

3BA5

Performance effect of Branch Prediction

$$p_r[p_r(1+d) + (1-p_r)1] + (1-p_r)[p_r(1+c) + (1-p_r)(1+c)]$$

$$\begin{aligned} \text{average number of cycles per branch instruction} = \\ p_r[p_r(1+d) + (1-p_r)1] + (1-p_r)[p_r(1+c) + (1-p_r)(1+c)] = \\ p_r(p_r d + 1) + (1-p_r)(1+c) \end{aligned}$$

Substituted into

$$t_{\text{ave}} = p_b(\text{average number of cycles per branch instruction}) + (1-p_b)(\text{average number of cycles per non-branch inst.})$$

with p_b = probability instruction is a branch

$$\begin{aligned} t_{\text{ave}} &= p_b[p_r(p_r d + 1) + (1-p_r)(1+c)] + (1-p_b)(1) \\ &= 1 + c p_b(1-p_r) + d p_b p_r p_r \end{aligned}$$

3BA5

Performance effect of Branch Prediction

$$t_{\text{ave}} = 1 + c p_b(1-p_r) + d p_b p_r p_r$$

For example $d = 1$ and $p_r = 0.9$ For a 5-stage pipe with $c=3$

$$t_{\text{ave}} = 1 + 3 * 0.2 * 0.1 + 1 * 0.2 * 0.65 * 0.9 = 1.12$$

$$\text{Efficiency} = 1 / t_{\text{ave}} * 100 = 89\%$$

3BA5

Multiple or Target Prefetching

- ⊕ On first encounter with a branch we initiate a prefetch of its target path code at ID in order to lower delay if the branch is taken.

3BA5

Delayed Branching

- ⊕ A compiler/assembler technique is to delay branching for k cycles after the branch instruction, providing time to fetch the target path.
- ⊕ Useful instructions or NOPs instructions are inserted

Example

Delay the next k sequential instructions

```

i1:  LOAD R1, A
i2:  LOAD R2, B
i3:  BRZR R2, i7  -- branch to i7 if R2=0
i4:  LOAD R3, C
i5:  ADD  R4, R2, R3
i6:  MUL  R5, R1, R2
i7:  ADD  R4, R1, R2

```

```

i2:  LOAD R2, B
i3:  BRZR R2, i7
i4:  LOAD R1, A
NOP
i4:  LOAD R3, C
i5:  ADD  R4, R2, R3
i6:  MUL  R5, R1, R2
i7:  ADD  R4, R1, R2

```

Assuming that $k=2$, the compiler modifies this code by moving the instruction i_4 and NOP instruction after the branch instruction i_3 .

Super-scalar vs. Super-pipelines

Super-scalar and Very Long Instruction Word (VLIW) design exploits the high circuit density of CMOS to provide multiple functional units and hence more parallelism.

Super-scalar vs. Super-pipelines

Example

For example with two fp-add pipes large vector additions can be distributed.

```

DO I=1 to 10000
  c(i) = a(i) + b(i)
END

DO I=1 to 10000 by 2
  c(i) = a(i) + b(i)   c(i+1) = a(i+1) + b(i+1)
  ADDER_1             ADDER_2
END

```

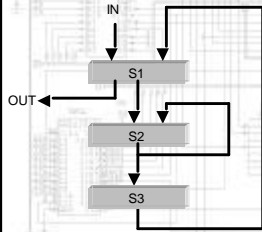
Super-pipelines

- Superpipelines rely on:
 - High speed technology
 - Small pipe stages
 - Hence long pipes to increase the performance.
- Such a design is susceptible to control hazards
 - Needs good branch prediction to perform well

Pipeline Control and Scheduling

Consider a re-circulating pipe

We consider only static pipes



	t_0	t_1	t_2	t_3	t_4
Stage 1	✓				✓
Stage 2		✓	✓		
Stage 3				X	

Reservation Table

Forbidden List
 0, 4
 0, 1
 0

Definitions

- ⊕ Latency:
 - ⊕ The interval between two initiations in a stage. Latencies between reservations are forbidden latencies.
- ⊕ Forbidden List:
 - ⊕ The vector forbidden latencies
 - ⊕ $F = (4, 1, 0)$
- ⊕ Collision Vector:
 - ⊕ $C = (c_n, \dots, c_1, c_0)$
 - ⊕ where $c_i = \{ 1 \in F, 0 \notin F \}$
 - ⊕ Eg. $C = (c_4, c_3, c_2, c_1, c_0)$
 $C = (1, 0, 0, 1, 1)$

Collision Vector

- ⊕ When $C = (1, 0, 0, 1, 1)$ it is seen that $c_i = 0$ means we may safely insert a new input to the pipe.
 - ⊕ Hence at 2,3 or 5,6
 - ⊕ If we do insert a new input at e.g. $t = 2$, then the forbidden latencies may alter, e.g.

	t_0	t_1	t_2	t_3	t_4	t_5	t_6
Stage 1	✓						
Stage 2		✓	✓				
Stage 3				X		Y	

Reservation Table

Forbidden List
 0, 2, 4
 0, 1, 2
 0, 2

Collision Vector

with altered forbidden latencies may

	t_0	t_1	t_2	t_3	t_4	t_5	t_6
Stage 1	✓						
Stage 2		✓	✓				
Stage 3				X		Y	

Reservation Table

Forbidden List
 0, 2, 4
 0, 1, 2
 0, 2

- ⊕ When it is seen that $c_i = 0$ means we may safely insert a new input to the pipe.
 - ⊕ $C_2 = (1, 0, 0, 1, 1)$ OR $C = (0, 0, 1, 0, 0)$
 $= (1, 0, 1, 1, 1)$

State Diagram

⊕ Hence we may describe the forbidden latencies of the pipe with a statediagram, starting at:

⊕ $C = (1, 0, 0, 1, 1)$

