Chapter 5: Superscalar Techniques Part 1: Instruction Flow

Modern Processor Design: Fundamentals of Superscalar Processors

# Instruction Flow Techniques

- Goal and Impediments
- Branch Types and Implementations
- What's So Bad About Branches?
- What are Control Dependences?
- Impact of Control Dependences on Performance
- Improving I-Cache Performance

# Instruction Flow in Context

# Goal and Impediments

- Goal of Instruction Flow
  - Supply processor with maximum number of <u>useful</u> instructions every clock cycle
- Impediments
  - Branches and jumps
  - I-Cache limitations
    - hit rate
    - width
    - alignment
    - etc.

# Branch Types and Implementation

#### 1. Types of Branches

- A. Conditional or Unconditional
- B. Save PC? Save other processor state?
- C. How is target computed?
  - constant target (immediate, PC-relative)
  - variable target (register, register + offset)
- 2. Branch Architectures
  - A. Condition code or condition registers
  - B. Register

# What's So Bad About Branches?

- Effects of Branches
  - Fragmentation of I-Cache lines
  - Need to determine branch outcome
  - Need to determine branch target
  - Use up execution resources

# What's So Bad About Branches?

Problem: Fetch stalls until outcome is determined

#### Solutions:

- Minimize delay
  - Move instructions determining branch condition away from branch
  - Make use of delay
  - Non-speculative:
    - Fill delay slots with useful safe instructions
    - Execute both paths (eager execution)
  - Speculative:
  - Predict branch outcome

# What's So Bad About Branches?

Problem: Fetch stalls until branch target is determined Solutions:

- <u>Minimize delay</u>
  - Generate branch target early
- Make use of delay: Predict branch target
  - Single target (constant—only need to compute once)
  - Multiple targets (variable—but may use previous target value as prediction for next time).





# Limits on Instruction Level Parallelism (ILP)

Weiss and Smith [1984]	1.58
Sohi and Vajapeyam [1987]	1.81
Tjaden and Flynn [1970]	1.86 (Flynn's bottleneck)
Tjaden and Flynn [1973]	1.96
Uht [1986]	2.00
Smith et al. [1989]	2.00
Jouppi and Wall [1988]	2.40
Johnson [1991]	2.50
Acosta et al. [1986]	2.79
Wedig [1982]	3.00
Butler et al. [1991]	5.8
Melvin and Patt [1991]	6
Wall [1991]	7 (Jouppi disagreed)
Kuck et al. [1972]	8
Riseman and Foster [1972]	51 (no control dependences)
Nicolau and Fisher [1984]	90 (Fisher's optimism)

# Riseman and Foster's Study • 7 benchmark programs on CDC-3600 Assume infinite machines

- Infinite memory and instruction stack
- Infinite register file
- Infinite functional units
- True dependencies only at dataflow limit
- · If lookahead bounded to single basic block,
- speedup is 1.72 (Flynn's bottleneck)
- If one can bypass branches (hypothetically), then: a theoretical ILP of 51 was determined by **Riseman and Foster**



- · Implicit Sequential Control Flow
  - Static Program Representation
    - Control Flow Graph (CFG)
    - Nodes = basic blocks
    - Edges = Control flow transfers
  - Physical Program Layout
    - Mapping of CFG to linear program memory
    - Implied sequential control flow
  - Dynamic Program Execution
    - Traversal of the CFG nodes and edges (e.g. loops) · Traversal dictated by branch conditions

  - Dynamic Control Flow
    - Deviates from sequential control flow • Disrupts sequential fetching
    - · Can stall IF stage and reduce I-fetch bandwidth

#### **Program Control Flow**

- Dynamic traversal of static CFG
- · Mapping CFG to linear memory



















#### **Branch Prediction Function**

Prediction function F(X1, X2, ...)
 X1 – opcode type

X2 – history

Prediction effectiveness based on opcode only, or history

	IBM1	IBM2	IBM3	IBM4	DEC	CDC
Opcode only	66	69	71	55	80	78
History 1	92	95	87	80	97	82
History 2	93	97	91	83	98	91
History 3	94	97	91	84	98	94
History 4	95	97	92	84	98	95
History 5	95	97	92	84	98	96







	% of each branch type % bc with p cycles				vith per	enalty	
Benchmark	b	bl	bc	bcr	3 сус	2 cyc	1 cyc
spice2g6	7.86	0.30	12.58	0.32	13.82	3.12	0.76
doduc	1.00	0.94	8.22	1.01	10.14	1.76	2.02
matrix300	0.00	0.00	14.50	0.00	0.68	0.22	0.20
tomcatv	0.00	0.00	6.10	0.00	0.24	0.02	0.01
gcc	2.30	1.32	15.50	1.81	22.46	9.48	4.85
espresso	3.61	0.58	19.85	0.68	37.37	1.77	0.31
li	2.41	1.92	14.36	1.91	31.55	3.44	1.37
eqntott	0.91	0.47	32.87	0.51	5.01	11.01	0.80

#### Exhaustive Search for Optimal 2-bit Predictor

- There are 2<sup>20</sup> possible state machines of 2-bit predictors
- Some machines are uninteresting, pruning them out reduces the number of state machines to 5248
- For each benchmark, determine prediction accuracy for all the predictor state machines



•



	Prediction Accuracy (Overall CPI Overhead)			
Benchmark	3 bit	2 bit	1 bit	0 bit
spice2g6	97.0 (0.009)	97.0 (0.009)	96.2 (0.013)	76.6 (0.031)
doduc	94.2 (0.003)	94.3 (0.003)	90.2 (0.004)	69.2 (0.022)
gcc	89.7 (0.025)	89.1 (0.026)	86.0 (0.033)	50.0 (0.128)
espresso	89.5 (0.045)	89.1 (0.047)	87.2 (0.054)	58.5 (0.176)
li	88.3 (0.042)	86.8 (0.048)	82.5 (0.063)	62.4 (0.142)
eqntott	89.3 (0.028)	87.2 (0.033)	82.9 (0.046)	78.4 (0.049)

In collisions, multiple branches share the same predictor

29

•

Constructive interference
 Destructive interference
 Marginal gains beyond 1K entries (for these programs)































#### IBM 360/91 FPU

- Multiple functional units (FU's)
  - Floating-point add
  - Floating-point multiply/divide
- Three register files (pseudo reg-reg machine in floating-point unit) - (4) floating-point registers (FLR)
  - (6) floating-point buffers (FLB)
  - (3) store data buffers (SDB)
- Out of order instruction execution:

  - After decode the instruction unit passes all floating point instructions (in order) to the floating-point operation stack (FLOS).
     In the floating point unit, instructions are then further decoded and issued from the FLOS to the two FU's
- Variable operation latencies:
  - Floating-point add: 2 cycles

  - Floating-point multiply: 3 cycles Floating-point divide: 12 cycles
- Goal: achieve concurrent execution of multiple floating-point instructions, in addition to achieving one instruction per cycle in instruction pipeline 45

#### **Dependence Mechanisms**

Two Address IBM 360 Instruction Format:

#### R1 <-- R1 op R2

- Major dependence mechanisms:
- Structural (FU) dependence = > virtual FU's
- Reservation stations
- True dependence = > pseudo operands + result forwarding
  - Register tags
  - Reservation stations
  - Common data bus (CDB)
- Anti-dependence = > operand copying Reservation stations
- Output dependence = > register renaming + result forwarding
  - Register tags
  - Reservation stations
  - Common data bus (CDB)

46



## **Reservation Stations**

- Used to collect operands or pseudo operands (tags).
- Associate more than one set of buffering registers (control, source, sink) with each FU, = > virtual FU's.
- Add unit: three reservation stations
- Multiply/divide unit: two reservation stations



 CDB is fed by all units that can alter a register (or supply register values) and it feeds all units which can have a register as an operand.

#### Sources of CDB:

- Floating-point buffers (FLB)
- Two FU's (add unit and the multiply/divide unit)
- Destinations of CDB:
  - Reservation stations
  - Floating-point registers (FLR)
  - Store data buffers (SDB)



**Operation of Dependence Mechanisms** 

#### 1. Structural (FU) dependence = > virtual FU's

- FLOS can hold and decode up to 8 instructions.
- Instructions are dispatched to the 5 reservation stations (virtual FU's) even though there are only two physical FU's.
- Hence, structural dependence does not stall dispatching.

# True dependence = > pseudo operands + result forwarding If an operand is available in FLR, it is copied to a res. station entry.

- If an operand is not available (i.e. there is pending write), then a tag is copied to the reservation station entry instead. This tag identifies the source of the pending write. This instruction then waits in its reservation station for the true dependence to be resolved.
- When the operand is finally produced by the source (ID of source = tag value), this source unit asserts its ID, i.e. its tag value, on the CDB followed by broadcasting of the operand on the CDB.
- All the reservation station entries and the FLR entries and SDB entries carrying this tag value in their tag fields will detect a match of tag values and latch in the broadcasted operand from the CDB.
- Hence, true dependence does not block subsequent independent instructions and does not stall a physical FU. Forwarding also minimizes delay due to true dependence.









- 3. Output dependence = > register renaming + result forwarding
  - If a register is waiting for a pending write, its tag field will contain the ID, or tag value, of the source for that pending write.
  - When that source eventually produces the result, that result will be written into the register via the CDB.
  - It is possible that prior to the completion of the pending write, another instruction can come along and also has that same register as its destination register.
  - If this occurs, the operands (or pseudo operands) needed by this instruction are still copied to an available reservation station. In addition, the tag field of the destination register of this instruction is updated with the ID of this new reservation station, i.e. the old tag value is overwritten. This will ensure that the said register will get the latest value, i.e. the late completing earlier write cannot overwrite a later write.
  - Hence, the output dependence is resolved without stalling a physical functional unit, not requiring additional buffers to ensure sequential write back to the register file.



#### Summary of Tomasulo's Algorithm

- Supports <u>out of order</u> execution of instructions.
- Resolves dependences dynamically using hardware.
- Attempts to delay the resolution of dependencies as late as possible.
- Structural dependence does not stall issuing; virtual FU's in the form
   of reservation stations are used.
- Output dependence does not stall issuing; copying of old tag to reservation station and updating of tag field of the register with pending write with the new tag.
- True dependence with a pending write operand does not stall the reading of operands; pseudo operand (tag) is copied to reservation station.
- Anti-dependence does not stall write back; earlier copying of operand awaiting read to the reservation station.
- Can support sequence of multiple output dependences.
- Forwarding from FU's to reservation stations bypasses the register file.

57

#### Tomasulo vs. Modern OOO

	IBM 360/91	Modern
Width	Peak IPC = 1	4+
Structural hazards	2 FPU	Many FU
	Single CDB	Many busses
Anti-dependences	Operand copy	Reg. Renaming
Output dependences	Renamed reg. tag	Reg. renaming
True dependences	Tag-based forw.	Tag-based forw.
Exceptions	Imprecise	Precise (ROB)
Implementation	3 x 66" x 15" x 78"	1 chip
	60ns cycle time	300ps
	11-12 gate delays	< \$100
	per pipe stage	
	>\$1 million	58

















![](_page_16_Figure_2.jpeg)

![](_page_16_Figure_3.jpeg)

	Regis	ter File	e Alterna	tives	
Pagiatar		Duration	Result stored	l where?	
Lifetime	Status	(cycles)	Future File	History File	Phys. RF
Dispatch	Unavail	≥1	N/A	N/A	N/A
Finish execution	Speculative	≥ 0	FF	ARF	PRF
Commit	Committed	≥ 0	ARF	ARF	PRF
Next def. Dispatched	Committed	≥ 1	ARF	HF	PRF
Next def. Committed	Discarded	≥0 ization	Overwritten	Discarded	Reclaimed

Future file (future updates buffered, later committed)

Rename register file

 History file (old versions buffered, later discarded) - Merged (single physical register file)

# Register Commit

**Register File Commit** 

- History file (only proposed)
  - · Copy previous value from ARF to HF at dispatch
  - · Use HF to reconstruct precise state if needed
- Future file: separate ARF & RRF (lecture notes, PPC 604/620, Pentium Pro)
  - · Copy committed value from RRF to ARF
  - Update rename table mapping
- Physical Register File: merged ARF & RRF (MIPS R10000 paper, Pentium 4, Alpha 21264, Power 4)
  - · No copy; simpler datapath (operand always in PRF)
  - · Simply "commit" rename table mapping as branches resolve 70

#### **Rename Table Implementation**

- MAP checkpointing
  - Recovery from branches, exceptions
  - Checkpoint granularity
    - · Every instruction
    - · Every branch, playback to get to exception
- boundary
- RAM Map
  - Just a lookup table; checkpoints nxm each
- CAM Map
  - Positional bit vectors; checkpoints a single column

71

69

# Summary

- Register dependences
  - True dependences
  - Antidependences
  - Output dependences
- **Register Renaming**
- Tomasulo's Algorithm
- **Reservation Station Implementation** .
- **Reorder Buffer Implementation** Register File Implementation
- - History file - Future file
  - Physical register file
- Rename Table Implementation

![](_page_18_Figure_0.jpeg)

75

#### Memory Data Dependences Besides branches, long memory latencies are one of the biggest performance challenges today. To preserve sequential (in-order) state in the data caches and external memory (so that recovery from exceptions is possible) stores are performed in order. This takes care of antidependences and output dependences to memory locations. However, loads can be issued out of order with respect to stores if the out-of-order loads check for data dependences with respect to previous, pending stores. RAW WAW WAR store X load X store X ÷ ÷ ε. store X store X load X 74

	Mer	mory Data Dependences
• " <u>Mer</u> mem	<mark>nory Alia</mark> ory locat	asing" = Two memory references involving the same ion (collision of two memory addresses).
• " <u>Mer</u> refer	nory Dis ences wil	ambiguation" = Determining whether two memory l alias or not (whether there is a dependence or not).
• Mem	ory Depe	endency Detection:
– N	Aust comp	ute effective addresses of both memory references
– E	effective ad	ddresses can depend on run-time data and other instructions
- 0	Comparisor	of addresses require much wider comparators
Example co	de:	
(1)	STORE	V
(2)	ADD	
(3)	LOAD	W
(4)	LOAD	х
(5)	LOAD	V
(6)	ADD	
(7)	STORE	W

#### Total Order of Loads and Stores

- Keep all loads and stores totally in order with respect to each other.
- However, loads and stores can execute out of order with respect to
   other types of instructions.
- Consequently, stores are held for all previous instructions, and loads are held for stores.
  - I.e. stores performed at commit point
  - Sufficient to prevent wrong branch path stores since all prior branches now resolved

![](_page_19_Figure_0.jpeg)

![](_page_19_Figure_1.jpeg)

![](_page_19_Figure_2.jpeg)

# Load Forwarding

- If a subsequent load has a dependence on a store still in the store buffer, it need not wait till the store is issued to the data cache.
- The load can be directly satisfied from the store buffer if the address is valid and the data is available in the store buffer.
- This avoids the latency of accessing the data cache.

![](_page_20_Figure_0.jpeg)

![](_page_20_Figure_1.jpeg)

![](_page_20_Figure_2.jpeg)

Optimizing Lo	oad/Store Disambiguatior	۱
<ul> <li>Non-specula</li> <li>1. Loads wait</li> <li>2. Full address</li> <li>3. Bypass if no</li> <li>(1) can limit</li> </ul>	ative load/store disambiguation for addresses of all prior stores s comparison o match, forward if match performance:	
load r5,MEM[r3] store r7, MEM[r5]  load r8, MEM[r9]	<ul> <li>← cache miss</li> <li>← RAW for agen, stalled</li> <li>← independent load stalled</li> </ul>	
		84

![](_page_21_Figure_0.jpeg)

![](_page_21_Figure_1.jpeg)

## Load/Store Disambiguation Discussion

- RISC ISA:
  - Many registers, most variables allocated to registers
  - Aliases are rare
  - Most important to not delay loads (bypass)
  - Alias predictor may/may not be necessary
- CISC ISA:
  - Few registers, many operands from memory
  - Aliases much more common, forwarding necessary
  - Incorrect load speculation should be avoided
  - If load speculation allowed, predictor probably necessary
     Address translation;
- Address translation:
  - Can't use virtual address (must use physical)
  - Wait till after TLB lookup is done
  - Or, use subset of untranslated bits (page offset)
  - Safe for proving inequality (bypassing OK)
  - Not sufficient for showing equality (forwarding not OK)

![](_page_21_Figure_19.jpeg)

![](_page_22_Figure_0.jpeg)

#### Memory Bottleneck Techniques

#### Dynamic Hardware (Microarchitecture):

Use Non-blocking D-cache (need missed-load buffers) Use Multiple Load/Store Units (need multiported D-cache) Use More Advanced Caches (victim cache, stream buffer) Use Hardware Prefetching (need load history and stride detection)

#### Static Software (Code Transformation):

Insert Prefetch or Cache-Touch Instructions (mask miss penalty) Array Blocking Based on Cache Organization (minimize misses) Reduce Unnecessary Load/Store Instructions (redundant loads) Software Controlled Memory Hierarchy (expose it to above DSI)

91

#### Advanced Memory Hierarchy

Easing The Memory Bottleneck

Dispatch

Integer

Complete

Retire

T

Dispatch Buffer

Branch Integer

Reorder Buff.

Store Buff.

RS

Reg. Write Back

Reg. File

Load/

Store

Data Cache

ł

Float.

Point

Ren. Reg.

Load/

Store

90

- Better miss rate: victim caches
- Reducing miss costs through software restructuring
- · Higher bandwidth: Lock-up free caches, superscalar caches
- Beyond simple blocks
- · Two level caches
- Prefetching, software prefetching
- Main Memory, DRAM
- Virtual Memory, TLBs
- Interaction of caches, virtual memory

![](_page_23_Figure_0.jpeg)

![](_page_23_Figure_1.jpeg)

![](_page_23_Figure_2.jpeg)

![](_page_24_Figure_0.jpeg)

- Increasing issue width => wider caches
- Parallel cache accesses are harder than parallel functional units
- Fundamental difference:
  - Caches have state, functional units don't
  - Operation thru one port affects future operations thru others
- Several approaches used
  - True multi-porting
  - Multiple cache copies
  - Virtual multi-porting
  - Multi-banking (interleaving)

99

![](_page_24_Figure_12.jpeg)

# True Multiporting of SRAM

- · Would be ideal
- Increases cache area
  - Array becomes wire-dominated
- Slower access
  - Wire delay across larger area
  - Cross-coupling capacitance between wires
- SRAM access difficult to pipeline

![](_page_24_Figure_21.jpeg)

![](_page_25_Figure_0.jpeg)

![](_page_25_Figure_2.jpeg)

#### **Combined Schemes**

- Multiple banks with multiple ports
- Virtual multiporting of multiple banks
- Multiple ports and virtual multiporting
- Multiple banks with multiply virtually multiported ports
- Complexity!
- No good solution known at this time
   Current generation superscalars get by with 1-3 ports
- Course project?

103

#### **Beyond Simple Blocks**

- · Break blocks into
  - Address block associated with tag
  - Transfer block to/from memory (subline, sub-block)
- · Large address blocks
  - Decrease tag overhead
  - But allow fewer blocks to reside in cache (fixed mapping)

#### Subline Valid Bits

Tag			Subline 0	Subline 1	Subline 2	Subline 3
Tag			Subline 0	Subline 1	Subline 2	Subline 3
Tag			Subline 0	Subline 1	Subline 2	Subline 3
Tag			Subline 0	Subline 1	Subline 2	Subline 3

## **Beyond Simple Blocks**

#### · Larger transfer block

- Exploit spatial locality
- Amortize memory latency
- But take longer to load
- Replace more data already cached (more conflicts)
- Cause unnecessary traffic
- Typically used in large L2/L3 caches to limit tag overhead
- Sublines tracked by MSHR during pending fill

#### Subline Valid Bits

Tag Subline 0 Sublin	ne 1 Subline 2 Subline
Tag Subline 0 Sublin	ne 1 Subline 2 Subline
Tag Subline 0 Sublin	ne 1 Subline 2 Subline

#### Latency vs. Bandwidth

- · Latency can be handled by
  - Hiding (or tolerating) it out of order issue, nonblocking cache
  - Reducing it better caches
- Parallelism helps to hide latency
  - MLP multiple outstanding cache misses overlapped
- · But increases bandwidth demand
- · Latency ultimately limited by physics

#### Latency vs. Bandwidth

- · Bandwidth can be handled by "spending" more (hardware cost)
  - Wider buses, interfaces

•

- Banking/interleaving, multiporting
- Ignoring cost, a well-designed system should never be bandwidth-limited
   Can't ignore cost!
- Bandwidth improvement usually increases latency
   No free lunch
- Hierarchies decrease bandwidth demand to lower levels
   Serve as traffic filters: a hit in L1 is filtered from L2
  - Parallelism puts more demand on bandwidth
- If average b/w demand is not met => infinite queues
   Bursts are smoothed by queues
- If burst is much larger than average => long queue
  - Eventually increases delay to unacceptable levels

107

#### Prefetching

- Even "demand fetching" prefetches other words in block
   Spatial locality
- Prefetching is useless
  - Unless a prefetch costs less than demand miss
- Ideally, prefetches should
  - Always get data before it is referenced
  - Never get data not used
  - Never prematurely replace data
  - Never interfere with other cache activity

![](_page_27_Figure_0.jpeg)

- Use compiler to try to
  - Prefetch early
  - Prefetch accurately
- Prefetch into
  - Register (binding)
    - · Use normal loads? ROB fills up, fetch stalls
  - What about page faults? Exceptions?
  - Caches (non-binding) preferred
    - Needs ISA support

#### Software Prefetching

- For example:
- do j= 1, cols
- do ii = 1 to rows by BLOCK
- prefetch (&(x[i,j])+BLOCK) # prefetch one block ahead
- do i = ii to ii + BLOCK-1
- sum = sum + x[i,j]
- How many blocks ahead should we prefetch?
  - Affects timeliness of prefetches
  - Must be scaled based on miss latency

![](_page_27_Figure_22.jpeg)

![](_page_27_Figure_23.jpeg)

## Stream or Prefetch Buffers

- Prefetching causes capacity and conflict misses (pollution)
   Can displace useful blocks
- · Aimed at compulsory and capacity misses
- Prefetch into buffers, NOT into cache
  - On miss start filling stream buffer with successive lines
  - Check both cache and stream buffer
    - Hit in stream buffer => move line into cache (promote)
    - Miss in both  $\Rightarrow$  clear and refill stream buffer
- Performance
  - Very effective for I-caches, less for D-caches
  - Multiple buffers to capture multiple streams (better for D-caches)
- · Can use with any prefetching scheme to avoid pollution

113

#### **Multilevel Caches**

- Ubiquitous in high-performance processors
  - Gap between L1 (core frequency) and main memory too high
  - Level 2 usually on chip, level 3 on or off-chip, level 4 off chip
- · Inclusion in multilevel caches
  - Multi-level inclusion holds if L2 cache is superset of L1
  - Can handle virtual address synonyms
  - Filter coherence traffic: if L2 misses, L1 needn't see snoop
  - Makes L1 writes simpler
    - · For both write-through and write-back

114

![](_page_28_Figure_24.jpeg)

# Multilevel Miss Rates

#### • Miss rates of lower level caches

- Affected by upper level filtering effect
- LRU becomes LRM, since "use" is "miss"
- Can affect miss rates, though usually not important

#### • Miss rates reported as:

- Miss per instruction
- Global miss rate
- Local miss rate
- "Solo" miss rate

• L2 cache sees all references (unfiltered by L1)