

Tehnici de prelucrare a instructiunilor de transfer al controlului

Prof.dr.ing. Gheorghe DODESCU
 Catedra de Informatica Economica, A.S.E. Bucuresti
 Lect.ing. Bogdan OANCEA
 Universitatea "Artifex" Bucuresti

Instructiunile de transfer al controlului au pus mari probleme proiectantilor de procesoare datorita penalitatilor mari introduse de acestea. Folosirea tehnicii de prelucrare pipeline si aparitia procesoarelor superscalare a impus gasirea unor modalitati de tratare eficiente a acestui tip de instructiuni. Articolul prezinta cele mai des utilizate tehnici de predictie a salturilor intalnite in procesoarele moderne.

Cuvinte cheie: procesor, instructiuni de salt, tehnici de predictie.

Instructiunile de salt reprezinta un caz special de instructiuni. Toate procesoarele moderne dispun de tehnici speciale pentru prelucrarea acestor instructiuni. Vom arata in primul rand de ce se impun astfel de tehnici speciale pentru instructiunile de salt. Pentru aceasta vom pleca de la constatarea ca in prezent toate procesoarele utilizeaza metoda de prelucrare in banda de asamblare (pipeline). Consideram un caz simplificat de procesor pipeline in care avem doar 4 nivele: aducerea instructiunii din cache (Fetch), decodificarea ei (D), executia instructiunii (EX) si in final scrierea rezultatului in

registrul destinatie (WB). Trebuie mentionat ca acest caz nu este departe de realitate, multe procesoare RISC avand implementat un astfel de pipeline la nivelul unitatii de prelucrare a instructiunilor. Presupunem ca avem de executat cu acest procesor fluxul de instructiuni I1, I2, Is, I3, I4, I5...It, unde Ii sunt instructiuni obisnuite, Is este o instructiune de salt neconditionat iar It este instructiunea unde se efectueaza saltul. Mai facem ipoteza ca la sfarsitul fazei de executie a instructiunii de salt este disponibila adresa la care se va efectua saltul. Putem urmari modul de executie al acestor instructiuni in figura 1.

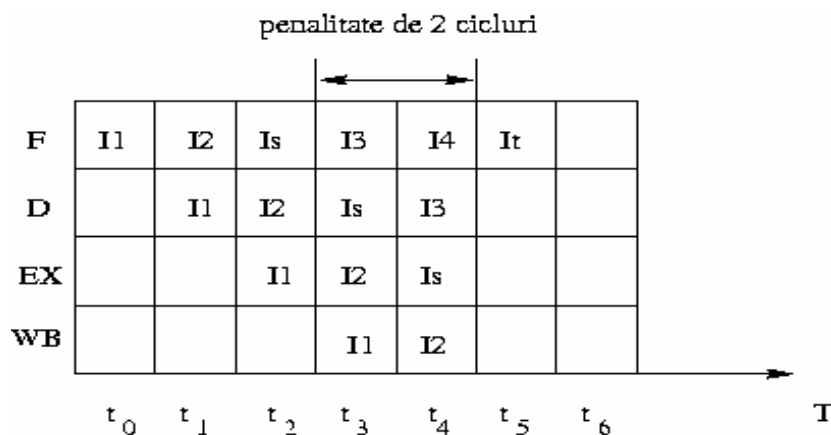


Fig. 1. Executia unei instructiuni de salt

Din momentul in care instructiunea de salt Is este in faza F si pana cand instructiunea de la adresa

de salt (It) intra in aceeaasi faza (F) trec doua cicluri de ceas care practic sunt irosite. Aceste

cicluri pierdute se numesc cicluri de penalitate la executarea unei instructiuni de salt. Dupa cum se poate usor observa din figura, instructiunile I3 si I4 care au fost deja admise în pipeline devin inutile si deci s-a pierdut timp cu prelucrarea lor. Situatiia se complica si mai mult daca banda de asamblare are mai multe nivele cum este cazul procesoarelor DEC Alfa 21164 si Pentium Pro care au 12 nivele sau MC68060 cu 9 astfel de nivele. O problema în plus apare în cazul instructiunilor de salt conditionat, unde evaluarea conditiei mai cere câtiva cicli procesor ceea ce duce la cresterea penalitatii. Pe de alta parte, în momentul executarii instructiunii de salt conditionat s-ar putea sa nu fie înca evaluata expresia care sta la baza conditiei, asa cum se întâmpla în procesoarele superscalare care posedea mai multe unitati functionale ce lucreaza independent. În acest caz spunem ca avem o instructiune de salt conditionat nerezolvata. Toate aceste aspecte pot conduce la o scadere sensibila a performantelor. Spre exemplu procesorul Intel 80386 are o penalitate de 8 cicli pentru instructiunile la care se executa salturile si 2 cicli pentru cele la care nu se executa saltul. Procesoarele moderne, care dispun de metode

a) Metoda "delayed branching" a fost printre primele încercari de a reduce penalitatile care apar în cazul instructiunilor de salt. Ideea de baza este urmatoarea: daca ne amintim situatia prezentata în figura 1, constam ca imediat dupa ce instructiunea de salt a intrat în pipeline se pierde un numar de cicli. Putem reutiliza acesti cicli pierduti daca procesorul dispune de posibilitatea de a lansa în executie mai întâi instructiunea de salt apoi cea care o precede.

Concret, sa consideram secventa de instructiuni:

```
...
add R1,R2,R3; R1<-R2+R3
jmp @1
```

...

@1: mul R5,R6,R4; R5<-R6*R4

speciale de tratare a unor astfel de situatii au reusit sa reduca aceste penalitati pâna la 0 cicli în anumite cazuri sau 2...4 cicli în cel mai rau caz.

Atentia acordata acestui tip de instructiuni este impusa si de faptul ca instructiunile de salt conditionat sau neconditionat apar foarte des în programe. Statisticile demonstreaza ca în programele de uz general ponderea instructiunilor de salt este de aproximativ 20-25% în timp ce în programele stiintifice este de 5-10%. Dintre acestea, aproximativ 75% dupa unii autori sau 83% dupa altii sunt instructiuni în care se efectueaza saltul iar în restul de 25% respectiv 17% din cazuri se continua pe calea secventiala. Procentul mare de instructiuni în care se efectueaza saltul conduce la cautarea unor solutii arhitecturale care sa minimizeze penalitatile în acest caz.

Putem rezuma tehnicile de tratare a instructiunilor de salt la urmatoarele :

- a) folosirea metodei "delayed branching"
- b) blocarea procesarii instructiunii de salt conditional nerezolvat
- c) executarea speculativa a instructiunilor de salt conditional nerezolvat
- d) executarea ambelor cai a unei instructiuni de salt (multiway branching)

Facem ipoteza ca adresa de salt va fi disponibila la sfârșitul fazei de decodificare, iar instructiunea de adunare se poate executa într-un singur ciclu. Atunci, daca înaintea instructiunii *add R1,R2,R3* se va executa instructiunea *jmp @1* se va putea recupera ciclul pierdut datorat saltului, dupa cum se poate observa si în figura 2.

Aceasta metoda presupune un ajutor din partea compilatorului care va trebui sa gaseasca o instructiune independenta care sa poata fi planificata pentru executie în modul prezentat în figura 2. Cu toate ca unele arhitecturi precum MIPS-X, MC88100, HP PA folosesc aceasta metoda, în prezent tendinta este de a se renunta la ea si de a folosi metoda de executie speculativa a instructiunilor.

b) Blocarea procesarii instructiunilor de salt nerezolvate este o metoda care s-a utilizat în primele procesoare pipeline (i386, MC68020) si presupune asteptarea rezolvarii conditiei de salt si abia apoi lansarea în executie a instructiunii de

transfer a con-trolului. Datorita ineficientei, aceasta metoda a fost abandonata.

c) Executia speculativa este deocamdata cea mai avansata tehnica de procesare a instructiunilor de salt.

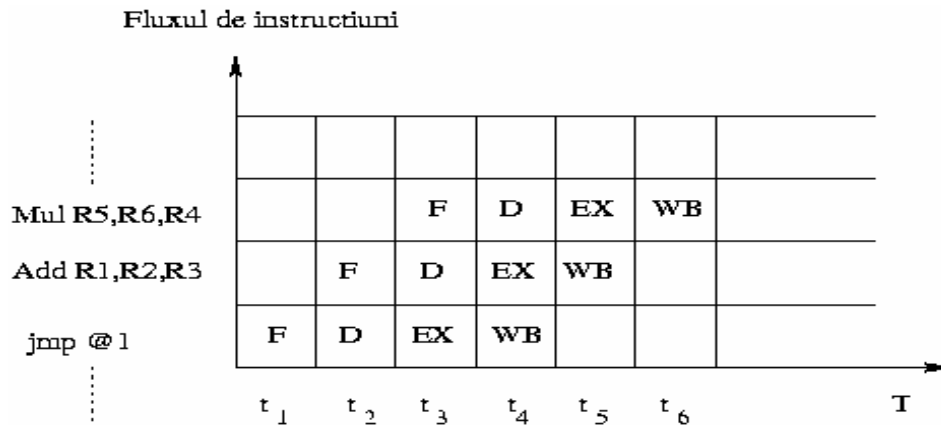


Fig. 2. Metoda “delayed branching”

Pentru a permite tratarea eficienta a instructiunilor de salt trebuie ca ele sa fie detectate din timp. Procesoarele din linia DEC Alfa precum si Power PC 601 detecteaza si decodifica salturile în paralel cu decodificarea obisnuita pe când Power1, Power2 si PowerPC 603 folosesc o tehnica denumita *look-ahead branch detection* în care instructiunile de salt sunt detectate si decodificate înaintea instructiunilor obisnuite. Astfel, daca pentru decodificare si planificare spre executie se aduc n instructiuni/ciclu din memoria cache, pentru detectarea instructiunilor de salt, vor fi inspectate ultimele $n+k$ intrari din memoria cache. Procesorul Power1 aduce 4 instructiuni/ciclu din memoria cache pentru decodificare, dar pentru detectarea salturilor inspecteaza ultimele 5 intrari din cache, în timp ce Power2 decodifica 6 instructiuni/ciclu si verifica ultimele 8 intrari din cache pentru detectarea salturilor.

Cea mai avansata metoda de detectare a instructiunilor de ramificare presupune ca la aducerea instructiunilor din memoria cache în buffer-ul de instructiuni sa fie detectat daca urmatoarea instructiune este de salt sau nu, iar în caz afirmativ se va aplica o metoda de prezicere

a saltului pentru a aduce din cache instructiunea de la adresa de salt sau cea care ar urma în ordinea secventiala. Aceasta metoda este utilizata de PowerPC 604, 620, R8000, PA8000.

Odata detectata o instructiune de salt, procesorul va încerca sa ghiceasca pe care ramura se va continua programul si în continuare se vor executa instructiunile de pe ramura prezisa. În momentul în care rezultatul evaluarii conditiei de salt devine disponibil, se va verifica daca prezicerea a fost corecta sau nu. Daca aceasta s-a dovedit a fi corecta atunci programul continua în mod normal, dar daca prezicerea a fost eronata, toate instructiunile executate speculativ vor fi anulate si executia continua pe calea corecta.

Metodele de predictie a instructiunilor de salt se pot clasifica în metode de predictie *fixe* si metode de predictie *reale*.

Metodele fixe se caracterizeaza prin faptul ca instructiunile de salt vor fi tratate întotdeauna în acelasi fel: fie se prezice ca saltul se va efectua mereu, cum este cazul procesorului MC68040, fie se prezice ca nu se efectueaza saltul niciodata ci se continua pe calea secventiala, metoda aplicata la INTEL 486, SuperSparc, Power1,

Power2, DEC Alfa 21064. Aceste metode sunt mai simplu de implementat, dar acuratetea predictiei s-a dovedit a fi necorespunzatoare. De aceea majoritatea procesoarelor din ultima generatie folosesc metode de predictie reala.

Metodele de predictie reale pot fi de doua tipuri: metode *statice* si metode *dinamice*.

Metodele statice le putem separa mai departe dupa informatia care sta la baza prezicerii saltului. Distingem astfel 3 metode statice:

Metode bazate pe codul operatiei. Unele procesoare verifica un anumit bit din codul unei operatii si în functie de acesta va lua decizia. Un exemplu în acest sens este procesorul MC88110 care seteaza bitul 21 din instructiunile de comparare când conditia din instructiune este $\neq 0$, > 0 sau ≥ 0 . În aceasta situatie se va efectua saltul. Bitul 21 se va reseta atunci când conditia este $= 0$, < 0 sau ≤ 0 si saltul nu se va efectua.

Metode bazate pe directia în care se efectueaza saltul: daca saltul se efectueaza înainte (adresa de salt mai mare decât a instructiunii curente) saltul nu se efectueaza, iar în cazul când adresa de salt este mai mica, deci saltul are loc înapoi se va efectua saltul. Procesoarele care implementeaza aceasta metoda sunt PowerPC 601, PowerPC 603, DEC Alfa 21064.

Metode bazate pe un bit din codul unei operatii care este setat sau nu de un compiler. În functie de o anumita constructie de limbaj, compilatorul poate prezice daca se va efectua

saltul sau nu. Exemple: PowerPC 601, PowerPC 603, HP PA8000.

Metodele dinamice iau în considerare istoria executiei instructiunilor de salt. În functie de comportarea în trecut a unei instructiuni de salt se prezice si comportarea ei viitoare. Pentru a putea tine cont de istoria executiei instructiunilor de salt, trebuie memorata comportarea acestor instructiuni. Deosebim în acest fel scheme de predictie cu 1 bit, cu 2 biti respectiv cu 3 biti de memorie.

Cea mai simpla, schema cu 1 bit memoreaza istoria ultimei executii: bitul va avea valoarea 1 daca saltul s-a executat si 0 în caz contrar. În momentul predictiei se verifica valoarea acestui bit: daca este 1 se va prezice saltul, daca este 0 se continua pe calea secventiala. În momentul în care conditia care a stat la baza saltului este evaluata valoarea bitului se va actualiza în mod corespunzator.

Metoda cu 2 biti poate fi asemanata cu functionarea unei masini cu stari finite. Cu ajutorul a 2 biti putem memora 4 stari : 00, 01, 10 si 11. Daca starea curenta în momentul când trebuie sa se prezica saltul este 00 sau 01 decizia este de a nu efectua saltul, iar în starea 10 sau 11 decizia este de a efectua saltul. În momentul când conditia de salt a fost evaluata, cei 2 biti se vor actualiza conform diagramei de tranzitie a starilor din figura 3.

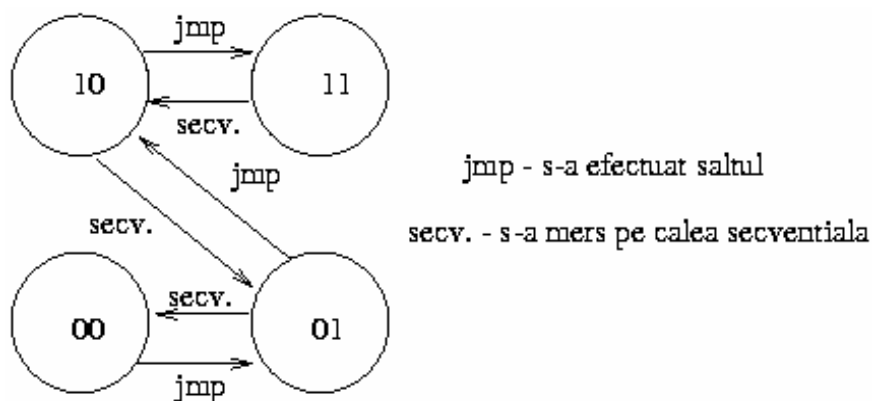


Fig. 3. Diagrama tranzitiilor de stare.

Simularile arata ca aceasta metoda are rezultate mult mai bune decât în situația folosirii unui singur bit. Probabilitatea de precizie corectă este în jurul valorii de 93% pentru 2 biti fata de numai 89% pentru 1 bit sau 68% în cazul metodelor statice.

Schema de predicție cu 3 biti presupune memorarea ultimelor 3 execuții ale unei instrucțiuni de salt. În momentul când trebuie

luată o decizie se numără bitii cu valoarea 1 și 0 iar decizia este luată pe baza majorității: mai mulți de 1 înseamnă salt, mai mulți 0 - se continuă ramura secvențială. Actualizarea acestor biti se face pe baza disciplinei FIFO.

În tabelul 1 sunt prezentate exemple de procesoare care implementează aceste tehnici de predicție a salturilor.

Predicție	Execută saltul întotdeauna	MC68040		
	fixa	Nu execută saltul niciodată	i486, R4000, SuperSparc, Power1, Power2, DEC Alfa 21064	
Predicție dinamica	Predicția statică	Bazată pe codul operației	PowerPC 601, PowerPC 603	
		Bazată pe deplasamentul adresei de salt	DEC Alfa 21064A	
		Bazată pe indicații date de compilator	PA8000, PowerPC 603	
	Predicția dinamica	1 bit	R8000	
		2 biti	MC68060, PowerPC604, 620, UltraSparc, DEC Alfa 21164, Pentium, R10000	
		3 biti	PA8000	

Tabelul 1. Clasificarea metodelor de predicție a salturilor

Execuția speculativă a instrucțiunilor poate conduce la următoarea situație: la întâlnirea unei instrucțiuni de transfer a controlului nerezolvate se prezice o cale de urmat, dar în timpul execuției instrucțiunilor de pe această cale apare o altă instrucțiune de salt. Problema este că astfel de instrucțiuni vor fi tratate speculativ. Cu cât sunt mai multe, cu atât vor fi mai mari performanțele procesorului, dar în cazul unei predicții greșite cu atât va fi mai grea sarcina de refacere a stării inițiale. Cele mai multe procesoare acceptă tratarea speculativă doar a unei singure instrucțiuni de salt (DEC Alfa 21064, 21064A, PowerPC 601, 603) dar există și procesoare care permit mai multor instrucțiuni de salt să fie tratate speculativ: 2 instrucțiuni pentru Power2, 4 pentru PowerPC 620 și MIPS R10000 și chiar 6 pentru DEC Alfa 21164.

Pentru a face față situației în care s-a făcut o predicție greșită, majoritatea procesoarelor

dispun de facilități pentru a reduce penalitățile în astfel de cazuri. Metoda cea mai simplă este de a salva adresa de continuare secvențială în cazul în care s-a prezis efectuarea saltului sau de a începe precalcularea adresei de salt atunci când s-a prezis urmarea căii secvențiale (PowerPC 601, 603, 620). O metodă mai avansată include chiar salvarea câtorva dintre instrucțiunile de pe calea secvențială în prima situație sau a instrucțiunilor de la adresa de salt în cea de-a doua situație. Acest lucru implică existența a mai multor buffere de instrucțiuni (Power1, Power2, Pentium). După cum am arătat mai sus, în aproximativ 75%-83% din instrucțiunile de salt se efectuează saltul. De aceea trebuie luate măsuri speciale pentru a micșora timpul de acces la instrucțiunea de la adresa de salt. Modalitatea cea mai simplă este ca la apariția unei instrucțiuni de salt, adresa să fie calculată de un sumator și apoi să fie aduse

în cache instrucțiunile de la aceasta adresa. Există însă și metode mai performante:

- **BTAC (Branch Target Address Cache)**

A.I.S.	A.S.

BTAC

PA 8000, PowerPC620
Pentium, MC68060

A.I.S. = Adresa instrucțiunii de salt

A.S. = Adresa de salt

I.A.S. = Instrucțiunea de la adresa de salt

A.U. = Adresa instrucțiunii ce urmează instrucțiunii de la adresa de salt

a)

A.I.S.	I.A.S.	A.U.

BTIC

MC88110

b)

A.I.S.	... I ...	index

cache care conține
indexul următoarei
instrucțiuni

K5, UltraSparc

c)

Fig. 4. Metode de accesare a caii de salt

Metoda BTAC constă în existența unei memorii cache suplimentare în care se păstrează adresa unei instrucțiuni de salt precum și adresa la care se efectuează saltul (figura 4 a). Această memorie va fi accesată asociativ de fiecare dată când următoarea instrucțiune este una de salt, iar dacă se găsește o intrare corespunzătoare în BTAC atunci se salvează timpul pierdut la calcularea adresei de salt. Mărimea BTAC variază între 32 până la 4k intrări. Metoda BTIC este asemănătoare metodei precedente cu deosebirea că în loc de a se mentine în memoria cache specială adresa saltului, se mentine chiar instrucțiunea de la adresa de salt (figura 4 b).

Ultima metoda amintită aici, implică existența unui câmp suplimentar în memoria cache pentru instrucțiuni, câmp care va puncta către următoarea instrucțiune care trebuie adusă în buffer-ul de prefetch (figura 4 c).

Aceste tehnici pot ajuta chiar la luarea unei decizii cu privire la tratarea unei instrucțiuni de salt nerezolvată: o intrare în BTAC sau BTIC înseamnă că respectiva instrucțiune a mai fost întâlnită în cursul executării programului și să

- **BTIC (Branch Target Instruction Cache)**

- existența în memoria cache a unui indice care să arate instrucțiunea imediat următoare.

efectuat saltul (deci seamănă cu metoda de predicție cu 1 bit). Pentru a mari performanțele, unele procesoare aplică schema de predicție BTAC sau BTIC în combinație cu predicția cu 2 biți (PowerPC 604, 620, Pentium).

d) executarea ambelor cai ale unei instrucțiuni de transfer a controlului este o metodă foarte pretentioasă, mare consumatoare de resurse și deocamdată nu este aplicată decât în unele mașini VLIW.

Bibliografie

1. Dezso Sima, Terence Fountain, Peter Kacsuk, *Advanced Computer Architecture*, Addison-Wesley 1997.
2. Kai Hwang, Faye A. Briggs, *Computer Architecture and Parallel Processing*, McGraw-Hill, 1985.
3. Kai Hwang, *Advanced Computer Architecture*, McGraw-Hill, 1993.
4. Hans Peter Messmer, *The Indispensable PC Hardware Book*, Addison-Wesley, 1995.
5. Marius Marcu, *Microprocesoarele INTEL*, PC Report, iunie 1998.

6. Marius Marcu, Procesoarele AMD, PC Report, septembrie 1998.