

Korszerű számítógép architektúrák I.

Simon Péter ¹

2021. március 25.

¹Durczy Levente előadásai alapján

1. fejezet

Első előadás

2. fejezet

Függőségek

2.1. Bevezetés

A függőségek gátolják a párhuzamos végrehajtást.

2.2. Típusai

- adat
- vezérlés
- erőforrás

2.3. Adat függőségek

Probléma: az utasítás végrehajtáshoz egy előző utasítás eredményére van szükség.

2.3.1. Csoportosítása

Jellege szerint

- utasítás szekvenciában (lineáris feldolgozás)
 - valós függőség - nem teljesen megszüntethető (RAW - Read After Write)
 - * műveleti adatfüggőség
 - * behívási adatfüggőség
 - ál függőség - teljesen megszüntethető
 - * WAR - Write After Read
 - * WAW - Write After Write
- ciklusban

Operandus típusa szerint

- regiszter
- memória

2.3.2. Műveleti adatfüggőségek

Probléma felvetés: feltételezzük, hogy

- a processzor 3 operandusos utasításokat használ
- 4 fokozatú futószalagos végrehajtás van (Fetch, Decode, Execute, WriteBack).

Ezekkel a feltételekkel két számot szeretnénk összeszorozni, az eredményt pedig megduplázni. Az utasításaink:

```
; I1
MUL r3, r2, r1 ; r3 = r1 * r2
; I2
SHL r3 ; r3 * 2
```

Az utasítások végrehajtásának időbeli sorrendje:

	t ₁	t ₂	t ₃	t ₄	t ₅
I ₁	F _{MUL}	D _{r1, r2}	E _{MUL}	W/B _{r3}	
I ₂		F _{SHL}	D _{r3}	E _{SHL}	W/B _{r3}

A probléma, hogy I₂ végrehajtása során, a dekódolási fázisban (t₃ időpillanat) szükség lenne az r3 regiszter értékére, viszont az csak t₄ időpillanatban áll elő (I₁ végrehajtásának writeback fázisában). Tehát a futószalagos módszerrel párhuzamosított végrehajtás során műveleti adatfüggőség keletkezett, mivel az utasítások lehívása és végrehajtása között átfedés van. Ilyenkor a műveletek elakadnak.

Megoldás: egy speciális utasítás, a NOP (No Operand) használata az alábbi módon:

	t ₁	t ₂	t ₃	t ₄	t ₅	t ₆	t ₇
I ₁	F _{MUL}	D _{r1, r2}	E _{MUL}	W/B _{r3}			
I ₂		F _{SHL}	NOP	NOP	D _{r3}	E _{SHL}	W/B _{r3}

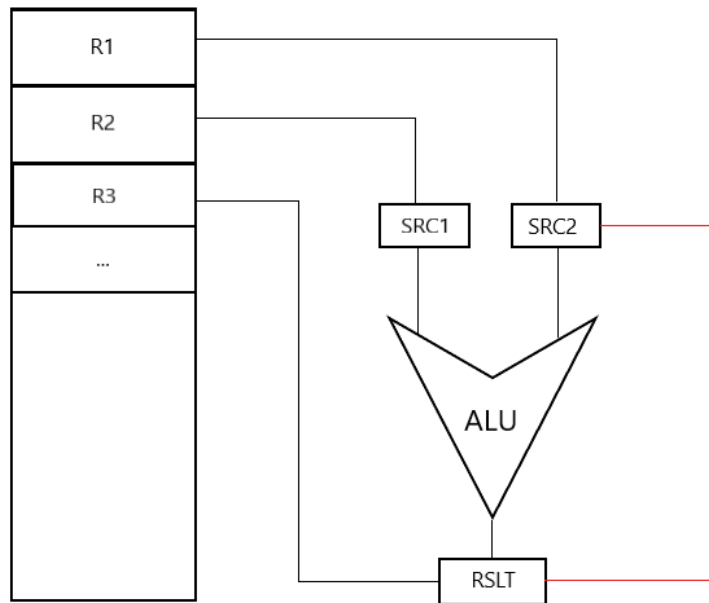
Következmény: a műveleteknek várakozniuk kell egymásra, két óraciklus késés keletkezik a futószalagon. Az ezeket követő utasításokhoz is be kell szűrni két NOP-ot, mivel a dekóder foglalt. Ezt a jelenséget teljesen nem lehet megszüntetni, viszont a fékező hatást csökkenthetjük.

Kezelés: operandus előrehozásával csökkenthető a fékező hatás, ez viszont extra hardvert igényel (hardveres, azaz dinamikus megoldás). Extra hardver nélkül csak szoftveresen, azaz statikusan kezelhetjük a problémát. Ilyenkor a compiler oldja meg a függőségek kezelését. Általában előnyösebb a dinamikus megoldás.

Dinamikus megvalósítás: az ALU-hoz tartozó rejtett regisztereket és az adatutakat a 2.1 ábra mutatja. Alapesetben a MUL utasítás végrehajtása során az adat az r₁ és r₂ regiszterekből az src₁, illetve src₂ regiszterekbe kerül, majd a művelet elvégzése után a rslt rejtett regiszteren keresztül visszaírásra kerül r₃-ba. Az adatút rövidítésének érdekében, extra hardver segítségével az rslt regiszter tartalmát közvetlenül visszavezethetjük az ALU egyik forrásregiszterébe. Ennek útja látható az ábrán pirossal. Ekkor az utasítások végrehajtása az alábbi módon valósul meg:

	t ₁	t ₂	t ₃	t ₄	t ₅	t ₆	t ₇
I ₁	F _{MUL}	D _{r1, r2}	E _{MUL}	W/B _{r3}			
I ₂		F _{SHL}	D	E _{SHL}	W/B _{r3}		

Mivel már t₃ időpillanatban is rendelkezésre áll az SHL utasítás operandusa, két óraciklussal hamarabb kezdhető meg a művelet végrehajtása. Ezzel megszüntettük a késést. Ezt a megoldást minden modern CPU használja.

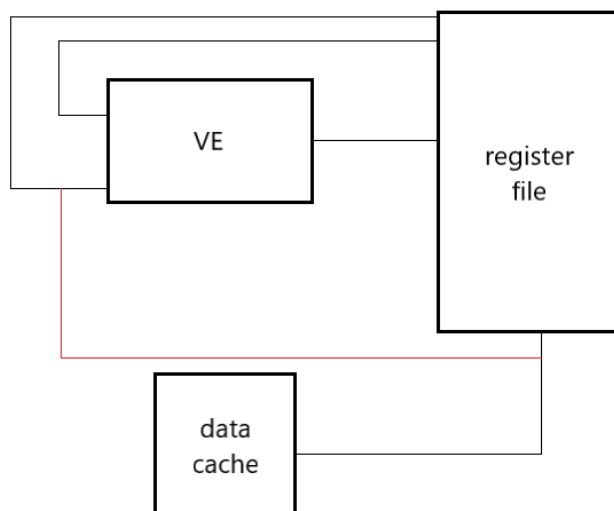


2.1. ábra. Az eredmény visszavezetése a forrás regiszterbe

2.3.3. Lehívási adatfüggőség

Probléma: a regiszterekben az operatív tárból (cache) töltjük be a szükséges adatokat, majd ezután a regiszterekből hívja le a végrehajtó egység (ALU). A cache elérése viszont sok időt vesz igénybe. Ennek látható az általános adatútja a 2.2 ábrán, fekete vonallal jelölve.

Kezelés: a folyamat gyorsítására extra hardvert alkalmazunk, amivel a cache-ből történő lehíváskor egyúttal a végrehajtó egységbe is betöltjük az adatot (piros vonal). Így egy óraciklust megspórolhatunk.



2.2. ábra. A lehívott adat bevezetése a műveletvégző egységbe

2.3.4. WAR - Write After Read

Probléma felvetés: egy MUL utasítást egy ADD követ az alábbi módon:

```
; I1
MUL r3 , r2 , 1    ; r3 = r1 * r2
; I2
ADD r2 , r4 , r5    ; r2 = r4 + r5
```

A szorzás (MUL) sokkal lassabb, mint az összeadás (ADD), ezért előfordulhat, hogy a párhuzamos végrehajtás során I_2 hamarabb lefut, mint hogy I_1 betöltse a forrás operandust. Mivel I_2 módosította I_1 bemeneti operandusát, a MUL utasítás hibás eredményt fog adni. Következménye, hogy sérül a szekvenciális konzisztencia.

Megoldás: r_2 tartalmát egy ideiglenes regiszterbe irányítjuk (pl. r_{23}). Ekkor az assembly utasítások így néznek ki:

```
; I1
MUL r3 , r2 , 1    ; r3 = r1 * r2
; I2
ADD r23 , r4 , r5   ; r23 = r4 + r5
```

Az $r_{23} \rightarrow r_2$ hozzárendelést nyilvántartjuk, majd amikor a MUL utasítás végzett, visszaírjuk r_{23} tartalmát r_2 -be. Az átmeneti (átnevezési) regiszterek tulajdonságai:

- új, önálló, de rejtett,
- saját címtartománnyal rendelkezik,
- a programozó számára traszparens,
- extra hardvernek számít.

Megjegyzés: a regiszterkészletek csoportosítása:

- architekturális: programozó használja,
- átnevezési: a vezérlés használja az álfüggőségek feloldására.

2.3.5. WAW - Write After Write

Probléma felvetés: egy MUL utasítást egy ADD követ az alábbi módon:

```
; I1
MUL r3 , r2 , 1    ; r3 = r1 * r2
; I2
ADD r3 , r4 , r5    ; r3 = r4 + r5
```

A szorzás (MUL) sokkal lassabb, mint az összeadás (ADD), ezért előfordulhat, hogy a párhuzamos végrehajtás során I_1 később fut le, mint I_2 . Mivel az eredményt ugyanabba a regiszterbe írják, ebben az esetben I_1 felülírja I_2 eredményét az r_3 -ban. Ezzel sérül a szekvenciális konzisztencia.

Megoldás: r_3 átirányítása egy átnevezési regiszterbe, az előbb leírt módon.

2.3.6. Ciklusbeli függőség

Probléma felvetés: egy ciklusban az előző iterációban kiszámolt adatot használunk fel, például:

```
for  $i = 2$  to  $n$  do
   $X_i \leftarrow A_i * X_{i-1} + B_i$ 
end for
```

Kezelés: ez egy erős függőség, hardveresen nehezen feloldható. Megoldás az algoritmus áttervezése.

2.4. Vezérlés függőségek

Elágazások esetén léphetnek fel. Itt a statikus és dinamikus kezelésnek eltérő jelentése van, mint az adat-függőségeknél. A statikus kezelés itt egy állandó, mindig alkalmazható eljárást jelent, míg a dinamikus kezelés az adott programtól függ.

Probléma felvetés: az alábbi utasítássorozatban a feltétlen ugrás (JMP) egy SHL utasításra mutat:

```
DIV
MUL
JMP ; SHL-re mutat
ADD
...
SHL
```

Ekkor a kritikus utasítások így követik egymást időben:

	t_1	t_2	t_3	t_4	t_5	t_6
MUL	F_{MUL}	D	E	W/B		
JMP		F_{JMP}	D	E	W/B	
ADD			F_{ADD}	D	E	W/B

A JMP utasítás az Execute fázisban állítja át a Program Countert, ezzel végzi el az ugrást. A futószalag végrehajtás miatt azonban ekkorra már a következő utasítás, az ADD is lehívásra került, sőt, előfordulhat, hogy az azt követő utasítás is. Ezek viszont fölösleges lépések. Ritkább esetekben a JMP-t követő utasítás be is fejeződhet, mire az ugrás végrehajtásra kerül, ami veszélyezteti az architektúrais regisztertartalmakat.

Megoldás: a probléma kezelése statikus, dinamikus, vagy spekulatív (branch prediction) módon történhet.

Kezelés utasítások átrendezésével (dinamikus): compiler segítségével történő optimalizálás. A compiler megpróbálja átrendezni az utasítások sorrendjét. Az előző kódrészlet optimalizált változata:

```
JMP ; SHL-re mutat
MUL
DIV
ADD
...
SHL
```

Az optimalizálás utáni végrehajtási sorrend:

	t_1	t_2	t_3	t_4	t_5	t_6	t_7
JMP	F_{JMP}	D	E	W/B			
MUL		F_{MUL}	D	E	W/B		
DIV			F_{DIV}	D	E	W/B	
SHL				F_{SHL}	D	E	W/B

A sorrend megváltoztatásával elértük, hogy amíg az ugrás végre nem hajtódik (t_3), csak olyan utasításokat hívunk le, amiknek még az ugrás előtt kell lefutniuk. Mire az ADD utasításhoz elérnénk, felülíródik a PC és a megfelelő utasítás hívódik le (SHL). A módszer hátránya, hogy hatékonysága a futószalag fokozatok számának növelésével rohamosan csökken.

Kezelés NOP utasításokkal (statikus): a JMP utasítás mögé egy vagy több NOP utasítás kerül be. Ez a futószalag várákoztatását jelenti, amíg elő nem áll az ugráshoz szükséges PC. Ez az ún. ugrási buborék, nagysága $n - 1$, ahol n a futószalag fokozatok száma.