# Pipelining, Branch Prediction, Trends

10.1-10.4

# Topics

- 10.1 Quantitative Analyses of Program Execution
- 10.2 From CISC to RISC
- 10.3 Pipelining the Datapath Branch Prediction, Delay Slots
- 10.4 Overlapping Register Windows





# Frequency of Instructions

 Frequency of occurrence of instruction types for a variety of languages/benchmark programs

Statement	Average Percent of Time
Assignment	47
If	23
Call	15
Loop	6
Goto	3
Other	7

Arithmetic and other "powerful" instructions only 7%

	Percentage of Number of Terms in Assignments	Percentage of Number of Locals in Procedures	Percentage of Number of Parameters in Procedure Calls
0	_	22	41
1	80	17	19
2	15	20	15
3	3	14	9
4	2	8	7
≥ 5	0	20	8

80% of assignments involve one term;

80% of procedures could be handled supported 4 locals









### **Instruction Prefetch**

- Simple version of Pipelining treating the instruction cycle like an assembly line
- Fetch accessing main memory
- Execution usually does not access main memory
- Can fetch next instruction during execution of current instruction
- Called instruction prefetch















### Dealing with Branches

- Multiple Streams
- Prefetch Branch Target
- Loop buffer
- Branch prediction
- Delayed branching

#### **Multiple Streams**

- Have two pipelines
- Prefetch each branch into a separate pipeline
- Use appropriate pipeline
- Leads to bus & register contention
- Still a penalty since it takes some cycles to figure out the branch target and start fetching instructions from there
- Multiple branches lead to further pipelines being needed
  - Would need more than two pipelines then
- More expensive circuitry

# Prefetch Branch Target

- Target of branch is prefetched in addition to instructions following branch
  - Prefetch here means getting these instructions and storing them in the cache
- Keep target until branch is executed
- Used by IBM 360/91

# Loop Buffer

- Very fast memory
- Maintained by fetch stage of pipeline
- Remembers the last N instructions
- Check buffer before fetching from memory
- Very good for small loops or jumps
- c.f. cache
- Used by CRAY-1

#### Branch Prediction (1)

- Predict never taken
  - Assume that jump will not happen
  - Always fetch next instruction
  - 68020 & VAX 11/780
  - VAX will not prefetch after branch if a page fault would result (O/S v CPU design)
- Predict always taken
  - Assume that jump will happen
  - Always fetch target instruction
  - Studies indicate branches are taken around 60% of the time in most programs







#### Dealing with Branches – RISC Approach

#### • Delayed Branch - used with RISC machines

- Requires some clever rearrangement of instructions
- Burden on programmers but can increase performance
- Most RISC machines: Doesn't flush the pipeline in case of a branch
- Called the Delayed Branch
  - This means if we take a branch, we'll still continue to execute whatever is currently in the pipeline, at a minimum the next instruction
  - Benefit: Simplifies the hardware quite a bit
  - But we need to make sure it is safe to execute the remaining instructions in the pipeline
  - Simple solution to get same behavior as a flushed pipeline: Insert NOP No Operation instructions after a branch
    - Called the Delay Slot

# **RISC** Pipeline with Delay Slot

Using a Five Stage pipeline:

IF = Fetch, ID = Decode, EX = Execute

MEM = Memory access, WB = Write back register values

In this example: CPU knows if branches are to be taken after the ID stage (implications if not known until after the EX stage?)

Untaken branch instruction	IF	ID	EX	MEM	WB				
Branch delay instruction (i + 1)		IF	ID	EX	MEM	WB			
Instruction $i + 2$			IF	ID	EX	MEM	WB		
Instruction i + 3				IF	ID	EX	MEM	WB	
Instruction i + 4					IF	ID	EX	MEM	WB
Taken branch instruction	IF	ID	EX	MEM	WB				
Branch delay instruction $(i + 1)$		IF	ID	EX	MEM	WB			
Branch target			IF	ID	EX	MEM	WB		
brunen talget									
Branch target + 1				IF	ID	EX	MEM	WB	

# Normal vs. Delayed Branch

Address	Normal	Delayed		
100	LOAD X,A	LOAD X,A		
101	ADD 1,A	ADD 1,A		
102	JUMP 105	JUMP 106		
103	ADD A,B	NOOP		
104	SUB C,B	ADD A,B		
105	STORE A,Z	SUB C,B		
106		STORE A,Z		

One delay slot - Next instruction is always in the pipeline. "Normal" path contains an implicit "NOP" instruction as the pipeline gets flushed. Delayed branch requires explicit NOP instruction placed in the code!

Optimized Delayed Branch But we can optimize this code by rearrangement! Notice we always Add 1 to A so we can use this instruction to fill the delay slot								
Address	Normal	Delayed	Optimized					
100	LOAD X,A	LOAD X,A	LOAD X,A					
101	ADD 1,A	ADD 1,A	JUMP 105					
102	JUMP 105	JUMP 106	ADD 1,A					
103	ADD A,B	NOOP	ADD A,B					
104	SUB C,B	ADD A,B	SUB C,B					
105	STORE A,Z	SUB C,B	STORE A,Z					
106		STORE A,Z						





### Other Pipelining Overhead

- Each stage of the pipeline has overhead in moving data from buffer to buffer for one stage to another. This can lengthen the total time it takes to execute a single instruction!
- The amount of control logic required to handle memory and register dependencies and to optimize the use of the pipeline increases enormously with the number of stages. This can lead to a case where the logic between stages is more complex than the actual stages being controlled.
- Need balance, careful design to optimize pipelining



486 Pipelining Examples										
Fetch	D1	D2	Ex	WB			MOV R1, M			
	Fetch	D1	D2	Ex	WB		MOV R1, R2			
L		Fetch	D1	D2	Ex	WB	MOV M, R1			
Fetch	D1	D2	Ex	WB			MOV R2, M			
	Fetch	D1		D2	Ex		MOV R1, (R2)			
Need R2 written back to use as addr for second instruction in stage D2										
Normally this data is not available until after the WB stage, but bypass circuitry allows us to send the proper data directly to EX of the next stage (this is called <b>forwarding</b> )										



# Pentium II/IV Pipelining

- Pentium II
  - 12 pipeline stages
  - Dynamic execution incorporates the concepts of out of order and speculative execution
  - Two-level, adaptive-training, branch prediction mechanism
- Pentium IV
  - 20 stage pipeline
  - Combines different branch prediction
    - mechanisms to keep the pipeline full







# **Register Windows**

- Parameters are "passed" by simply updating the window pointer
  - All parameter access in registers, very fast
  - In the rare event we exceed the number of registers available, can use main memory for overflow