    d = a * b         (i4)
```

În banda de asamblare se aduc din memorie la rand instrucțiunile `i1,i2,i3`, însă dacă în urma execuției instrucțiunii `if` (tradusă în assembler printr-o instrucțiune de comparare și apoi una de salt condiționat (`jmp - jump`)) se sare la ramura de `else` apare un hazard de control. Instrucțiunile `i2,i3` nu mai trebuie executate, și în schimb trebuie adusă instrucțiunea `i4`, deci se face 'flush' la pipeline și se aduce instrucțiunea `i4`, pierzându-se astfel din eficiența prelucrării în paralel a instrucțiunilor.

3.2 Soluții

- O soluție simplă ar fi adăugarea de către programator sau compilator sau automat, de către UC a procesorului, a unor instrucțiuni `nop` pentru întârzierea execuției instrucțiunilor cu dependențe; așa numitele *pipeline bubbles* sau *pipeline stalls*.
- **Forwarding** - dacă unul sau mai mulți operanzi ai unei instrucțiuni depind de o instrucțiune aflată înainte în pipeline, rezultatul acesteia poate fi forwardat către stagiile anterioare. Aceasta soluție necesită hardware suplimentar pentru logica de control.

Exemplu:

```
instr k:           add r1,r2
instr k+1:         add r1,5
```

Rezultatul obținut pentru instrucțiunea `k` în unitatea aritmetico-logică, în stagiul de Execute este forwardat stagiului Decode al instrucțiunii `k+1`.

- **Out-of-order execution(OOE)** este o tehnică de planificare dinamică(dynamic scheduling) prin care instrucțiunile sunt executate în altă ordine decât cea a programului, însă rezultatele sunt ordonate ca și cum ar fi fost executate normal. Prin această rearanjare, se evită adăugarea de întârzieri, intrând în pipeline instrucțiuni fără dependențe de date între ele. Aceasta soluție necesită hardware suplimentar(registre suplimentare, buffers pentru cozi de așteptare etc), iar algoritmul cel mai cunoscut este cel al lui [Tomasulo](#).

- Deoarece registrele generale sunt limitate, iar codul unui program îi refolosește des, generând hazarde de date, putem folosi tehnica de **register renaming**. Aceasta minimizează apariția hazardelor, prin adaugarea unor registre suplimentare, ce nu sunt vizibile programatorului. Redenumirea registrelor rezolvă hazardele de date WAR, WAW însă nu și RAW, fiind folosită și în algoritmul Tomasulo de planificare OOE.

Exemplu:

```
instr k      add r1,r2,r3
instr k+1    mul r3,r2,r1
instr k+2    sub r1,r2,r4 //hazard WAW intre instr k si k+2
```

După redenumirea registrelor:

```
add l1, l2, l3
mul l5, l2, l4
add l6, l2, l7
```

- În cazul hazardelor de control, putem amâna aducerea instrucțiunilor din branch prin adăugarea de întârzieri în pipeline sau intercalarea altor instrucțiuni independente de acel branch, însă soluția cea mai bună este predicția ramificațiilor (**branch prediction**) care minimizează numărul de cicli pierduți. Această predicție poate fi statică, făcută de către compilator prin optimizări de tipul *loop unrolling* (de exemplu în bucla *for*, a cărei condiție este neîndeplinită doar o dată pentru a se ieși din buclă) sau profile-driven (se creează profile ale programelor în urma mai multor execuții). Predicția dinamică se face de către procesor prin intermediul unor circuite (*branch predictors*) și în general se bazează pe menținerea unor tabele cu informații despre fiecare dintre ramificații.

Referințe

Patterson, David A.; Hennessey, John L, *Computer Architecture: A Quantitative Approach*, 4th Edition

[RISC vs CISC](#)

[RISC vs CISC wars](#)

[Pipeline. IPL. Hazarde](#)