

as a general-purpose register, and using R14 is very difficult because it interferes with instructions that manipulate the stack frame. Condition codes are used for branching and are set by all arithmetic and logical operations and by the move instruction. The move instruction transfers data between any two addressable locations and subsumes load, store, register–register moves, and memory–memory moves as special cases.

VAX Addressing Modes

The addressing modes include most of those we discussed in Chapter 3: literal, register (operand is in a register), register deferred (register indirect), autodecrement, autoincrement, autoincrement deferred, byte/word/long displacement, byte/word/long displacement deferred, and scaled (called “indexed” in the VAX architecture). Scaled addressing mode may be applied to any general addressing mode except register or literal. Register is an addressing mode no different from any other in the VAX. Thus, a 3-operand VAX instruction may include from zero to three operand memory references, each of which may be any of the memory addressing modes. Since the memory indirect modes require an additional memory access, up to 6 memory accesses may be required for a 3-operand instruction. When the addressing modes are used with R15 (the PC), only a few are defined, and their meaning is special. The defined addressing modes with R15 are as follows:

- *Immediate*—an immediate value is in the instruction stream; this mode is encoded as autoincrement on PC.
- *Absolute*—a 32-bit absolute address is in the instruction stream; this mode is encoded as autoincrement deferred with PC as the register.
- *Byte/word/long displacement*—the same as the general mode, but the base is the PC, giving PC-relative addressing.
- *Byte/word/long displacement deferred*—the same as the general mode, but the base is the PC, giving addressing that is indirect through a memory location that is PC-relative.

A VAX instruction consists of an opcode followed by zero or more operand specifiers. The opcode is almost always a single byte that specifies the operation, the data type, and the operand count. Almost all operations are fully orthogonal with respect to addressing modes—any combination of addressing modes works with nearly every opcode, and many operations are supported for all possible data types.

Operand specifiers may vary in length from one byte to many, depending on the information to be conveyed. The first byte of each operand specifier consists of two 4-bit fields: the type of address specifier and a register that is part of the addressing mode. If the operand specifier requires additional bytes to specify a

displacement, additional registers, or an immediate value, it is extended in 1-byte increments. The name, assembler syntax, and number of bytes for each operand specifier are shown in Figure 4.3. The total instruction length and format are easy to state: Simply add up the sizes of the operand specifiers and include one byte (or rarely two) for the opcode.

Example

How long is the following instruction?

```
ADDL3 R1, 737(R2), #456
```

Answer

The opcode length is 1 byte, as is the first operand specifier (R1). The second operand specifier has two parts: the first part is a byte that specifies the addressing mode and base register; the second part is the 2-byte long displacement. The third operand specifier also has two parts: the first byte specifies immediate mode, and the second part contains the immediate. Because the data type is long (ADDL3), the immediate value takes 4 bytes.

Thus, the total length of the instruction is $1 + 1 + (1+2) + (1+4) = 10$ bytes.

Addressing mode	Syntax	Length in bytes
Literal	#value	1 (6-bit signed value)
Immediate	#value	1 + length of the immediate
Register	Rn	1
Register deferred	(Rn)	1
Byte/word/long displacement	Displacement (Rn)	1 + length of the displacement
Byte/word/long displacement deferred	@displacement (Rn)	1 + length of the displacement
Scaled (Indexed)	Base mode [Rx]	1 + length of base addressing mode
Autoincrement	(Rn)+	1
Autodecrement	-(Rn)	1
Autoincrement deferred	@(Rn)+	1

FIGURE 4.3 Length of the VAX operand specifiers. The length of each addressing mode is 1 byte plus the length of any displacement or immediate field that is in the mode. Literal mode uses a special 2-bit tag and the remaining 6 bits encode the constant value. The data we examined in Chapter 3 on constants showed the heavy use of small constants; the same observation motivated this optimization. The length of an immediate is dictated by the data type indicated in the opcode, not the value of the immediate.

Type	Example	Instruction meaning
Data transfers		Move data between byte, halfword, word, or doubleword operands; * is the data type
	MOV*	Move between two operands
	MOVZB*	Move a byte to a halfword or word, extending it with zeroes
	MOVA*	Move address of operand; data type is last
	PUSH*	Push operand onto stack
Arithmetic, logical		Operations on integer or logical bytes, halfwords (16 bits), words (32 bits); * is the data type
	ADD*_	Add with 2 or 3 operands
	CMP*	Compare and set condition codes
	TST*	Compare to zero and set condition codes
	ASH*	Arithmetic shift
	CLR*	Clear
	CVTB*	Sign extend byte to size of data type
Control		Conditional and unconditional branches
	BEQL, BNEQ	Branch equal/not equal
	BCC, BCS	Branch carry set, branch carry clear
	BRB, BRW	Unconditional branch with an 8-bit or 16-bit offset
	JMP	Jump using any addressing mode to specify target
	AOBLEQ	Add one to operand; branch if result \leq second operand
	CASE	Jump based on case selector
Procedure		Call/return from procedure
	CALLS	Call procedure with arguments on stack (see Section 3.9)
	CALLG	Call procedure with FORTRAN-style parameter list
	JSB	Jump to subroutine, saving return address
	RET	Return from procedure call
Bit-field character decimal		Operate on variable-length bit fields, character strings, and decimal strings, both in character and BCD format
	EXTV	Extracts a variable-length bit field into a 32-bit word
	MOVC3	Move a string of characters for given length
	CMPC3	Compare two strings of characters for given length
	MOVC5	Move string of characters with truncation or filling
	ADDP4	Add decimal string of the indicated length
	CVTPT	Convert packed-decimal string to character string
Floating point		Floating-point operations on D, F, G, and H formats
	ADDD_	Add double-precision D-format floating numbers
	SUBD_	Subtract double-precision D-format floating numbers
	MULF_	Multiply single-precision F-format floating point
	POLYF	Evaluate a polynomial using table of coefficients in F format
System		Change to system mode, modify protected registers
	CHMK, CHME	Change mode to kernel/executive
	REI	Return from exception or interrupt
Other		Special operations
	CRC	Calculate cyclic redundancy check
	INSQUE	Insert a queue entry into a queue

FIGURE 4.4 (Adjoining page) Classes of VAX instructions with examples. The asterisk stands for multiple data types—B, W, L, and usually D, F, G, H, and Q; remember how these VAX data types relate to the names used in the text (see Figure 4.2 on page 143). For example, a `MOVW` moves the VAX data-type word, which is 16 bits and is called a halfword in this text. The underline, as in `ADDD2`, means there are 2-operand (`ADDD2`) and 3-operand (`ADDD3`) forms of this instruction. The operand count is explicit in the opcode.

Operations on the VAX

What types of operators does the VAX provide? VAX operations can be divided into classes, as shown in Figure 4.4. (Detailed lists of the VAX instructions are included in Appendix B.) Figure 4.5 gives examples of typical VAX instructions and their meanings. Most instructions set the VAX condition codes according to their result; instructions without results, such as branches, do not. The condition codes are N (Negative), Z (Zero), V (oVerflow), and C (Carry).

Example assembly instruction	Length	Meaning
<code>MOVL @40(R4), 30(R2)</code>	5	$M[M[40+R4]] \leftarrow_{32} M[30+R2]$
<code>MOVAW R2, (R3)[R4]</code>	4	$R2 \leftarrow_{32} R3 + (R4 * 2)$
<code>ADDL3 R5, (R6)+, (R6)+</code>	4	$i \leftarrow M[R6]; R6 \leftarrow R6 + 4; R5 \leftarrow i + M[R6]; R6 \leftarrow R6 + 4$
<code>CMPL -(R6), #100</code>	7	$R6 \leftarrow R6 - 4$; Set the condition code using: $M[R6] - 100$
<code>CVTBW R10, (R8)</code>	3	$R10_{16..31} \leftarrow_{16} (M[R8]_0)^8 \text{ ## } M[R8]$
<code>BEQL name</code>	2	if equal(CC) {PC←name} $PC - 128 \leq \text{name} < PC + 128$
<code>BRW name</code>	3	PC←name $PC - 32768 \leq \text{name} < PC + 32768$
<code>EXTZV (R8), R5, R6, -564(R7)</code>	7	$t \leftarrow_{40} M[R7 - 564 + (R5 >> 3)]$; $i \leftarrow R5 \& 7$; $j \leftarrow$ if $R6 \geq 32$ then 32 else if $R6 < 0$ then 0 else $R6$; $M[R8] \leftarrow_{32} 0^{32-j} \text{ ## } t_{39-i-j+1..39-i}$
<code>MOV C3 @36(R9), (R10), 35(R11)</code>	6	$R1 \leftarrow 35 + R11$; $R3 \leftarrow M[36 + R9]$; for $(R0 \leftarrow M[R10]; R0 \neq 0; R0 --)$ $\{M[R3] \leftarrow_8 M[R1]; R1 ++; R3 ++\}$ $R2 = 0$; $R4 = 0$; $R5 = 0$
<code>ADD D3 R0, R2, R4</code>	4	$(R0 \text{ ## } R1) \leftarrow_{64} (R2 \text{ ## } R3) + (R4 \text{ ## } R5)$ register contents are type D floating point.

FIGURE 4.5 Some examples of typical VAX instructions. VAX assembly language syntax puts the result operand last; we have put it first for consistency with other machines. Instruction length is given in bytes. The condition `equal` (CC) is true if the condition-code setting reflects equality after a compare. Remember that most instructions set the condition code; the only function of compare instructions is to set the condition code. The names `t`, `i`, `j` are used as a temporaries in the instruction descriptions; `t` is 40 bits in length, while `i` and `j` are 32 bits. The `EXTZV` instruction may appear mysterious. Its purpose is to extract a variable-length field (0 to 32 bits) and zero extend it to 32 bits. The source operands to the `EXTZV` are the starting bit position (which may be any distance from the starting byte address), the length of the field, and the starting address of the bit string to extract the field from. The VAX numbers its bits from low order to high order, but we number bits in the reverse order. Thus, the subscripts adjust the bit offsets accordingly (which makes `EXTV` look more mysterious!). Although the result of the variable bit string operations are always 32 bits, the `MOV C3` changes the values of registers `R0` through `R5` as shown (although any of `R0`, `R2`, `R4`, and `R5` could be used to hold the count). A discussion of why `MOV C3` uses the GPRs as working registers appears in Section 5.6 of the next chapter.

4.3 The 360/370 Architecture

The IBM 360 was introduced in 1964. Its official goals included the following:

1. Exploit storage—large main storage, storage hierarchies (ROM used for microcode).
2. Support concurrent I/O—up to 5 MB/second with a standard interface on all machines.
3. Create a general-purpose machine with new OS facilities and many data types.
4. Maintain strict upward and downward machine-language compatibility.

The System/370, first introduced in 1970, was a successor to System/360. System/370 is fully upward compatible with System/360, even in system mode. The major extensions over the 360 included

- Virtual memory and dynamic address translation (see Chapter 8, Section 8.5)
- A few new instructions: synchronization support, long string instructions (long move and long compare), additional instructions for manipulating bytes in registers, and some additional decimal instructions
- Removal of data alignment requirements

In addition, several important implementation differences were introduced in the 370 implementations, including MOS main memory rather than core, and writeable control store (see Chapter 5).

In 1983, IBM introduced 370-XA, the eXtended Architecture. Until this extension, first used in the 3080 series, the 360/370 architecture had a 24-bit address space. Additional bits were added to the program status word so that the program counter could be extended. Unfortunately, it was common programming practice on the 360 to use the high-order byte of an address for status. Thus, old 24-bit programs cannot be run in 32-bit mode (actually a 31-bit address), while new and recompiled programs can take advantage of the larger address space. The I/O structure was also changed to permit higher levels of multiprocessing.

The latest extension to the architecture was ESA/370, introduced with the 3090 model in 1986. ESA/370 added additional instruction formats, called the Extended formats, with 16-bit opcodes. ESA/370 includes support for a Vector Facility (including a set of vector registers) and an extended (128-bit) floating-point format. The address space was extended by adding segments on top of the 31-bit address space (see Chapter 8, Sections 8.5 and 8.6); a new and more powerful protection model was added as well.

The remainder of this section surveys the IBM 360 architecture and presents measurements for the workload. First, let's examine the basics of the 360 architecture, then look at the instruction set formats and some sample instructions.

The 360/370 Instruction Set Architecture

The IBM System/360 is a 32-bit machine with byte addressability and support for a variety of data types: byte, halfword (16 bits), word (32 bits), doubleword (double-precision real), packed decimal, and unpacked character strings. The System/360 had alignment restrictions, which were removed in the System/370 architecture.

The internal state of the 360 has the following components:

- Sixteen 32-bit, general-purpose registers; register 0 is special when used in an addressing mode, where a zero is always substituted.
- Four double-precision (64-bit) floating-point registers.
- Program status word (PSW) holds the PC, some control flags, and the condition codes.

Later versions of the architecture extended this state with additional control registers.

Addressing Modes and Instruction Formats

The 360/370 has five instruction formats. Each format is associated with a single addressing mode and has a set of operations defined for that format. While some operations are defined in multiple formats, most are not. The instruction formats are shown in Figure 4.6 (page 150). While many instructions follow the paradigm of operating on sources and putting the result in a destination, other instructions (such as the control instructions BAL, BALR, BC) do not follow this paradigm, but use the same fields for other purposes. The associated addressing modes are as follows.

RR (*register-register*)—Both operands are simply contents of registers. The first source operand is also the destination.

RX (*register-indexed*)—The first operand and destination are a register. The second operand is the contents of the memory location given by the sum of a 12-bit displacement field D2, the contents of the register B2, and the contents of the register X2. This format is used when an index register is needed (and for most loads and stores).

RS (*register-storage*)—The first operand is a register that is the destination. The third operand is a register that is used as the second source. The second operand is the contents of the memory location given by the sum of the 12-bit displacement field D2 and the contents of the register B2. RS mode differs from RX in that a 3-operand form is supported, but the index register is eliminated. This instruction format is used for only a small number of instructions.

SI (storage-immediate)—The destination is a memory operand given by the sum of the contents of register B1 and the value of displacement D1. The second operand, an 8-bit immediate field, is the source.

SS (storage-storage)—The addresses of the two memory operands are the sum of the contents of a base register Bi and a displacement Di. The first operand is the destination. This storage-to-storage operation is used for decimal operations and for character strings. The length field can specify a single length of 1 to 256, or two lengths, each from 1 to 16. A single length is used for string instructions, while decimal instructions specify a length for each operand.

The displacement in the RS, RX, SI, and SS formats is 12 bits and is unsigned.

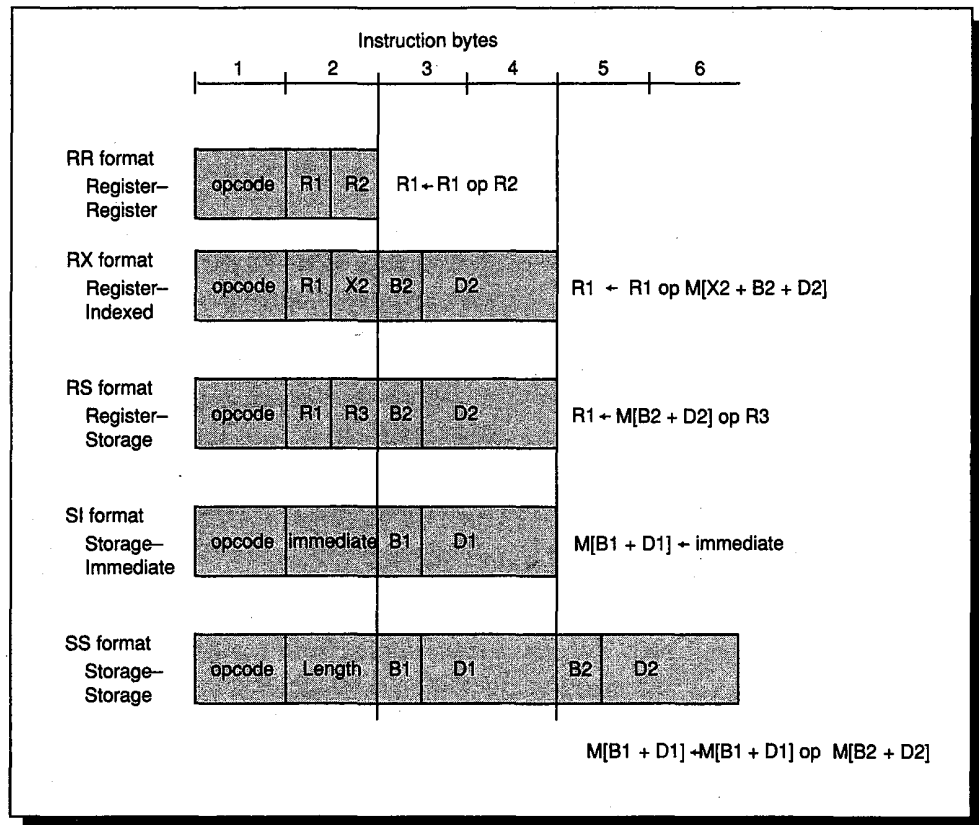


FIGURE 4.6 The 360/370 instruction formats. The possible instruction operands are a register (R1, R2, or R3), an 8-bit immediate, or a memory location. The opcode specifies where the operands reside and the addressing mode. The effective addresses for memory operands are formed using the sum of one or two registers (called B1, B2, or X2) and a 12-bit unsigned displacement field (called D1 or D2). In addition, the storage-storage instructions, which are all string-oriented, specify an 8-bit length field. Other instruction formats have been added in later architectural extensions. These formats allowed the opcode space to be extended and new data types to be added. For loads, stores, and moves only one source operand is used and the operation only moves the data (see Figure 4.8 on page 152). For SS instructions, the length is one greater than the value in the instruction.

Operations on the 360/370

Just as on the VAX, the instructions on the 360 can be divided into classes. Four basic types of operations on data are supported:

1. *Logical operations on bits, character strings, and fixed words.* These are mostly RR and RX formats with a few RS instructions.
2. *Decimal or character operations on strings of characters or decimal digits.* These are SS format instructions.
3. *Fixed-point binary arithmetic.* This is supported in both RR and RX formats.
4. *Floating-point arithmetic.* This is supported primarily with RR and RX instructions.

Branches use the RX instruction format with the effective address specifying the branch target. Since branches are not PC-relative, a base register may need to be loaded to specify the branch target. This has a rather substantial impact: in general, it means that there must be registers that point to every region containing a branch target. The condition codes are set by all arithmetic and logical operations. Conditional branches test the condition codes under a mask to determine whether or not to branch.

Some example instructions and their formats are shown in Figure 4.7. When an operation is defined for more than one format, separate opcodes are used to specify the instruction format. For example, the opcode AR (add register) says that the instruction type is RR; thus, the operands are in registers. The opcode A (add) says the format is RX; thus, one operand is in memory, accessed with the RX addressing mode. Figure 4.8 (page 152) has a longer listing of operations, including all the most common ones; a full table of instructions appears in Appendix B.

Type	Instruction example	Meaning
RR	AR R4, R5	$R4 \leftarrow R4 + R5$
RX	A R4, 10(R5, R6)	$R4 \leftarrow R4 + M[R5 + R6 + 10]$
RX	BC Mask, 20(R5, R6)	if (CC & Mask) != 0 {PC ← 20 + R5 + R6}
RS	STM 20(R14), R2, R8	for (i=2; i<=8; i++) {M[R14+20+(i-2)*4] ← ₃₂ Ri}
SI	MVI 20(R5), #40	$M[R5+20] \leftarrow \text{8 } 40$
SS	MVC 10(R2), Len, 20(R6)	for (i=0; i<Len+1; i++) {M[R2+10+i] ← ₈ M[R6+20+i]}

FIGURE 4.7 Typical IBM 360 instructions with their meanings. The MVC instruction is shown with the length as the second operand. The length field is a constant in the instruction; standard 360 assembly language syntax includes the length with the first operand. The variable *i* used in the MVC and STM is a temporary.

Class or instruction	Format	Instruction meaning
Control		
Change the PC		
BC_	RX,RR	Test the condition and conditionally branch
BAL_	RX,RR	Branch and link (address of next instruction is placed in R15)
Arithmetic, logical		
Arithmetic and logical operations		
A_	RX,RR	Add
S_	RX,RR	Subtract
SLL	RS	Shift left logical; shifts a register by an immediate amount
LA	RX	Load address—put effective address into destination
CLI	SI	Compare storage byte against immediate
NI	SI	AND immediate into storage byte
C_	RX,RR	Compare and set condition codes
TM	RS	Test under mask—perform a logical AND of the operand and an immediate field; set condition codes based on the result
MH	RX	Multiply halfword
Data transfer		
Moves between registers or register and memory		
L_	RX,RR	Load a register from memory or another register
MVI	SI	Store an immediate byte in memory
ST	RX	Store a register
LD	RX	Load a double-precision floating-point register
STD	RX	Store a double-precision floating-point register
LPDR	RR	Move a double-precision floating-point register to another
LH	RX	Load a halfword from memory into a register
IC	RX	Insert a memory byte into low-order byte of a register
LTR	RR	Load a register and set condition codes
Floating point		
Floating-point operations		
AD_	RX,RR	Double-precision floating-point add
MD_	RX,RR	Double-precision FP multiply
Decimal, string		
Operations on decimal and character strings		
MVC	SS	Move characters
AP	SS	Add packed-decimal strings, replacing first with sum
ZAP	SS	Zero and add packed—replace destination with source
CVD	RX	Convert a binary word to decimal doubleword
MP	SS	Multiply two packed-decimal strings
CLC	SS	Compare two character strings
CP	SS	Compare two packed-decimal strings
ED	SS	Edit—convert packed-decimal to character string

FIGURE 4.8 Most frequently used IBM 360 instructions. The underline means that the opcode is two distinct opcodes with an RX format and an RR format. For example A_ stands for AR and A. The full instruction set is shown in Appendix B.

4.4 The 8086 Architecture

The Intel 8086 architecture was announced in 1978 as an upward-compatible extension of the then-successful 8080. Whereas the 8080 was a straightforward accumulator machine, the 8086 extended the architecture with additional registers. The 8086 fails to be a truly general-purpose register machine, however, because nearly every register has a dedicated use. Thus, its architecture falls somewhere between an accumulator machine and a general-purpose register machine. The 8086 is a 16-bit architecture; all internal registers are 16 bits. To obtain addressability greater than 16 bits the designers added segments to the architecture. This allowed a 20-bit address space, broken into 64-KB fragments. Chapter 8 discusses segmentation in detail, while this chapter will focus only on the implications for a compiler.

The 80186, 80286, 80386, and 80486 are “compatible” extensions of the 8086 architecture and are collectively referred to as the 80x86 processors. They are compatible in the sense that they all belong to the same architectural family. There are more instances of this architectural family than of any other in the world. The 80186 added a small number of extensions (about 16) to the 8086 architecture in 1981. The 80286, introduced in 1982, extended the 80186 architecture by creating an elaborate memory-mapping and protection model and by extending the address space to 24 bits (see Chapter 8, Section 8.6). Because 8086 programs needed to be binary compatible, the 80286 offered a real addressing mode to make the machine look just like an 8086.

The 80386 was introduced in 1985. It is a true 32-bit machine when running in native mode. Like the 80286, a real addressing mode is provided for 8086 compatibility. There is also a virtual 8086 mode that provides for multiple 20-bit 8086 address partitions within the 80386’s memory. In addition to a 32-bit architecture with 32-bit registers and a 32-bit address space, the 80386 has a new set of addressing modes and additional operations. The added instructions make the 80386 nearly a general-purpose register machine—for most operations any register can be used as an operand. The 80386 also provides paging support (see Chapter 8). The 80486 was introduced in 1989 and added only a few new instructions, while substantially increasing performance.

Since 8086 compatibility mode is the dominant use of all 80x86 processors, we will take a detailed look in this section at the 8086 architecture. We will begin by summarizing the architecture and then discuss its usage by typical programs.

8086 Instruction Set Summary

The 8086 provides support for both 8-bit (byte) and 16-bit (called word) data types. The data type distinctions apply to register operations as well as memory accesses.

The address space on the 8086 is a total of 20 bits; however, it is broken into 64-KB segments addressable with 16-bit offsets. A 20-bit address is formed by taking a 16-bit effective address—as an offset within a segment—and adding it to a 16-bit segment base address. The segment base address is obtained by shifting the contents of a 16-bit segment register 4 bits to the left.

Class	Register	Purposes of class or register
Data		Used to hold and operate on data
	AX	Used for multiply, divide, and I/O; sometimes an implicit operand; AH and AL also have dedicated uses in byte multiply, divide, decimal arithmetic
	BX	Can also be used as address-base register
	CX	Used for string operations and loop instructions; CL is the dynamic shift count
	DX	Used for multiply, divide, and I/O
Address		Used to form 16-bit effective memory addresses (within segment)
	SP	Stack pointer
	BP	Base register—used in based-addressing mode
	SI	Index register, and also used as string source base register
	DI	Index register, and also used as string destination base register
Segment		Used to form 20-bit real memory addresses
	CS	Code segment—used with instruction access
	SS	Stack segment—used for stack references (SP) or when BP is base register
	DS	Data segment—used when a reference is not for code (CS used), to the stack (SS used), or a string destination (ES used)
	ES	Extra segment—used when operand is string destination
Control		Used for status and program control
	IP	Instruction pointer—provides the offset of the currently executing instruction (this is the lower 16-bits of the effective PC)
	FLAGS	Contains six condition code bits—carry, zero, sign, borrow, parity, and overflow—and three status control bits

FIGURE 4.9 The 14 registers on the 8086. The table divides them into four classes that have restricted uses. In addition, many of the individual registers are required for certain instructions. The data registers have an upper and lower half: xL refers to lower byte and xH to upper byte of register x.

The 8086 provides a total of 14 registers broken into four groups—data registers, address registers, segment registers, and control registers—as shown in Figure 4.9. The segment register for a memory access is usually implied by the base register used to form the effective address within the segment.

The addressing modes for data on the 8086 use the segment registers implied by the addressing mode or specified in the instruction with an override of the default mode. We will discuss how branches and jumps deal with segmentation in the section on operations.

Addressing Modes

Most of the addressing modes for forming the effective address of a data operand are among those discussed in Chapter 3. The arithmetic, logical, and data-transfer instructions are two-operand instructions that allow the combinations shown in Figure 4.10.

Source/destination operand type	Second source operand
Register	Register
Register	Immediate
Register	Memory
Memory	Register
Memory	Immediate

FIGURE 4.10 Instruction types for the arithmetic, logical, and data-transfer instructions. The 8086 allows the combinations shown. Immediates may be 8 or 16 bits in length; a register is any one of the 12 major registers in Figure 4.9 (not one of the control registers). The only restriction is the absence of memory–memory mode.

The memory addressing modes supported are absolute (16-bit absolute address), register indirect, based, indexed, and based indexed with displacement (not mentioned in Chapter 3). Although a memory operand can use any addressing mode, there are restrictions on what registers can be used in a mode. The registers usable in specifying the effective address are as follows:

- *Register indirect*—BX, SI, DI.
- *Based mode with 8-bit or 16-bit displacement*—BP, BX, SI, DI. (Intel gives two names to this addressing mode, Based and Indexed, but they are essentially identical and we combine them.)
- *Indexed*—address is sum of two registers. The allowable combinations are BX+SI, BX+DI, BP+SI, and BP+DI. This mode is called Based Indexed on the 8086.

- *Based indexed with 8-bit or 16-bit displacement*—the address is sum of displacement and contents of two registers. The same restrictions on registers apply as in indexed mode.

Operations on the 8086

The 8086 operations can be divided into four major classes:

1. Data movement instructions, including move, push, and pop
2. Arithmetic and logic instructions, including logical operations, test, shifts, and integer and decimal arithmetic operations
3. Control flow, including conditional and unconditional branches, calls, and returns
4. String instructions, including string move and string compare

Instruction	Function
JE name	if equal(CC) {IP←name}; IP-128 ≤ name < IP+128
JMP name	IP←name
CALLF name, seg	SP←SP-2; M[SS:SP]←CS; SP←SP-2; M[SS:SP]←IP+5; IP←name; CS←seg;
MOVW BX, [DI+45]	BX← ₁₆ M[DS:DI+45]
PUSH SI	SP←SP-2; M[SS:SP]←SI
POP DI	DI←M[SS:SP]; SP←SP+2
ADD AX, #6765	AX←AX+6765
SHL BX, 1	BX←BX _{1..15} ## 0
TEST DX, #42	Set CC flags with DX & 42
MOVSB	M[ES:DI]← ₈ M[DS:SI]; DI←DI+1; SI←SI+1

FIGURE 4.11 Some typical 8086 instructions and their functions. A list of the most frequent operations appears in Figure 4.12 (page 158). We use the abbreviation SR:X to indicate the formation of an address with segment register SR and offset X. This effective address corresponding to SR:X is (SR<<4)+X. The CALLF saves the IP of the next instruction and the current CS on the stack.

In addition, there is a repeat prefix that may precede any string instruction, which says that the instruction should be repeated using the value in the CX register for the number of repetitions. Figure 4.11 shows some typical 8086 instructions and their functions.

Control-flow instructions must be able to address destinations in another segment. This is handled by having two types of control-flow instructions: “near” for intrasegment (within a segment) and “far” for intersegment (between segments) transfers. In far jumps, which must be unconditional, two 16-bit quantities follow the opcode. One of these is used as the instruction pointer, while the other is loaded into CS and becomes the new code segment. Calls and returns work similarly—a far call pushes the return instruction pointer and return segment on the stack and loads both the instruction pointer and code segment. A far return pops both the instruction pointer and the code segment from the stack. Programmers or compiler writers must be sure to always use the same type of call **and** return for a procedure—a near return does not work with a far call, and vice versa.

Figure 4.12 (page 158) summarizes the most popular 8086 instructions. Many of the instructions are available in both byte and word formats. A full listing of instructions appears in Appendix B.

The encoding of instructions in the 8086 is complex, and there are many different instruction formats. Instructions may vary from one byte, when there are no operands, up to six bytes, when the instruction contains a 16-bit immediate and uses 16-bit displacement addressing. Figure 4.13 (page 159) shows the instruction format for several of the example instructions in Figure 4.11 (page 156). The opcode byte usually contains a bit saying whether the instruction is a word or byte instruction. For some instructions the opcode may include the addressing mode and the register; this is true in many instructions that have the form “register←register op immediate.” For other instructions a “postbyte” or extra opcode byte contains the addressing mode information. This postbyte is used for many of the instructions that address memory. The encoding of the postbyte is shown in Figure 4.14 (page 160). Finally, there is a byte prefix that is used for three different purposes. It can override the default-segment usage of instructions, and it can be used to repeat a string instruction by a count provided in CX. (This latter function is useful for string instructions that operate on a single byte or word at a time and use autoincrement addressing.) Third, it can be used to generate an atomic memory access for use in implementing synchronization.

Instruction	Meaning
Control	Conditional and unconditional branches
JNZ, JZ	Jump if condition to IP + 8-bit offset; JNE (for JNZ), JE (for JZ) are alternative names
JMP, JMPF	Unconditional jump—8-bit or 16-bit offset intrasegment (near), and intersegment (far) versions
CALL, CALLF	Subroutine call—16-bit offset; return address pushed; near and far versions
RET, RETF	Pops return address from stack and jumps to it; near and far versions
LOOP	Loop branch—decrement CX; jump to IP + 8-bit displacement if CX ≠ 0
Data transfer	Move data between registers or between register and memory
MOV	Move between two registers or between register and memory
PUSH	Push source operand on stack
POP	Pop operand from stack top to a register
LES	Load ES and one of the GPRs from memory
Arithmetic, logical	Arithmetic and logical operations using the data registers and memory
ADD	Add source to destination; register–memory format
SUB	Subtract source from destination; register–memory format
CMP	Compare source and destination; register–memory format
SHL	Shift left
SHR	Shift logical right
RCL	Rotate right with Carry as fill
CBW	Convert byte in AL to word in AX
TEST	Logical AND of source and destination sets flags
INC	Increment destination; register–memory format
DEC	Decrement destination; register–memory format
OR	Logical OR; register–memory format
XOR	Exclusive OR; register–memory format
String instructions	Move between string operands; length given by a repeat prefix
MOVS	Copies from string source to destination; may be repeated
LODS	Loads a byte or word of a string into the A register

FIGURE 4.12 Some typical operations on the 8086. Many operations use register–memory format, where either the source or the destination may be memory and the other may be a register or immediate operand.

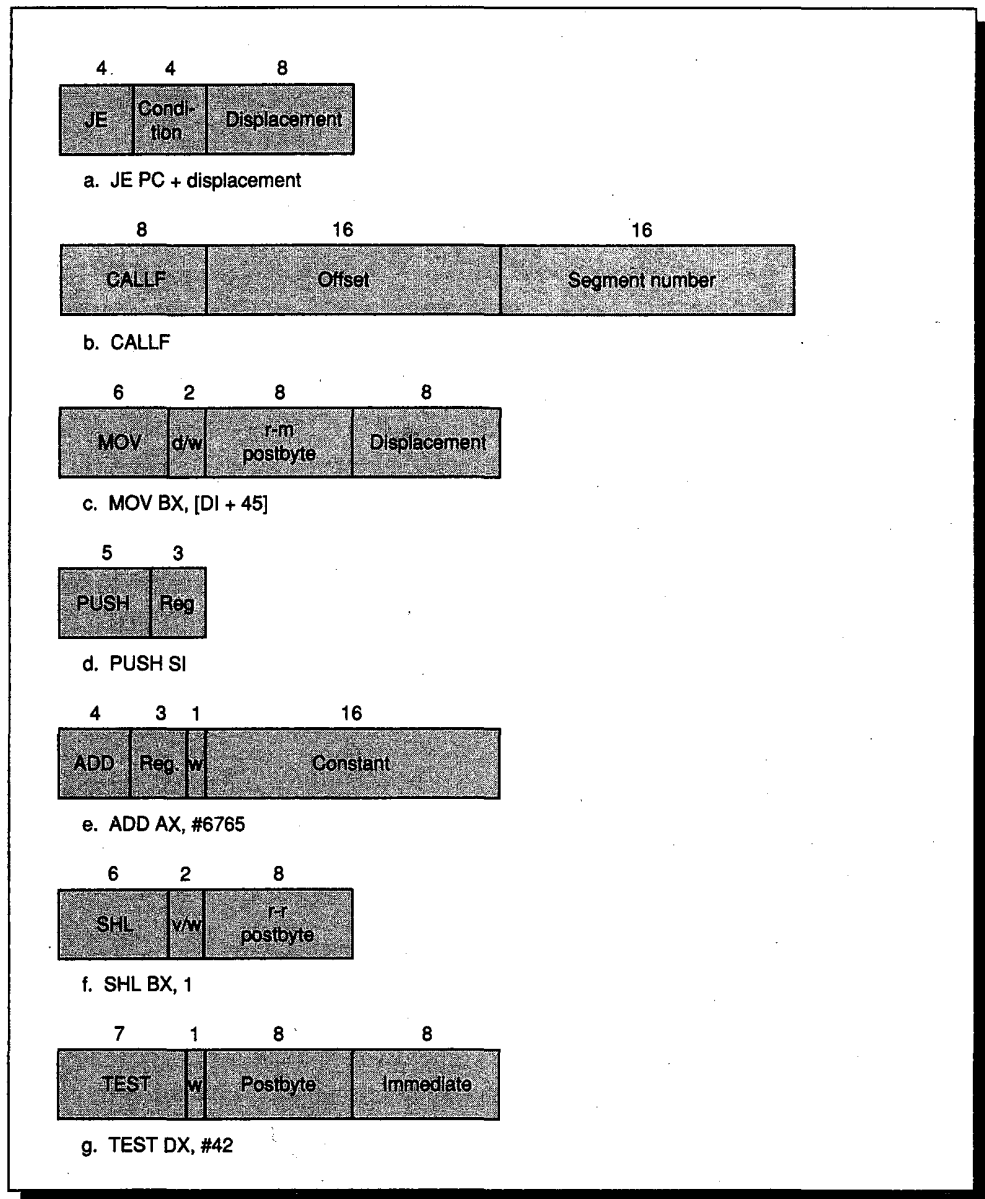


FIGURE 4.13 Typical 8086 instruction formats. The encoding of the postbyte is shown in Figure 4.14. Many instructions contain the 1-bit field *w*, which says whether the operation is a byte or word. Fields of the form *v/w* or *d/w* are a *d*-field or *v*-field followed by the *w*-field. The *d*-field in *MOV* is used in instructions that may move to or from memory and shows the direction of the move. The field *v* in the *SHL* instruction indicates a variable-length shift; variable-length shifts use a register to hold the shift count. The *ADD* instruction shows a typical optimized short encoding usable only when the first operand is *AX*. Overall instructions may vary from one to six bytes in length.

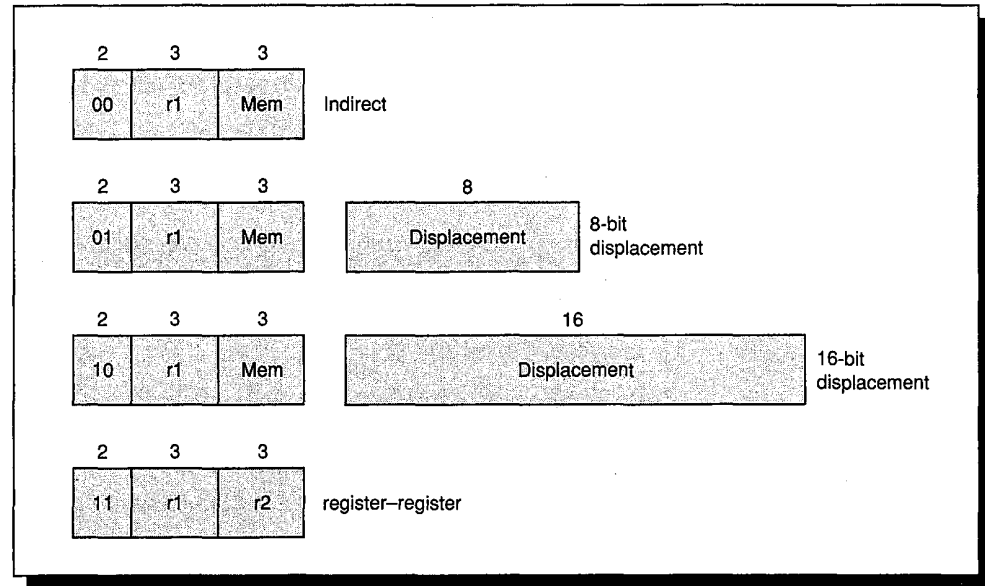


FIGURE 4.14 There are four postbyte encodings on the 8086 designated by a 2-bit tag. The first three indicate a register–memory instruction, where Mem is the base register. The fourth form is register–register.

4.5 The DLX Architecture

In many places throughout this book we will have occasion to refer to a computer's "machine language." The machine we use is a mythical computer called "MIX." MIX is very much like nearly every computer in existence, except that is, perhaps, nicer ... MIX is the world's first polyunsaturated computer. Like most machines, it has an identifying number—the 1009. This number was found by taking 16 actual computers which are very similar to MIX and on which MIX can be easily simulated, then averaging their number with equal weight:

$$\lfloor (360 + 650 + 709 + 7070 + U3 + SS80 + 1107 + 1604 + G20 + B220 + S2000 + 920 + 601 + H800 + PDP-4 + 11) / 16 \rfloor = 1009.$$

The same number may be obtained in a simpler way by taking Roman numerals.

Donald Knuth, The Art of Computer Programming. Volume I: Fundamental Algorithms

In this section we will describe a simple load/store architecture called DLX (pronounced “Deluxe”). The authors believe DLX to be the world’s second polyunsaturated computer—the average of a number of recent experimental and commercial machines that are very similar in philosophy to DLX. Like Knuth, we derived the name of our machine from an average expressed in Roman numerals:

(AMD 29K, DECstation 3100, HP 850, IBM 801, Intel i860, MIPS M/120A, MIPS M/1000, Motorola 88K, RISC I, SGI 4D/60, SPARCstation-1, Sun-4/110, Sun-4/260) / 13 = 560 = DLX.

The architecture of DLX was chosen based on observations about the most frequently used primitives in programs. More sophisticated (and less performance-critical) functions are implemented in software with multiple instructions. In Section 4.9 we discuss how and why these architectures became popular.

Like most recent load/store machines, DLX emphasizes

- A simple load/store instruction set
- Design for pipelining efficiency (discussed in Chapter 6)
- An easily decoded instruction set
- Efficiency as a compiler target

DLX provides a good architectural model for study, not only because of the recent popularity of this type of machine, but also because it is an easy architecture to understand.

DLX—Our Generic Load/Store Architecture

In this section, the DLX instruction set is defined. We will use this architecture again in Chapters 5 through 7, and it forms the basis for a number of exercises and programming projects.

- The architecture has thirty-two 32-bit general-purpose registers (GPRs); the value of R0 is always 0. Additionally, there are a set of floating-point registers (FPRs), which can be used as 32 single-precision (32-bit) registers, or as even-odd pairs holding double-precision values. Thus, the 64-bit floating-point registers are named F0, F2, ..., F28, F30. Both single- and double-precision operations are provided. There are a set of special registers used for accessing status information. The FP status register is used for both compares and FP exceptions. All movement to/from the status register is through the GPRs; there is a branch that tests the comparison bit in the FP status register.

- Memory is byte addressable in Big Endian mode with a 32-bit address. All memory references are through loads or stores between memory and either the GPRs or the FPRs. Accesses involving the GPRs can be to a byte, to a halfword, or to a word. The FPRs may be loaded and stored with single-precision or double-precision words (using a pair of registers for DP). All memory accesses must be aligned. There are also instructions for moving between a FPR and a GPR.
- All instructions are 32 bits and must be aligned.
- There are also a few special registers that can be transferred to and from the integer registers. An example is the floating-point status register, used to hold information about the results of floating-point operations.

Operations

There are four classes of instructions: loads and stores, ALU operations, branches and jumps, and floating-point operations.

Example instruction	Instruction name	Meaning
LW R1, 30 (R2)	Load word	$R1 \leftarrow_{32} M[30+R2]$
LW R1, 1000 (R0)	Load word	$R1 \leftarrow_{32} M[1000+0]$
LB R1, 40 (R3)	Load byte	$R1 \leftarrow_{32} (M[40+R3]_0)^{24} \#\# M[40+R3]$
LBU R1, 40 (R3)	Load byte unsigned	$R1 \leftarrow_{32} 0^{24} \#\# M[40+R3]$
LH R1, 40 (R3)	Load halfword	$R1 \leftarrow_{32} (M[40+R3]_0)^{16} \#\#M[40+R3] \#\#M[41+R3]$
LHU R1, 40 (R3)	Load halfword unsigned	$R1 \leftarrow_{32} 0^{16} \#\#M[40+R3] \#\#M[41+R3]$
LF F0, 50 (R3)	Load float	$F0 \leftarrow_{32} M[50+R3]$
LD F0, 50 (R2)	Load double	$F0 \#\#F1 \leftarrow_{64} M[50+R2]$
SW 500 (R4), R3	Store word	$M[500+R4] \leftarrow_{32} R3$
SF 40 (R3), F0	Store float	$M[40+R3] \leftarrow_{32} F0$
SD 40 (R3), F0	Store double	$M[40+R3] \leftarrow_{32} F0; M[44+R3] \leftarrow_{32} F1$
SH 502 (R2), R3	Store half	$M[502+R2] \leftarrow_{16} R3_{16..31}$
SB 41 (R3), R2	Store byte	$M[41+R3] \leftarrow_8 R2_{24..31}$

FIGURE 4.15 The load and store instructions in DLX. All use a single addressing mode and require that the memory value be aligned. Of course, both loads and stores are available for all the data types shown.

Any of the general-purpose or floating-point registers may be loaded or stored, except that loading R0 has no effect. There is a single addressing mode, base register + 16-bit signed offset. Halfword and byte loads place the loaded object in the lower portion of the register. The upper portion of the register is filled with either the sign extension of the loaded value or zeros, depending on the opcode. Single-precision floating-point numbers occupy a single floating-point register, while double-precision values occupy a pair. Conversions between single and double precision must be done explicitly. The floating-point format is IEEE 754 (see Appendix A). Figure 4.15 gives an example of the load and store instructions. A complete list of the instructions appears in Figure 4.18 (page 165).

All ALU instructions are register–register instructions. The operations include simple arithmetic and logical operations: add, subtract, AND, OR, XOR, and shifts. Immediate forms of all these instructions, with a 16-bit sign-extended immediate, are provided. The operation LHI (load high immediate) loads the top half of a register, while setting the lower half to 0. This allows a full 32-bit constant to be built in two instructions. (We sometimes use the mnemonic LI, standing for Load Immediate, as an abbreviation for an add immediate where one of the source operands is R0; likewise, the mnemonic MOV is sometimes used for an ADD where one of the sources is R0.)

There are also compare instructions, which compare two registers ($=, \neq, <, >, \leq, \geq$). If the condition is true, these instructions place a 1 in the destination register (to represent true); otherwise they place the value 0. Because these operations “set” a register they are called set-equal, set-not-equal, set-less-than, and so on. There are also immediate forms of these compares. Figure 4.16 gives some examples of the arithmetic/logical instructions.

Control is handled through a set of jumps and a set of branches. The four jump instructions are differentiated by the two ways to specify the destination address and by whether or not a link is made. Two jumps use a 26-bit signed offset added to the program counter (of the instruction sequentially following the jump) to determine the destination address; the other two jump instructions specify a register that contains the destination address. There are two flavors of jumps: plain jump, and jump and link (used for procedure calls). The latter places the return address in R31.

Example instruction	Instruction name	Meaning
ADD R1, R2, R3	Add	$R1 \leftarrow R2 + R3$
ADDI R1, R2, #3	Add immediate	$R1 \leftarrow R2 + 3$
LHI R1, #42	Load high immediate	$R1 \leftarrow 42 \# 0^{16}$
SLLI R1, R2, #5	Shift left logical	$R1 \leftarrow R2 \ll 5$
SLT R1, R2, R3	Set less than	if $(R2 < R3)$ $R1 \leftarrow 1$ else $R1 \leftarrow 0$

FIGURE 4.16 Examples of arithmetic/logical instructions on DLX, both with and without immediates.

Example instruction	Instruction name	Meaning
J name	Jump	$PC \leftarrow \text{name}; ((PC+4) - 2^{25}) \leq \text{name} < ((PC+4) + 2^{25})$
JAL name	Jump and link	$R31 \leftarrow PC+4; PC \leftarrow \text{name}; ((PC+4) - 2^{25}) \leq \text{name} < ((PC+4) + 2^{25})$
JALR R2	Jump and link register	$R31 \leftarrow PC+4; PC \leftarrow R2$
JR R3	Jump register	$PC \leftarrow R3$
BEQZ R4, name	Branch equal zero	if $(R4 == 0)$ $PC \leftarrow \text{name}; ((PC+4) - 2^{15}) \leq \text{name} < ((PC+4) + 2^{15})$
BNEZ R4, name	Branch not equal zero	if $(R4 \neq 0)$ $PC \leftarrow \text{name}; ((PC+4) - 2^{15}) \leq \text{name} < ((PC+4) + 2^{15})$

FIGURE 4.17 Typical control-flow instructions in DLX. All control instructions, except jumps to an address in a register, are PC-relative. If the register operand is R0, the branch is unconditional, but the compiler will usually prefer to use a jump with a longer offset over this "unconditional branch."

All branches are conditional. The branch condition is specified by the instruction, which may test the register source for zero or nonzero; this may be a data value or the result of a compare. The branch target address is specified with a 16-bit signed offset that is added to the program counter. Figure 4.17 gives some typical branch and jump instructions.

Floating-point instructions manipulate the floating-point registers and indicate whether the operation to be performed is single or double precision. Single-precision operations can be applied to any of the registers, while double-precision operations apply only to an even-odd pair (e.g., F4, F5), which is designated by the even register number. Load and store instructions for the floating-point registers move data between the floating-point registers and memory both in single and double precision. The operations MOVF and MOVD copy a single-precision (MVF) or double-precision (MOVD) floating-point register to another register of the same type. The operations MOVFP2I and MOVI2FP move data between a single floating-point register and an integer register; moving a double-precision value to two integer registers require two instructions. Integer multiply and divide that work on 32-bit floating-point registers are also provided, as are conversions from integer to floating point and vice versa.

The floating-point operations are add, subtract, multiply, and divide; a suffix D is used for double precision and a suffix F is used for single precision (e.g., ADDD, ADDE, SUBD, SUBF, MULTD, MULTF, DIVD, DIVF). Floating-point compares set a bit in the special floating-point status register that can be tested with a pair of branches: BFPT and BFPE, branch floating point true and branch floating point false.

Figure 4.18 contains a list of all operations and their meaning.

Instruction type / opcode	Instruction meaning
Data transfers	Move data between registers and memory, or between the integer and FP or special registers; only memory address mode is 16-bit displacement + contents of a GPR
LB, LBU, SB	Load byte, load byte unsigned, store byte
LH, LHU, SH	Load halfword, load halfword unsigned, store halfword
LW, SW	Load word, store word (to/from integer registers)
LF, LD, SF, SD	Load SP float, load DP float, store SP float, store DP float
MOVI2S, MOVS2I	Move from/to GPR to/from a special register
MOVFP, MOVDP	Copy one floating-point register or a DP pair to another register or pair
MOVFP2I, MOVI2FP	Move 32 bits from/to FP registers to/from integer registers
Arithmetic / Logical	Operations on integer or logical data in GPRs; signed arithmetics trap on overflow
ADD, ADDI, ADDU, ADDUI	Add, add immediate (all immediates are 16 bits); signed and unsigned
SUB, SUBI, SUBU, SUBUI	Subtract, subtract immediate; signed and unsigned
MULT, MULTU, DIV, DIVU	Multiply and divide, signed and unsigned; operands must be floating-point registers; all operations take and yield 32-bit values
AND, ANDI	And, and immediate
OR, ORI, XOR, XORI	Or, or immediate, exclusive or, exclusive or immediate
LHI	Load high immediate—loads upper half of register with immediate
SLL, SRL, SRA, SLLI, SRLI, SRAI	Shifts: both immediate (S__I) and variable form (S__); shifts are shift left logical, right logical, right arithmetic
S__, S__I	Set conditional: “__” may be LT, GT, LE, GE, EQ, NE
Control	Conditional branches and jumps; PC-relative or through register
BEQZ, BNEZ	Branch GPR equal/not equal to zero; 16-bit offset from PC+4
BFPT, BFPF	Test comparison bit in the FP status register and branch; 16-bit offset from PC+4
J, JR	Jumps: 26-bit offset from PC (J) or target in register (JR)
JAL, JALR	Jump and link: save PC+4 to R31, target is PC-relative (JAL) or a register (JALR)
TRAP	Transfer to operating system at a vectored address; see Chapter 5
RFE	Return to user code from an exception; restore user mode; see Chapter 5
Floating point	Floating-point operations on DP and SP formats
ADDD, ADDF	Add DP, SP numbers
SUBD, SUBF	Subtract DP, SP numbers
MULTD, MULTF	Multiply DP, SP floating point
DIVD, DIVF	Divide DP, SP floating point
CVTF2D, CVTF2I, CVTD2F, CVTD2I, CVTI2F, CVTI2D	Convert instructions: CVT _x 2 _y converts from type x to type y, where x and y are one of I (Integer), D (Double precision), or F (Single precision). Both operands are in the FP registers
__D, __F	DP and SP compares: “__” may be LT, GT, LE, GE, EQ, NE; sets comparison bit in FP status register

FIGURE 4.18 Complete list of the instructions in DLX. The formats of these instructions are shown in Figure 4.19. This list can also be found in the back inside cover.

Instruction Format

All instructions are 32 bits with a 6-bit primary opcode. Figure 4.19 shows the instruction layout.

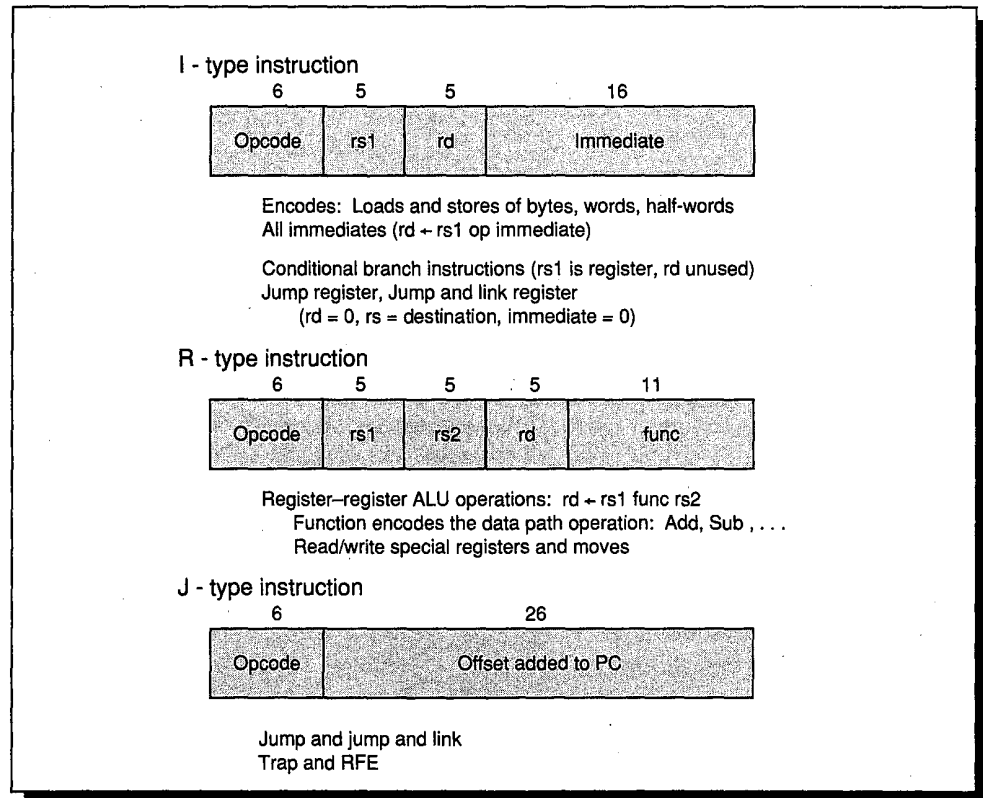


FIGURE 4.19 Instruction layout for DLX. All instructions are encoded in one of three types.

Machines Related to DLX

Between 1985 and 1990 many load/store machines were announced that are similar to DLX. Figure 4.20 describes the major features of these machines. All have 32-bit instructions and are load/store architectures; the figure lists their differences. These machines are all very similar—if you're not convinced, try making a table such as this one comparing these machines to the VAX or 8086.

DLX bears a close resemblance to all the other load/store machines shown in Figure 4.20. (See Appendix E for a detailed description of four load/store machines closely related to DLX.) Thus, the measurements in the next section will be reasonable approximations of the behavior of any of the machines. In fact, some studies suggest that compiler differences are more significant than architectural differences among these machines.

Machine	Registers	Addressing modes	Operations
DLX	32 integer; 16 DP or 32 SP FP	16-bit displacement; 16-bit immediates	See Figure 4.18.
AMD 29000	192 integer with stack cache; 8 DP FP	Register deferred only; 8-bit immediates	Integer multiply/divide trap to software. Branches =, ≠ 0 only.
HP Precision Architecture	32 GPRs 32 DP or 64 SPFP	5-bit, 14-bit, and 32-bit displacements; scaled mode (load only); autoincrement; autodecrement	Every ALU operation can skip the next in- struction. Many special bit-manipulation instructions. 32-bit immediates; decimal- support instructions: integer multiply/divide not single instructions. Stores of partial word. 64-bit addresses possible through segmentation.
Intel i860	32 integer; 16 DP or 32 SP FP	16-bit displacement; indexed mode; autoincrement; 16-bit immediates	Branch compares two registers for equality. Conditional traps are supported. FP reciprocal rather than divide. Some support for 128-bit loads and stores.
MIPS R2000 / R3000	32 integer; 16 FP	16-bit displacement; 16-bit immediates	Floating-point load/store moves 32 bits to upper or lower half of FP register. Branch condition can compare two registers. Integer multiply/divide in GPRs. Special instructions for partial word load/store.
Motorola 88000	32 GPRs	16-bit displacement; indexed mode	Special bit-manipulation instructions. Branches can test for zero and also test bits set by compares.
SPARC	Register windows with 32 integer registers available per procedure; 16 DP or 32 SP FP	13-bit offset and 13-bit immediates; indexed addressing mode	Branches use condition code, set selectively by instructions. Integer multiply/divide not instructions. No moves between integer and FP registers.

FIGURE 4.20 Comparison of the major features of a variety of recent load/store architectures. All the machines have a basic instruction size of 32 bits, though some provisions for shorter or longer are supported. For example, the Precision Architecture uses 2-word instructions for long immediates. Register windows and stack caches, which are used in the SPARC and AMD 29000 architectures, are discussed in Chapter 8. The MIPS R2000 is used in the DECstation 3100, the machine benchmarked in Chapter 2, and used as the load/store machine in Chapter 3. The number of double-precision floating-point registers is indicated if they are separate from the integer registers. Appendix E has a detailed comparative description of DLX, the MIPS R2000, SPARC, the i860, and the 88000 architectures. Both the MIPS and SPARC architectures have extensions that were not supported in hardware in the first implementation. These are discussed in Appendix E. In several of these machines $R0=0$, so they really have one less register available.

4.6

Putting It All Together: Measurements of Instruction Set Usage

In this section we examine the dynamic use of the four instruction sets presented in this chapter. All instructions responsible for 1.5% or more of the instruction

executions in a set of benchmarks are included in measurements of each architecture. In the interest of conciseness, fractional percents are rounded so that all entries in the graphs of opcode frequency will be at least 2%.

To facilitate comparisons among dynamic instruction set measurements, the measurements are organized by class of application. Figure 4.21 shows these application classes and the programs used to obtain instruction-use data on each of the machines discussed. We sometimes compare data for different architectures running the same type of application (e.g., a compiler) **but** different programs. The reader is cautioned that such comparisons must be made cautiously and with substantial limitations. Despite the fact that both programs may be the same type of application, differences in programming language, coding style, compilers, and so on, could substantially affect the results.

Machines	Compilers	Floating point	General integer	Business data processing
VAX	GCC	Spice	TeX	COBOLX
360	PL/I	FORTGO	PLIGO	COBOLGO
8086	Turbo C		Assembler	Lotus 1-2-3
DLX	GCC	Spice	TeX	US Steel

FIGURE 4.21 Programs used for reporting information about instruction mixes.

There are four types of workloads, and each workload type has a representation program—except that there is no floating-point program for the 8086. The inputs to GCC, Spice, and TeX used for the VAX were purposely shortened because the measurement process is very time intensive. (Readers who obtain measurements for the 360 or 8086 running GCC, Spice, or TeX and who are willing to share their data are asked to contact the publisher.)

In this section we present the instruction-mix measurements using a chart for each machine. The chart shows the average use of an instruction across the programs measured for that architecture. The detailed individual measurements for each program can be found in Appendix C. This appendix will be needed as a reference to do the exercises and examples in the chapter.

Remember that these measurements depend on the benchmarks chosen and the compiler technology used. While the authors feel that the measurements in this section are reasonably indicative of the usage of these four architectures, other programs may behave differently from any of the benchmarks here, and different compilers may yield different results. In doing a real instruction set study, the architect would want to have a much larger set of benchmarks, spanning as wide an application range as possible. He would also want to consider the operating system and its usage of the instruction set. Single-user benchmarks like those measured here do not necessarily behave in the same fashion as the operating system.

VAX Instruction Set Measurements

The data on VAX instruction set usage in this section come primarily from measurements on our three benchmark programs. We add the data reported in another study for COBOL when we discuss opcode distributions. For these measurements, Spice and TeX were compiled with the globally optimizing versions of the VAX compilers originally developed for VMS (called VCC and fort). GCC cannot be compiled by the vcc compiler and hence uses the standard VAX cc compiler, which performs only peephole optimization. Once compiled, these programs were run with the Trace bit turned on. This causes the program to trap on every instruction execution, allowing a measurement program to collect data. Because this slows the program down by a factor of between 1,000 and 10,000 times, smaller inputs were used for the programs GCC, TeX, and Spice.

Addressing Mode Usage

Let's begin by looking at the VAX addressing modes, since the choice addressing modes and operations are orthogonal. First, we break the references into three broad classes: register, immediate (including short literal), and memory addressing modes. Figure 4.22 shows the breakdown into these three classes for our benchmarks. In all three programs, more than half the operand references are to registers.

About one-third of the operands on the VAX are memory references. How are those memory locations specified? The VAX memory addressing modes fall

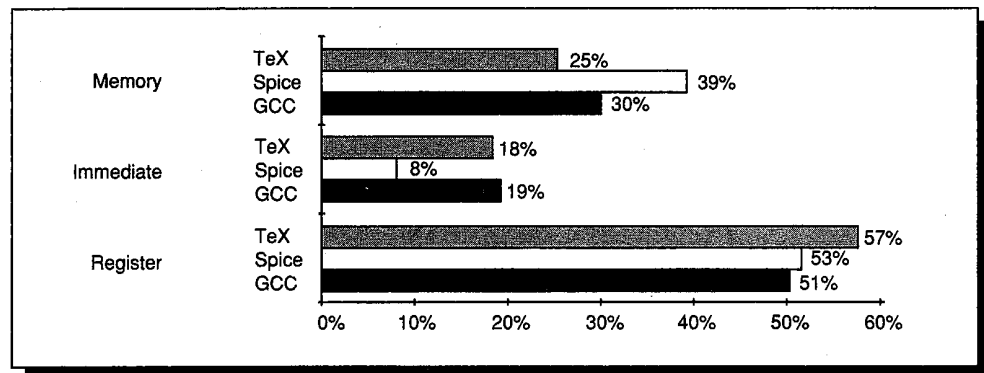


FIGURE 4.22 Breakdown of basic operand types for the three benchmarks on the VAX. The frequencies are very similar across programs, except for the low usage of immediates by Spice and its correspondingly higher use of memory operands. This probably arises because few floating-point constants are stored as immediates, but are instead accessed from memory. An operand is counted by the number of times it appears in an instruction, rather than by the number of references. Thus, the instruction `ADDL2 R1, 45 (R2)` counts as one memory reference and one register reference. The memory address modes in Figure 4.23 are counted in the same fashion. Wiecek [1982] reports that about 90% of the operand accesses are either a read or a write, and only about 10% of the accesses both read and write the same operand (such as R1 in the `ADDL2`).

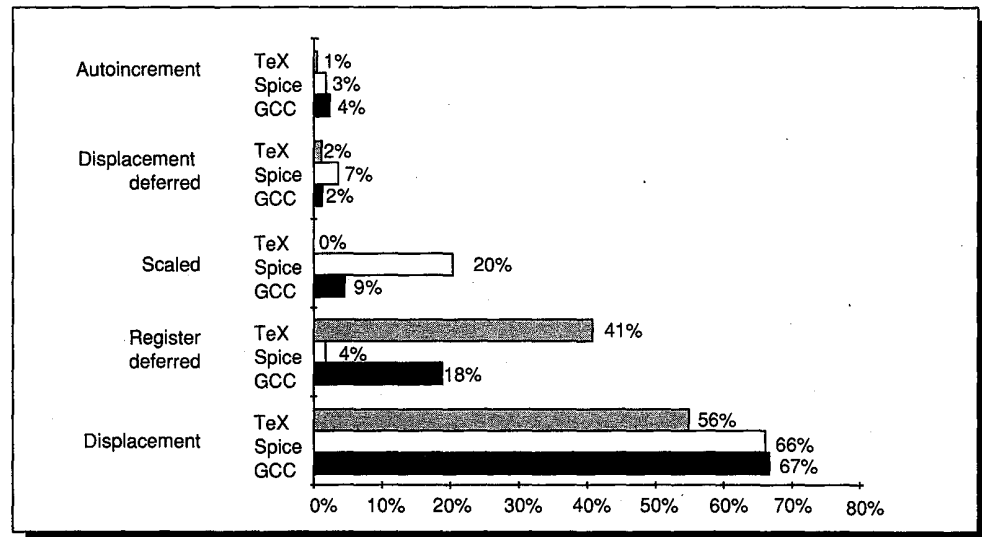


FIGURE 4.23 Use of VAX memory addressing modes, which account for about 31% of the operand references, in the three programs. Spice again stands out because of the low frequency of register deferred. In Spice, nonzero displacement values occur much more frequently. The use of arrays rather than pointers probably influences this. Likewise, Spice uses the scaled mode to access array elements. The displacement deferred mode is used to access actual parameters in a FORTRAN subroutine. Remember that PC-based addressing is not included here—use of PC-based addressing can be measured by branch frequency.

into three separate classes: PC-based addressing, scaled addressing, and the other addressing modes (sometimes called the general addressing modes). The primary use of PC-based addressing is to specify branch targets, rather than data operands; thus, we do not include this addressing mode here. Scaled mode is counted as a separate addressing mode, and the based mode on which it is built is counted as well. Figure 4.23 shows the use of addressing modes in the three benchmark programs. Not surprisingly, displacement mode dominates. Taken together, displacement and register deferred, which is essentially a special case of displacement with a zero constant value, constitute from 70% to 96% of the dynamically occurring addressing modes.

The size of a VAX instruction is almost always one byte for the opcode plus the number of bytes in the addressing modes. From these data the average size of an instruction can be estimated. Architects often do this type of estimating when they do not have exact measurements available. This is particularly true when data collection is expensive. Collecting the VAX data in this chapter, for example, took from one to several days of running time for each program.

Example

The average VAX instruction has 1.8 operands. Use this fact and the data on displacement sizes in Figure 3.13 (on page 100 of Chapter 3) to estimate the average size of a VAX instruction. Such an estimate is useful for determining memory bandwidth per instruction, a critical design parameter.

Answer

From the above data we know that literal and register modes, which each take 1 byte, dominate the mix. The most heavily used addressing mode, displacement mode, can vary from 2 bytes to 5 bytes—the register byte plus 1 or more offset bytes. Based on the length information in Figure 3.13 we guess that the average displacement is 1.5 bytes, for a total size of 2.5 bytes for the addressing mode. For this example, we assume that literal, register, and displacement modes make up all the accesses.

This means there is 1 byte for the opcode, 1 byte for register or literal mode, and about 2.5 bytes for displacement mode. Using 1.8 operands per instruction and the average frequencies of accesses from Figure 4.22 (page 169), we obtain $1 + 1.8 * (0.54 + 0.15 + 0.31 * 2.5)$ or 3.64 bytes.

Wiecek [1982] measured 3.8 bytes per instruction. Direct measurements of our three programs showed the average sizes to be 3.6, 4.9, and 4.2 for GCC, Spice, and TeX, respectively.

Instruction Mixes

Now let's look at the distribution for instruction operations, using our three benchmarks plus the COBOLX program from the study published by Clark and Levy [1982]. COBOLX is a synthetic, internal DEC benchmark that was compiled by the VAX VMS COBOL compiler and uses decimal instructions. However, the new DEC compilers for the VAX avoid using the decimal instruction set, since most of that portion of the architecture is emulated in software—and is therefore much slower—on the newer VLSI-based VAXes.

The data in this section are presented in chart form, but detailed tables for each machine and benchmark appear in Appendix C. The data here focus on instruction frequency, but frequency distributions and time distributions do not always match. We will see an example of this in the next section. Appendix D contains a set of detailed measurements based on time-distribution measurements.

Figure 4.24 shows all instructions responsible for more than 1.5% of the dynamic instruction executions across all the benchmarks. Each complete bar shows an average instruction mix over the four programs, and how the programs make up that mix.

GCC and TeX are very similar in behavior; the largest difference is the higher frequency of data transfers in TeX. Spice and COBOLX look very different. Each of these executes more than 20% of its instructions using a portion of the instruction set that is essentially unused by the other benchmarks. Both COBOLX and Spice do many fewer integer arithmetic operations, instead using decimal or floating-point operations. COBOLX makes small use of the data transfer instructions (4% versus an average of 20% for the other three programs); instead, 38% of the instructions it executes are decimal or string instructions.

These 27 instructions in Figure 4.24 correspond to an average of 88% of the instructions executed in the four benchmarks. However, the tail of the

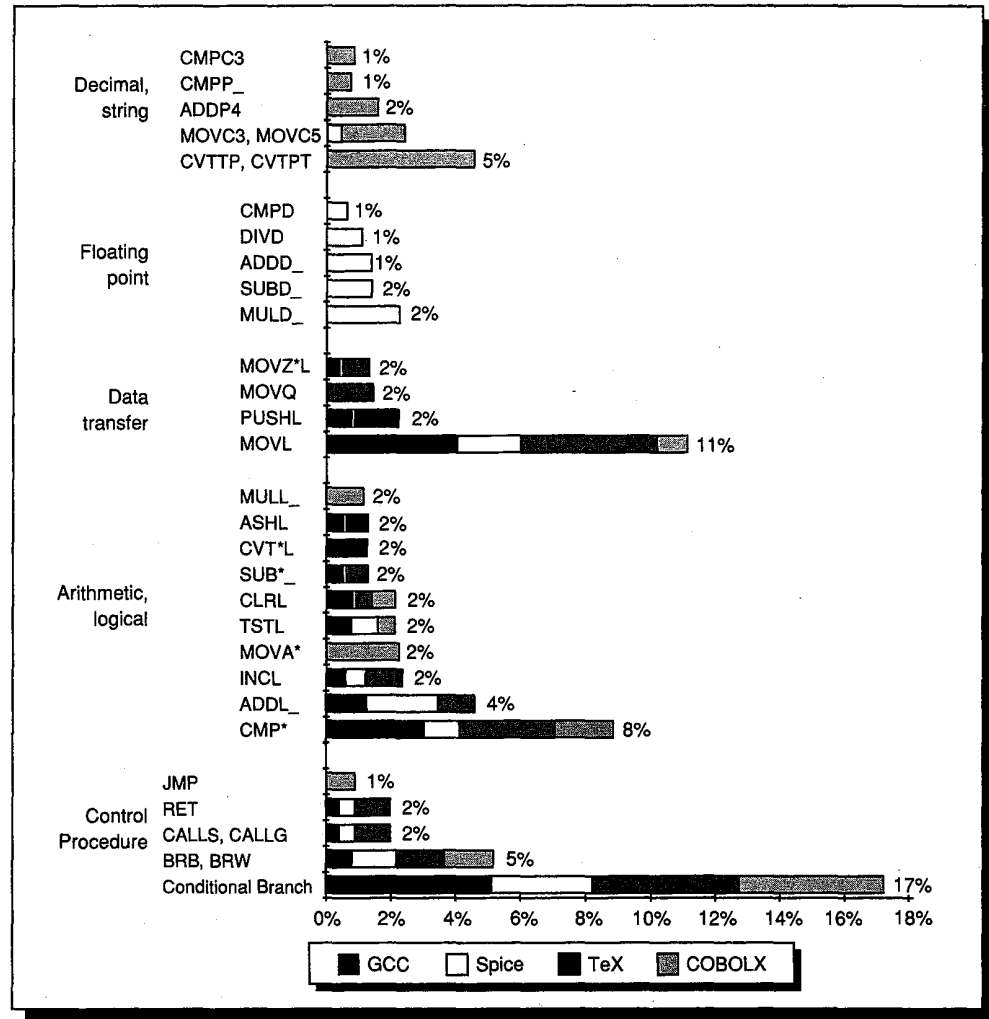


FIGURE 4.24 The VAX instruction frequencies combined graphically. The total size of each bar shows the behavior that would be seen on a machine that ran these four programs with equal frequency. The segments of the bar show what percentage of the usage of that instruction would come from each of the programs. This illustrates that some portions of the instruction set need to be there for only one class of applications. Overall, only a small number of instructions outside of the control, data transfer, and integer arithmetic instructions are heavily used.

distribution is long and there are many instructions executed with a frequency of 1/2 to 1%. In Spice, for example, the top 15 instructions make up 90% of the executions, and the top 26 make up 95%. However, there are 149 different VAX instructions executed at least once!

Measurements of 360 Instruction Set Usage

The measurements in this section are taken from those made by Shustek in his Ph.D. thesis [1978]. His work includes a study of the dynamic characteristics of

seven large programs on the IBM 360 architecture. He collected his data by building an interpreter for the 360 architecture. The four programs described in Figure 4.25 are used in this section to examine characteristics of 360 instruction set usage.

Program	Benchmark class	Instruction count	Program function
COBOLGO	Business D.P.	3,559,533	COBOL usage report formatter
PLIGO	General integer	23,863,497	PL/I computer usage accounting
FORTGO	Floating point	11,719,853	FORTRAN linear systems solver
PLIC	Compiler	24,338,101	PL/I compile

FIGURE 4.25 Four programs used to measure the IBM 360. The suffix "GO" indicates an execution of a program, while the suffix "C" indicates a compile. We chose the PL/I compiler because it is the largest and most representative; it is also written in PL/I. Shustek's thesis used two FORTRAN executions. We chose to use LINSYS2 to represent the FORTRAN execution, since it is a more typical FORTRAN program; we refer to the execution of LINSYS2 as FORTGO.

Addressing Modes and Instruction Types

Figure 4.26 shows the frequency of data accesses by addressing mode. The COBOL program has a very high frequency of data accesses. Movements of character data and use of decimal data, which always reside in memory, probably account for this. FORTGO has a substantially lower number of memory references. This may arise because of allocation of variables to registers in the tight inner loops of the program.

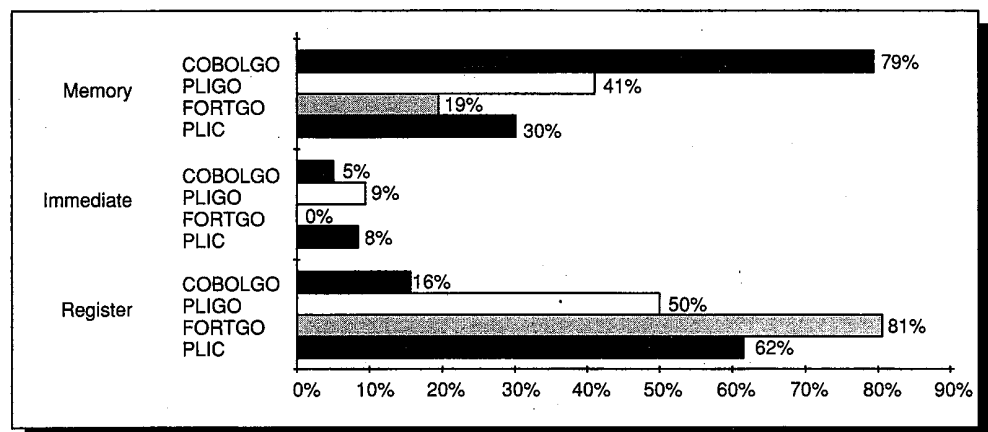


FIGURE 4.26 Distribution of operand accesses made by 360 instructions. Limited support for immediates is the chief reason that immediates see so little use.

There are only two memory addressing modes on the 360: base register + displacement (RS format, SI format, and SS format) and base register + displacement + index register (RX format). However, the operations available in the instructions that address memory typically appear in only one format. Therefore, it is probably most useful to look at instruction format usage, as shown in Figure 4.27. Most instructions are RX format with RR following behind that. The high usage of RX format should not lead you to conclude that the displacement + base register + index register addressing mode is heavily used, because in 85% of the RX instructions the index register is zero. COBOL displays a high percentage of SS-format instructions, and this is to be expected because the decimal and string instructions are all SS format. The FORTRAN execution displays a large percentage of RR format, 2-byte instructions. This makes sense in a program that makes heavy use of registers in its optimized inner loops.

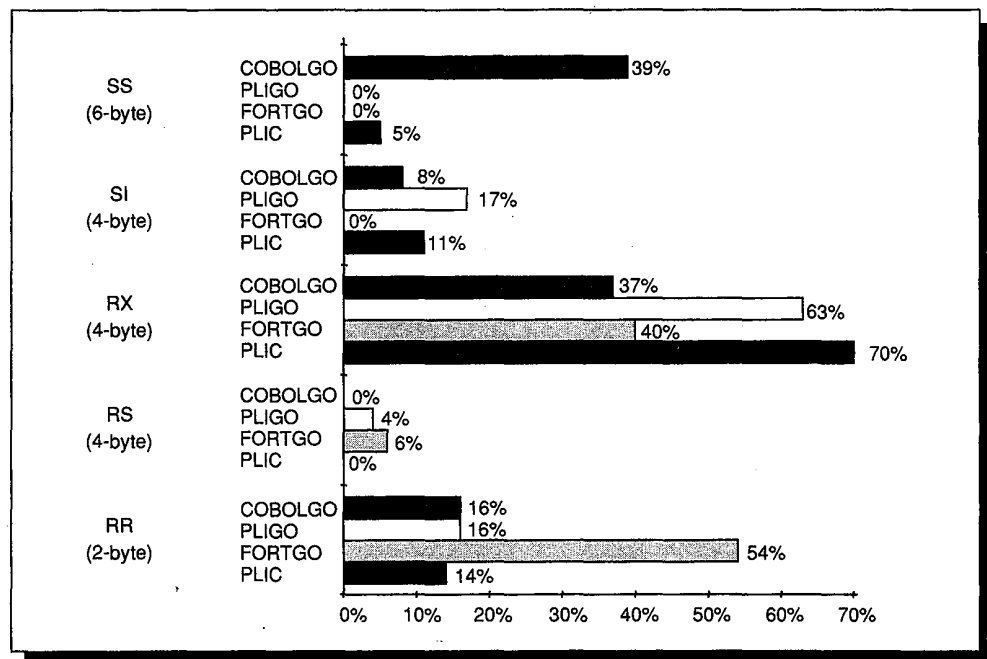


FIGURE 4.27 Percentage of 6-, 4-, and 2-byte instructions for the four 360 programs. The majority of the instructions are 4 bytes, and almost none are 6 bytes, except when running COBOLGO.

Example

Given the data in Figure 4.27 compute the average instruction length for the PLIGO program.

Answer

The average instruction length is

$$\begin{aligned}
 & 6 * \% \text{ SS} + 4 * (\% \text{ RX} + \% \text{ RS} + \% \text{ SI}) + 2 * \% \text{ RR} \\
 & = 0 + 4 * (0.63 + 0.04 + 0.17) + 2 * 0.16 = 3.68 \text{ bytes}
 \end{aligned}$$

Across all the four programs the average measured length is 3.7 bytes.

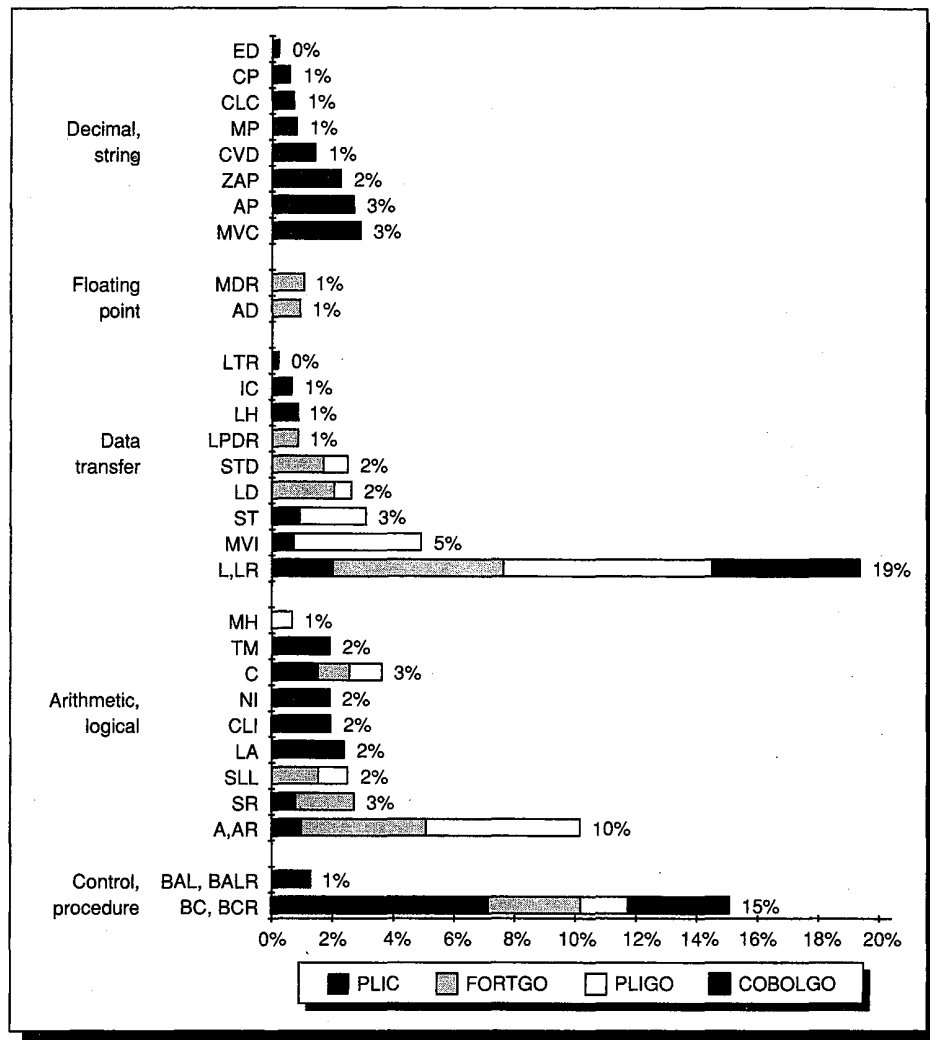


FIGURE 4.28 Combined data for the four programs on the 360. Compare this with Figure 4.24, where the data for the VAX are graphed.

Instruction Mixes

Now let's examine the data for the instruction mixes. Figure 4.28 shows the most heavily used instructions in the four 360 benchmarks. As Figure 4.28 illustrates, variations among the programs are very large. The PL/I compiler has an extraordinarily large number of branches, while the PL/I execution has very few. The use of arithmetic and logical operators is fairly uniform with the exception of the COBOL program, which uses decimal operations instead.

Comparing these programs to the VAX, the much lower frequency of branches—16% on the 360 versus 23% on the VAX—stands out. The number of branches in a program is largely fixed by the program, except for some architectural anomalies and possible compiler optimizations (such as loop unrolling—discussed in Chapter 6—but not used by these compilers). Thus, the

percentage of branches is an indirect measure of instruction power or density, since it says how many other instructions are required for each branch. We would expect the VAX with its more powerful addressing modes and multiple memory operands per instruction to have a high instruction density and a higher branch frequency. We see further evidence of greater instruction density of the VAX in the higher frequency of data transfers on the 360—more data is moved explicitly on the 360 rather than used as memory operands, as on the VAX. However, we cannot draw any specific quantitative conclusions about instruction density because the measured programs and compilers are different.

Also very different is the percentage of character and string operations used by the 360 versus the VAX for the two COBOL applications. Finally, the FORTRAN execution uses a much larger number of integer operations on the 360; this may be traceable to differences arising when the VAX uses an addressing mode but the 360 must use explicit instructions for address calculations.

As we have seen, the differences in instruction usage on the 360 and VAX are fairly significant. The next two architectures differ from these first two even more dramatically.

Measurements of 8086 Usage

The data in this section were collected by Adams and Zimmerman [1989] in a study of seven programs running on an IBM PC under MS DOS 3.1. They collected the data by single-stepping the programs and collecting data after every instruction execution, just as was done for the VAX. The three programs used here, a brief description, and the number of instructions executed are shown in Figure 4.29. As with the VAX and 360, we will begin by examining operand access and addressing modes, and then progress to instruction mixes.

Addressing Modes and Instruction Length

Our first measurement on the 8086, shown in Figure 4.30, graphs the origins of operands. immediates play a small role, while register access slightly dominates memory access. Compared to the VAX, these programs on the 8086 use a higher frequency of memory operands. The limited register set of the 8086 probably

Program	Benchmark class	Instruction count	Program function
Lotus	Business	2,904,931	Lotus 1-2-3 calculating a 128-cell worksheet four times
MASM	General integer	2,365,711	Microsoft Macro Assembler assembling a 500-line program
Turbo C	Compiler	1,806,143	Turbo C compiling Dhrystone

FIGURE 4.29 Three programs used for 8086 measurements. The benchmarks are written in a combination of 8086 Assembler and in C.

plays a role in increasing the memory traffic, which substantially exceeds that of the 360, if we ignore the COBOL program (which must use SS instructions) on the 360.

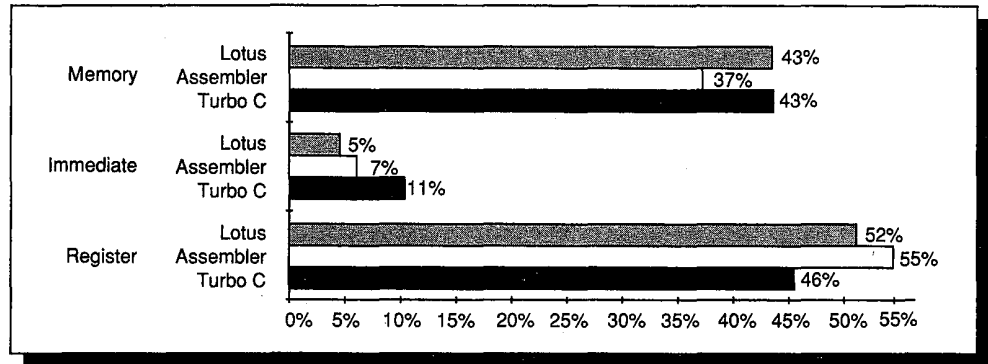


FIGURE 4.30 Three classes of basic operand access on the 8086 and their distribution. The implied use of the accumulator register (AX), which occurs in a number of instructions, is counted as a register access.

In the above programs 41% of the operand references are memory accesses. Figure 4.31 shows the distribution of addressing modes for these memory references.

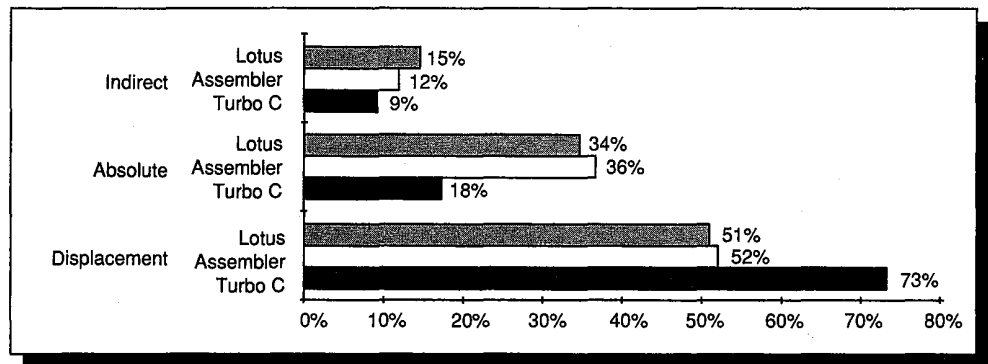


FIGURE 4.31 The 8086 memory addressing modes shown in this graph account for almost all the memory references in the three programs. Memory addressing modes indexed and based have been combined, since their effective address calculations are the same. Register indirect mode is in effect based with a zero offset, equivalent to the VAX register-deferred mode. If register indirect were counted as a based mode with zero offset, about two-thirds of the memory references would be displacement mode. The other two remaining modes are essentially unused in the three programs.

The variable-length instructions, use of implicit registers, and small size of the register specifier combine to yield a fairly short average instruction. For these three programs the average instruction length is approximately 2.5 bytes.

Instruction Mixes on the 8086

The instructions responsible for greater than 1.5% of the executions for the 8086 running the three programs are shown graphically in Figure 4.32. The displayed subset of the instruction set accounts for a higher proportion of all instruction executions (90%) than it does on the VAX or 360. As we might suspect, the architectures with smaller instruction repertoires use a higher percentage of their opcodes.

The major distinguishing characteristic among the programs is the shift from data transfer instructions to control instructions in Lotus. Lotus makes heavy use of the LOOP instruction, which may account for that shift.

The overall frequency of move instructions is much larger on the 8086 than on the VAX. This difference probably arises because the 8086 has fewer general-purpose registers. Other possible explanations include the use of string instructions that generate a sequence of move instructions, and explicit movement of data among segments to ease processing. The total branch frequency is not very different between the 8086 and VAX, though the distribution of different types of control instructions is very different. The percentage of arithmetic operations on the 8086 is much smaller, due at least partially to the larger number of move instructions.

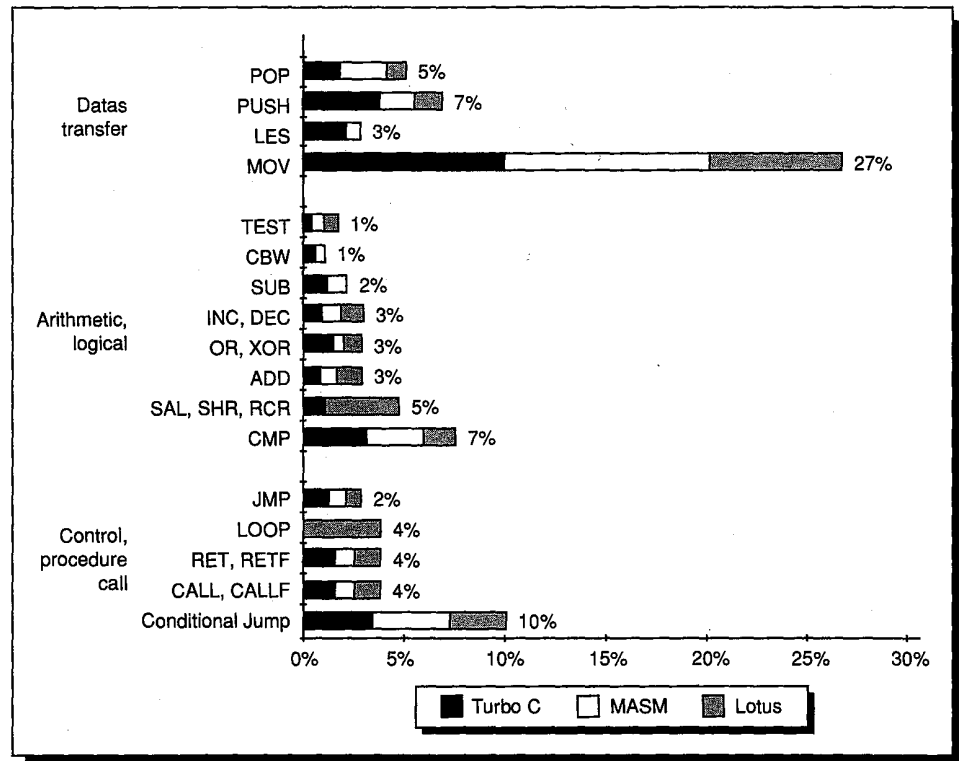


FIGURE 4.32 Distribution of instruction frequencies on the 8086 shown in the same format used for the VAX and 360.

In this and the preceding two sections we saw machines designed in the 1960s (the 360) and the 1970s (the VAX and the 8086). In the next section we will talk about a machine typical of those designed in the 1980s and its usage.

Instruction Set Usage Measurements on DLX

As with the other architectures we have looked at thus far, we start our examination of instruction set usage on DLX with measurements of operand location and move from there to instruction mixes. The DLX data throughout the book was measured using the MIPS R2000/3000 architecture and adjusting the data to reflect the differences between DLX and the MIPS architecture. The MIPS compiler technology with optimization level 2, which does full global optimization with register allocation, was used to compile the programs. A special program called *pixie* was used to instrument the object module. The instrumented object module produces a monitoring file that is used to produce detailed execution statistics.

Addressing Mode Usage

Operand usage is shown in Figure 4.33. This data is very uniform across the applications on DLX. Compared to the VAX, a much higher percentage of the operand references are to registers: On the VAX, only about half the references are to registers, while roughly three-quarters are on DLX. This probably occurs because of the larger number of registers available on DLX and greater emphasis on register allocation by the DLX compiler.

Since DLX has only a single addressing mode, it makes no sense to ask what the distribution of addressing modes is. However, we noticed earlier that on the VAX and 8086 the deferred addressing mode, which is equivalent to displacement addressing with a zero displacement, was the second or third most popular. Would it be useful to add this mode to DLX?

Example

Using the data on offset values from Figure 3.13 on page 100, determine how often on average deferred mode would be used for the three programs if the case of a zero-offset displacement were made a special mode. In particular, what percentage of the memory references would use it? How much memory bandwidth would be saved if we had a 16-bit instruction for this addressing mode?

Answer

The frequencies of zero-offset displacement values are

GCC: 27%

Spice: 4%

TeX: 17%

The average frequency for a zero-offset value is then $(27\% + 4\% + 17\%)/3 = 16\%$. Thus, the mode would be used by 16% of the loads and stores, which average 32% of the executions. The decrease in instruction bandwidth would be about $\frac{1}{2} * 32\% * 16\%$, or about 3%.

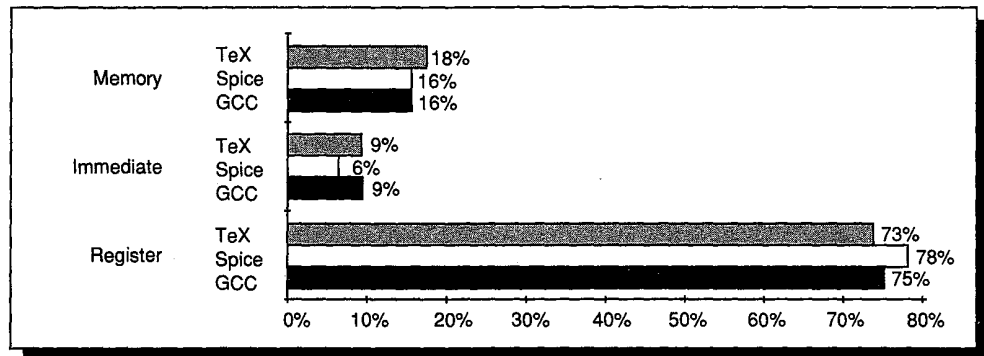


FIGURE 4.33 Distribution of operand accesses for the three benchmarks on DLX.

Only accesses for operands—not for effective address calculations—are included. The fact that DLX has only 3-operand register formats probably increases the frequency of register operand access slightly, since some instructions probably have only two unique register operands and use one register as both a source and destination. On a machine like the VAX, such an operation might use a 2-operand instruction and thus be counted as having only 2 register operands. This effect has not been measured.

The other two addressing modes used with some frequency are scaled on the VAX and absolute on the 8086. Scaled addressing mode is synthesized on DLX with a separate add; the presence of this address mode is significantly affected by the compiler technology. Better optimizers use the indexed mode less often because the optimization of induction variable elimination obviates the need for indexed addressing and for scaling (see the discussion in Section 3.7). The direct mode is synthesized by dedicating a register to point to a global area and accessing variables with a displacement from that register. Because only scalar variables (i.e., not structures or arrays) need to be accessed in this way, this works very well for most programs.

Instruction Mixes on DLX

Figure 4.34 shows the instruction mixes for our three programs plus the U.S. Steel COBOL benchmark—the most widely employed COBOL benchmark. The benchmark is a synthetic program of about 1,000 lines in length. It is included here because its behavior is substantially different from FORTRAN and C programs. Measurements on COBOL are also interesting because they reflect what changes in instruction set usage occur when decimal arithmetic is not

directly supported by decimal instructions. Let's first look at the differences among the programs before we consider how these mixes compare to the VAX.

The significant differences among these programs are surprising. Both Spice and TeX stand out as having very low branch frequency. The effect of translating the decimal arithmetic of COBOL into binary arithmetic is clearly seen in the large percentage of arithmetic operations in US Steel. Shifts, logical operators, and load immediates, which are all used to do fast decimal-binary conversion, occur in significant frequencies. US Steel's low frequency of data transfer is certainly affected by this increase in arithmetic and logical operations. Interestingly, only the call frequency of US Steel is high enough to account for more than 1% of the instruction executions (the frequency of JAL is about 1% for the other three benchmarks).

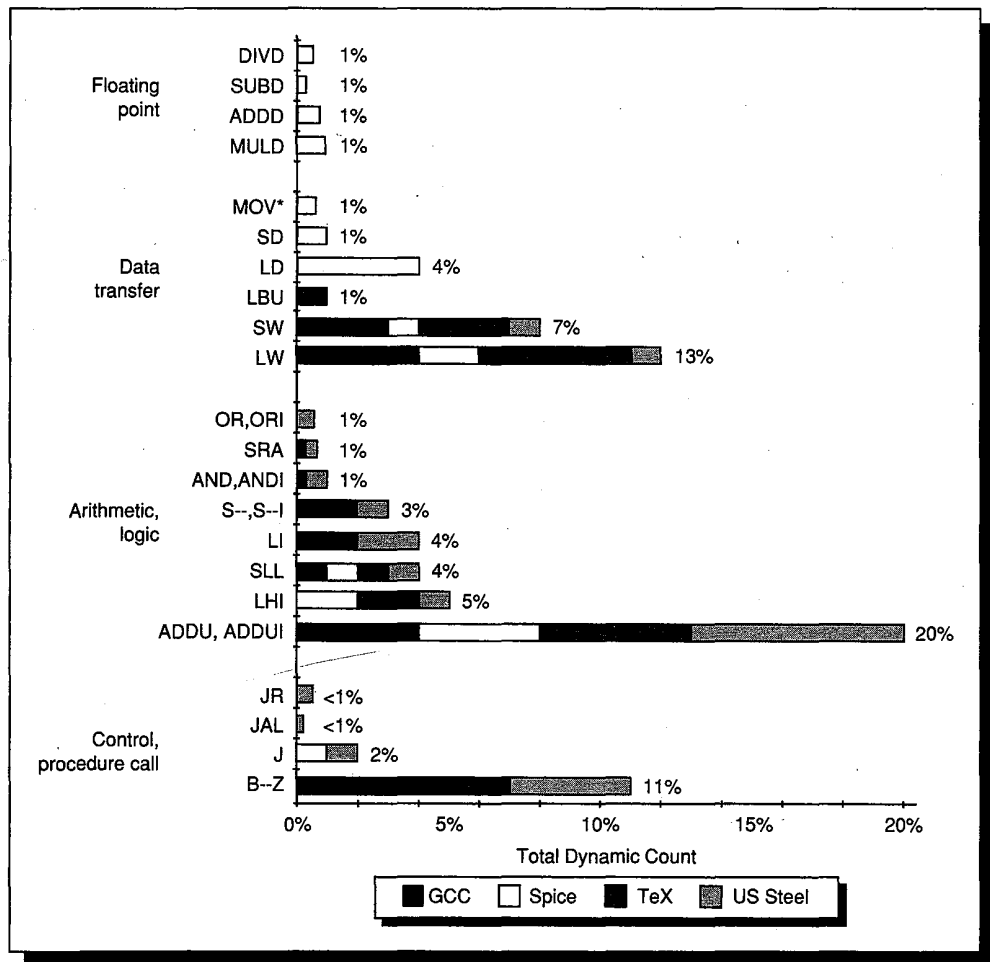


FIGURE 4.34 The DLX instruction mix visible over four programs with breakdown showing each program's contribution. What is remarkable is how a small number of instructions—conditional branch, add, load, and store—dominate across all four programs. The opcode LI is really an ADDUI with R0 as an operand; the high frequency of ADDU and ADDUI is discussed below.

These mixes differ dramatically from the VAX (or other machines in this section). One difference is the very high percentage of ADDU and ADDUI instructions. These instructions are used for a variety of purposes where the other machines may use a different instruction or a more powerful addressing mode or instruction. Among the most frequent uses for ADDU and ADDUI are: register-register copies (coded as ADDU with R0), synthesizing an address mode such as scaled, and incrementing the stack pointer on a procedure call.

It is interesting to compare the branch frequency between DLX and the VAX, since the absolute branch count should be approximately equal (for reasons discussed earlier), and the ratio of branch frequencies should be about the same as the ratio of overall instruction counts. However, the compilers may affect the type of branch used—conditional branch versus jump—so we need to combine all the branches and jumps, except those used in procedure calls, to make a comparison.

Example

Find the ratio of absolute branches on the VAX versus DLX for the three common benchmarks. The ratio of instruction counts, measured in Section 3.8 is

$$\frac{\text{Instructions}_{\text{DLX}}}{\text{Instructions}_{\text{VAX}}} = 2.0$$

Use the data in Appendix C for exact percentages of branches.

Answer

From Appendix C, we find that the average branch frequency on DLX is $\frac{19\%+2\%+7\%}{3} = 9.3\%$, while the average for the VAX is 17.3%. Thus, the ratio of the branch counts is

$$\begin{aligned} \frac{\text{Branches}_{\text{DLX}}}{\text{Branches}_{\text{VAX}}} &= \frac{9.3\% * \text{Instructions}_{\text{DLX}}}{17.3\% * \text{Instructions}_{\text{VAX}}} \\ &= \frac{9.3 * 2.0 * \text{Instructions}_{\text{VAX}}}{17.3 * \text{Instructions}_{\text{VAX}}} \\ &= \frac{18.6}{17.3} = 1.08 \end{aligned}$$

So DLX does about 8% more branches.

In the arithmetic and logical instructions, GCC and US Steel are the most different between the VAX and DLX. We know US Steel differs because of the absence of decimal instructions—it would be interesting to see what the instruction mix on such a program would look like with the new VAX compilers that avoid the decimal instructions. Another major difference between the two machines is the lower frequency of compare instructions and test instructions on DLX. The use of compare with zero in the branch instruction is responsible for

this. Because the set instructions are also used to set logical variables, we cannot know exactly what percentage of conditional branches on DLX do not need a compare, but we can guess that it is between 75% and 80%.

The difference in data transfers has been discussed extensively at the end of Chapter 3 (for a machine very close to DLX) and in the previous subsection. We know that the larger number of registers (at least twice as many) and more ambitious register allocator mean that the load and store frequency is lower on DLX than on the VAX.

We have now seen instruction mixes for four very different machines. In Appendix D we can see how these mixes differ when we look at time distributions rather than frequency of occurrence, and in the next section we will review some of our key observations and point out some additional pitfalls using data we have examined in this and earlier sections.

4.7 Fallacies and Pitfalls

Fallacy: There is such a thing as a typical program.

Many people would like to believe that there is a single “typical” program that could be used to design an optimal instruction set. For example, see the synthetic benchmarks discussed in Section 2.2. The data in this chapter clearly show that programs can vary significantly in how they use an instruction set. For example, the frequency of control-flow instructions on DLX varied from 5% to 23%. The variations are even larger on an instruction set that has specific features for supporting a class of applications, such as decimal or floating-point instructions that are unused by other applications. There is a related pitfall.

Pitfall: Designing an architecture on the basis of small benchmarks or large benchmarks from a restricted application domain when the machine is intended to be general purpose.

Many programs exhibit somewhat biased behavior or do not use a particular aspect of an architecture. Obviously, choosing TeX or GCC benchmarks to design the instruction set might result in a machine that wouldn’t do well on a program like Spice or COBOLX. A more subtle example arises when choosing a representative, but synthetic, benchmark. For example, Dhrystone (see Section 2.2) does a procedure call approximately every 40 instructions on a machine like DLX—the number of procedure calls is more than half the number of conditional branches! By comparison, in GCC a call occurs about once every 100 instructions, and branches are 15 times more frequent than procedure calls.

Fallacy: An architecture with flaws cannot be successful.

The IBM 360 is often criticized in the literature—the branches are not PC-relative, and the offset is too small in based addressing. Yet, the machine has

been an enormous success because it did several new things properly. First, the architecture has a big enough address space. Second, it is byte addressed and handles bytes well. Third, it is a general-purpose register machine. Finally, it is simple enough that it can be efficiently implemented across a wide performance and cost range.

The 8086 provides an even more dramatic example. The 8086 architecture is the only widespread architecture in existence today that is not truly a general-purpose register machine. Furthermore, the segmented address space of the 8086 causes major problems both for programmers and compiler writers. Despite these major difficulties, the 8086 architecture—because of its selection as the microprocessor in the IBM PC—has been enormously successful.

Fallacy: One can design a flawless architecture.

All architecture design involves tradeoffs made in the context of a set of hardware and software technologies. Over time those technologies are likely to change, and decisions that may have been correct at the time they were made look like mistakes. For example, in 1975 the VAX designers overemphasized the importance of code-size efficiency and underestimated how important ease of decoding and pipelining would be ten years later. Almost all architectures eventually succumb to the lack of sufficient address space. However, avoiding this problem in the long run would probably mean compromising the efficiency of the architecture in the short run.

Fallacy: In instruction mixes, time distribution and frequency distribution will be close.

Appendix D shows the time distributions for our benchmark programs and compares the time and frequency distributions. A simple example of where these distributions are very different is in the COBOLGO program on the 360. Figure 4.35 shows the top instructions by frequency and by time. The two highest occurring instructions are responsible for 33% of the instruction executions in COBOLGO, but only 4% of the execution time! Remember that time distributions are dependent on both the architecture and the **implementation** used for the measurement. Hence, time distributions may differ from model to model, while frequency distributions will be the same, provided neither the software nor the program changes. This large difference between time and frequency distributions does not exist for simpler load/store architectures, such as DLX.

Pitfall: Examining only the worst-case or average behavior of an instruction as design input.

The best example of this comes from the use of MVC on an IBM 360. The instruction can move overlapped fields of characters, but this occurs less than 1% of the time, and then usually to clear a field. The average length of a move

as measured by Shustek was ten bytes, but more than three-quarters of the moves were either one byte or four bytes in length. Assuming worst-case behavior (overlapping strings) or average length can each lead to suboptimal design decisions.

Top instructions by frequency	Frequency	Top instructions by time distribution	Percentage of time
L, LR	19%	ZAP	16%
BC, BCR	14%	AP	16%
AP	11%	MP	13%
ZAP	9%	MVC	9%
MVC	7%	CVD	5%

FIGURE 4.35 The top five instructions by frequency and by time for the COBOLGO benchmark run on the 360. The actual frequency or percentage of time is also shown. Further data appears in Appendix D.

4.8 Concluding Remarks

We have seen that instruction sets can vary quite dramatically, both in how they access operands and in the operations that can be performed by a single instruction. The comparison of opcode usage across architectures by instruction frequency is summarized in Figure 4.36. This figure shows that even very different architectures behave similarly in their use of instruction classes. However, this should also remind us that performance may be only distantly related to instruction usage—the execution-time distributions for these architectures in Appendix D look very different indeed.

Dramatic though the variation in instruction usage is across architectures, it is equally dramatic across applications. We have seen that floating-point programs, COBOL programs, and C systems programs differ in how they use a machine. Large segments of the instruction set are unused by some programs. When such application-specific features are not part of the instruction set—for example, the absence of decimal instructions in DLX—the impact is a shift in the use of other parts of the instruction. Even across two programs written in the same language—GCC and TeX, or PLIC and PLIGO—the differences in instruction usage can be significant.

Instruction-usage data are an important input for the architect, but they do not necessarily tell us what are the most time-consuming instructions. The next several chapters will help explain why the difference arises by quantifying the CPI difference among instructions and machines.

Machine	Program	Control	Arithmetic, logical	Data transfer	Floating point	Decimal, string	Totals
VAX	GCC	30%	40%	19%			89%
VAX	Spice	18%	23%	15%	23%		79%
VAX	TeX	30%	33%	28%			91%
VAX	COBOLX	25%	24%	4%		38%	91%
360	PLIC	32%	29%	17%		4%	82%
360	FORTGO	13%	35%	40%	7%		95%
360	PLIGO	5%	29%	56%			90%
360	COBOLGO	16%	9%	20%		40%	85%
8086	Turbo C	21%	23%	49%			93%
8086	MASM	20%	24%	46%			90%
8086	Lotus	32%	26%	30%			88%
DLX	GCC	24%	35%	27%			86%
DLX	Spice	4%	29%	35%	15%		83%
DLX	TeX	10%	41%	33%			84%
DLX	US Steel	23%	49%	10%			82%

FIGURE 4.36 The frequency of instruction distribution for each benchmark broken into five classes of instructions. Because only instructions with frequencies greater than 1.5% have been included in previous figures, the totals are less than 100%.

4.9 Historical Perspective and References

Although a large number of machines have been developed in the same time frames as the four machines covered in this chapter, the discussion here is confined to these machines and measurements of them.

The IBM 360 was introduced in 1964 with six models and a 25:1 performance ratio. Amdahl, Blaauw, and Brooks [1964] discuss the architecture of the IBM 360 and the concept of permitting multiple object-code-compatible implementations. The notion of an instruction set architecture as we understand it today was the most important aspect of the 360. The architecture also introduced several important innovations, now in wide use:

1. 32-bit architecture
2. Byte-addressable memory with 8-bit bytes
3. 8-, 16-, 32-, and 64-bit data sizes

In 1971, IBM shipped the first System/370 (models 155 and 165), which included a number of significant extensions of the 360, as discussed by Case and Padege [1978], who also discuss the early history of System/360. The most important addition was virtual memory, though virtual memory 370s did not ship until 1972 when a virtual-memory operating system was ready. By 1978,

the high-end 370 was several hundred times faster than the low-end 360s shipped ten years earlier. In 1984, the 24-bit addressing model built into the IBM 360 needed to be abandoned, and the 370-XA (eXtended Architecture) was introduced. While old 24-bit programs could be supported without change, several instructions could not function in the same manner when extended to a 32-bit addressing model (31-bit addresses supported) because they would not produce 31-bit addresses. Converting the operating system, which was written mostly in assembly language, was no doubt the biggest task.

Several studies of the IBM 360 and instruction measurement have been made. Shustek's thesis [1978] is the best known and most complete study of the 360/370 architecture. He made several observations about instruction set complexity that were not fully appreciated until some years later. Another important study of the 360 is the Toronto study by Alexander and Wortman [1975] done on an IBM 360 using 19 XPL programs.

In the mid-1970s, DEC realized that the PDP-11 was running out of address space. The 16-bit space had been extended in several creative ways. However, as Strecker and Bell [1976] observed, the small address space was a problem that could not be overcome, but only postponed.

In 1978, DEC introduced the VAX. Strecker [1978] described the architecture and called the VAX "a Virtual Address eXtension of the PDP-11." One of DEC's primary goals was to keep the installed base of PDP-11 customers. Thus, the customers were to think of the VAX as a 32-bit successor to the PDP-11. A 32-bit PDP-11 was possible—there were three designs—but Strecker reports that they were "overly compromised in terms of efficiency, functionality, programming ease." The chosen solution was to design a new architecture and include a PDP-11 compatibility mode that would run PDP-11 programs without change. This mode also allowed PDP-11 compilers to run and to continue to be used. The VAX-11/780 was made similar to the PDP-11 in many ways. These are among the most important:

1. Data types and formats are mostly equivalent to those on the PDP-11. The F and D floating formats came from the PDP-11. G and H formats were added later. The use of the term "word" to describe a 16-bit quantity was carried from the PDP-11 to the VAX.
2. The assembly language was made similar to the PDP-11's.
3. The same buses were supported (Unibus and Massbus).
4. The operating system, VMS, was "an evolution" of the RSX-11M/IAS OS (as opposed to the DECsystem 10/20 OS, which was a more advanced system).
5. The file system was basically the same.

The VAX-11/780 was the first machine announced in the VAX series. It is one of the most successful and heavily studied machines ever built. The cornerstone of DEC's strategy was a single architecture, VAX, running a single

operating system, VMS. This strategy worked well for over ten years. The large number of papers reporting instruction mixes, implementation measurements, and analysis of the VAX make it an ideal case study.

Wiecek [1982] reported on the use of various architectural features in running a workload consisting of six compilers. Emer did a set of measurements (reported by Clark and Levy [1982]) on the instruction set utilization of the VAX when running four very different programs and when running the operating system. A good detailed description of the architecture, including memory management and an examination of several of the VAX implementations, can be found in Levy and Eckhouse [1989].

The first microprocessors were produced late in the first half of the 1970s. The Intel 4004 and 8008 were extremely simple 4-bit and 8-bit accumulator-style machines. Morse et al. [1980] describe the evolution of the 8086 from the 8080 in the late 1970s in an attempt to provide a 16-bit machine with better throughput. At that time almost all programming for microprocessors was done in assembly language—both memory and compilers were in short supply. Intel wanted to keep its base of 8080 users, so the 8086 was designed to be “compatible” with the 8080. The 8086 was **never** object-code compatible with the 8080, but the machines were close enough that translation of assembly language programs could be done automatically.

In early 1980, IBM selected a version of the 8086 with an 8-bit external bus, called the 8088, for use in the IBM PC. (They chose the 8-bit version to reduce the cost of the machine.) This choice, together with the tremendous success of the IBM PC and its clones (made possible because IBM opened the architecture of the PC), has made the 8086 architecture ubiquitous. While the 68000 was chosen for the popular Macintosh, the Macintosh was never as pervasive as the PC (partly because Apple did not allow clones), and the 68000 did not acquire the same software leverage that the 8086 enjoys. The Motorola 68000 may have been more significant **technically** than the 8086, but the impact of the selection by IBM and IBM’s open architecture strategy dominated the technical advantages of the 68000 in the market. As discussed in Section 4.4, the 80186, 80286, 80386, and 80486 have extended the architecture and provided a series of performance enhancements.

There are numerous descriptions of the 80x86 architecture that have been published—Wakerly’s [1989] is both concise and easy to understand. Crawford and Gelsinger [1988] is a thorough description of the 80386. The work of Adams and Zimmerman [1989] represents the first detailed, published study of the dynamic use of the architecture that we are aware of; the data on the 8086 used in this book come from their study.

The simple load/store machines from which DLX is derived are commonly called RISC (*reduced instruction set computer*) architectures. The roots of RISC architectures go back to machines like the 6600, where Thornton, Cray, and others recognized the importance of instruction set simplicity in building a fast machine. Cray continued his tradition of keeping machines simple in the CRAY-1. However, DLX and its close relatives are built primarily on the work of three

research projects: the Berkeley RISC processor, the IBM 801, and the Stanford MIPS processor. These architectures have attracted enormous industrial interest because of claims of a performance advantage of anywhere from two to five times over other machines using the same technology.

Begun in the late 1970s, the IBM project was the first to start but was the last to become public. The IBM machine was designed as an ECL minicomputer, while the university projects were both MOS-based microprocessors. John Cocke is considered to be the father of the 801 design. He received both the Eckert-Mauchly and Turing awards in recognition of his contribution. Radin [1982] describes the highlights of the 801 architecture. The 801 was an experimental project, but was never designed to be a product. In fact, to keep down cost and complexity, the machine was built with only 24-bit registers.

In 1980, Patterson and his colleagues at Berkeley began the project that was to give this architectural approach its name (see Patterson and Ditzel [1980]). They built two machines called RISC-I and RISC-II. Because the IBM project was not widely known or discussed, the role played by the Berkeley group in promoting the RISC approach was critical to the acceptance of the technology. In addition to a simple load/store architecture, this machine introduced register windows—an idea that has been adopted by several commercial RISC machines (this concept is discussed further in Chapter 8). The Berkeley group went on to build RISC machines targeted toward Smalltalk, described by Ungar et al. [1984], and LISP, described by Taylor et al. [1987].

In 1981, Hennessy and his colleagues at Stanford published a description of the Stanford MIPS machine. Efficient pipelining and compiler-assisted scheduling of the pipeline were both key aspects of the original MIPS design.

These three early RISC machines had much in common. Both the university projects were interested in designing a simple machine that could be built in VLSI within the university environment. All three machines—the 801, MIPS, and RISC-II—used a simple load/store architecture, fixed-format 32-bit instructions, and emphasized efficient pipelining. Patterson [1985] describes the three machines and the basic design principles that have come to characterize what a RISC machine is. Hennessy [1984] is another view of the same ideas, as well as other issues in VLSI processor design.

In 1985, Hennessy published an explanation of the RISC performance advantage and traced its roots to a substantially lower CPI—under two for a RISC machine and over ten for a VAX-11/780 (though not with identical workloads). A paper by Emer and Clark [1984] characterizing VAX-11/780 performance was instrumental in helping the RISC researchers understand the source of the performance advantage seen by their machines.

Since the university projects finished up, in the 1983-84 timeframe, the technology has been widely embraced by industry. Many of the early computers (before 1986) laid claim to being RISC machines. However, these claims were often born more of marketing ambition than of engineering reality.

In 1986, the computer industry began to announce processors based on the technology explored by the three RISC research projects. Moussoris et al. [1986]

describe the MIPS R2000 integer processor; while Kane [1987] is a complete description of the architecture. Hewlett-Packard converted their existing minicomputer line to RISC architectures; the HP Precision Architecture is described by Lee [1989]. IBM never directly turned the 801 into a product. Instead, the ideas were adopted for a new, low-end architecture that was incorporated in the IBM RT-PC and is described in a collection of papers [Waters 1986]. In 1990, IBM announced a new RISC architecture (the RS 6000), which is the first super scalar RISC machine (see chapter 6). In 1987, Sun Microsystems began delivering machines based on the SPARC architecture, a derivative of the Berkeley RISC-II machine; SPARC is described in Garner et al. [1988]. Starting in 1987, semiconductor manufacturers began to become suppliers of RISC microprocessors. With its announcement of the AMD 29000, AMD was the first major semiconductor manufacturer to deliver a RISC machine. In 1988, Motorola announced the availability of its RISC machine, the 88000.

Prior to the RISC architecture movement, the major trend had been highly microcoded architectures aimed at reducing the semantic gap. DEC, with the VAX, and Intel, with the iAPX 432, were among the leaders in this approach. In 1989, DEC and Intel both announced RISC products—the DECstation 3100 (based on the MIPS Computer Systems R2000) and the Intel i860, a new RISC microprocessor. With these announcements (and the IBM RS6000), RISC technology has achieved very broad acceptance. In 1990 it is hard to find a computer company without a RISC product.

References

- ADAMS, T. AND R. ZIMMERMAN [1989]. "An analysis of 8086 instruction set usage in MS DOS programs," *Proc. Third Symposium on Architectural Support for Programming Languages and Systems* (April) Boston, 152–161.
- ALEXANDER, W. G. AND D. B. WORTMAN [1975]. "Static and dynamic characteristics of XPL programs," *Computer* 8:11 (November) 41–46.
- AMDAHL, G., G. BLAAUW, AND F. BROOKS [1964]. "Architecture of the IBM System/360," *IBM J. of Research and Development* 8:2 (April) 87–101.
- CASE, R. AND A. PADEGS [1978]. "Architecture of the IBM System/370," *Comm. ACM* 21:1 (January) 73–96.
- CHOW, F., M. HIMELSTEIN, E. KILLIAN, AND L. WEBER [1986]. "Engineering a RISC compiler system," *Proc. COMPCON* (March), San Francisco, 132–137.
- CLARK, D. AND H. LEVY [1982]. "Measurement and analysis of instruction set use in the VAX-11/780," *Proc. Ninth Symposium on Computer Architecture* (April), Austin, Tex., 9–17.
- CRAWFORD, J. AND P. GELSINGER [1988]. *Programming the 80386*, Sybex Books, Alameda, Calif.
- EMER, J. S. AND D. W. CLARK [1984]. "A characterization of processor performance in the VAX-11/780," *Proc. 11th Symposium on Computer Architecture* (June), Ann Arbor, Mich., 301–310.
- GARNER, R., A. AGARWAL, F. BRIGGS, E. BROWN, D. HOUGH, B. JOY, S. KLEIMAN, S. MUNCHNIK, M. NAMJOO, D. PATTERSON, J. PENDLETON, AND R. TUCK [1988]. "Scalable processor architecture (SPARC)," *COMPCON, IEEE* (March), San Francisco, 278–283.

- HENNESSY, J. [1984]. "VLSI processor architecture," *IEEE Trans. on Computers* C-33:11 (December) 1221-1246.
- HENNESSY, J. [1985]. "VLSI RISC processors," *VLSI Systems Design* VI:10 (October) 22-32.
- HENNESSY, J., N. JOUPPI, F. BASKETT, AND J. GILL [1981]. "MIPS: A VLSI processor architecture," *Proc. CMU Conf. on VLSI Systems and Computations* (October), Computer Science Press, Rockville, Md.
- KANE, G. [1986]. *MIPS R2000 RISC Architecture*, Prentice Hall, Englewood Cliffs, N.J.
- LEE, R. [1989]. "Precision architecture," *Computer* 22:1 (January) 78-91.
- LEVY, H. AND R. ECKHOUSE [1989]. *Computer Programming and Architecture: The VAX*, Digital Press, Boston.
- MORSE, S., B. RAVENAL, S. MAZOR, AND W. POHLMAN [1980]. "Intel Microprocessors—8008 to 8086," *Computer* 13:10 (October).
- MOUSSOURIS, J., L. CRUDELE, D. FREITAS, C. HANSEN, E. HUDSON, S. PRZYBYLSKI, T. RIORDAN, AND C. ROWEN [1986]. "A CMOS RISC processor with integrated system functions," *Proc. COMPCON, IEEE* (March), San Francisco.
- PATTERSON, D. [1985]. "Reduced Instruction Set Computers," *Comm. ACM* 28:1 (January) 8-21.
- PATTERSON, D. A. AND D. R. DITZEL [1980]. "The case for the reduced instruction set computer," *Computer Architecture News* 8:6 (October), 25-33.
- RADIN, G. [1982]. "The 801 minicomputer," *Proc. Symposium Architectural Support for Programming Languages and Operating Systems* (March), Palo Alto, Calif. 39-47.
- SHUSTEK, L. J. [1978]. "Analysis and performance of computer instruction sets," Ph.D. Thesis (May), Stanford Univ., Stanford, Calif.
- STRECKER, W. [1978]. "VAX-11/780: A virtual address extension to the DEC PDP-11 family," *Proc. AFIPS NCC* 47, 967-980.
- STRECKER, W. D. AND C. G. BELL [1976]. "Computer structures: What have we learned from the PDP-11?," *Proc. Third Symposium on Computer Architecture*.
- TAYLOR, G., P. HILFINGER, J. LARUS, D. PATTERSON, AND B. ZORN [1986]. "Evaluation of the SPUR LISP architecture," *Proc. 13th Symposium on Computer Architecture* (June), Tokyo.
- UNGAR, D., R. BLAU, P. FOLEY, D. SAMPLES, AND D. PATTERSON [1984]. "Architecture of SOAR: Smalltalk on a RISC," *Proc. 11th Symposium on Computer Architecture* (June), Ann Arbor, Mich., 188-197.
- WAKERLY, J. [1989]. *Microcomputer Architecture and Programming*, J. Wiley, New York.
- WATERS, F., ED. [1986]. *IBM RT Personal Computer Technology*, IBM, Austin, Tex., SA 23-1057.
- WIECEK, C. [1982]. "A case study of the VAX 11 instruction set usage for compiler execution," *Proc. Symposium on Architectural Support for Programming Languages and Operating Systems* (March), IEEE/ACM, Palo Alto, Calif., 177-184.

EXERCISES

In these exercises you will often need to know the frequency of individual instructions in a mix. Figures C.1 through C.4 supply the data corresponding to Figures 4.24, 4.28, 4.32, and 4.34. Additionally, some problems involve the execution-time distribution rather than the frequency distribution. The information on instruction-time distribution appears in Appendix D; problems that require data from Appendices C or D include the letter C or D within the brackets, e.g., <C,D>.

In doing these exercises you will need to work with measurements that may not total 100%. In some cases you will need to normalize the data to the actual total. For example, if we were asked to find the frequency of MOV_ instructions in Spice running on the VAX, we would proceed as follows (using data from Figure C.1):

$$\text{Frequency of measured MOV_ in table} = 9\% + 6\% = 15\%$$

$$\text{Fraction of all instructions executed included in Figure C.1 for Spice} = 79\%$$

We now normalize the 15%. This is equivalent to assuming that the unmeasured 21% of the instruction mix behaves the way as the measured portion. Since there are unmeasured MOV_ instructions this is the most logical approach.

$$\text{Frequency of MOV_ in Spice on VAX} = \frac{15\%}{79\%} = 19\%$$

If, however, we were asked to find the frequency of MOVL in Spice, we know that it is exactly 9%, since we have a complete measurement for this instruction type.

4.1 [20/20] <4.2,4.6,C> You are being interviewed by Digital Equipment Corporation for a job as lead computer designer of future VAX computers. To see if you know what you are talking about, before they hire you they want to ask you a few questions. They have allowed you to bring your notes, including Section 4.6 and Appendix C.

You remember an example in Chapter 4 where you were told that the average VAX instruction had 1.8 operands. You also recall that opcodes are almost always 1 byte long.

- a. [20] They ask you to derive the average size of a VAX instruction for the TeX benchmark. Use the addressing-mode frequency data in 4.22 and 4.23, the information on sizes of displacements in Figure 3.35 (page 133), the information on immediate sizes in Figure 3.15 (page 102), and the length of the VAX addressing modes shown in Figure 4.3. (This should be a more accurate estimate than the example that appears on page 170, but ignore addressing modes that account for less than 5% of the occurrences.)
- b. [20] They then ask you to evaluate the performance of their new machine with a 100-MHz clock. They tell you that the average CPI for everything except instruction fetch and operand fetch is 3 clocks. They also tell you that
 - each data memory specifier and access takes an additional 2 clocks, and
 - every 4 bytes of instructions fetched by the instruction fetch unit take one clock.

Can you find the effective native MIPS?

4.2 [20/22/22] <4.2,4.3,4.5> Consider the following fragment of C code:

```
for (i=1; i<=100; i++)
  {A[i] = B[i] + C;}
```

Assume that A and B are arrays of 32-bit integers, and C and i are 32-bit integers. Assume that all data values are kept in memory (at addresses 0, 5000, 1500, and 2000 for A, B, C, and i, respectively) except when they are operated on.

- a. [20] Write the code for DLX; how many instructions are required dynamically? How many memory data references will be executed? What is the code size?

- b. [22] Write the code for the VAX; how many instructions are required dynamically? How many memory data references will be executed? What is the code size?
- c. [22] Write the code for the 360; how many instructions are required dynamically? How many memory data references will be executed? What is the code size? For simplicity, you may assume that register R1 contains the address of the first instruction in the loop.

4.3 [20/22/22] <4.2,4.3,4.5> For this question use the code sequence of problem 4.2, but put the scalar data—the value of *i* and the address of the array variables (but not the actual array)—in registers and keep them there whenever possible.

- a. [20] Write the code for DLX; how many instructions are required dynamically? How many memory-data references will be executed? What is the code size?
- b. [22] Write the code for the VAX; how many instructions are required dynamically? How many memory data references will be executed? What is the code size?
- c. [22] Write the code for the 360; how many instructions are required dynamically? How many memory data references will be executed? What is the code size? Assume R1 is set-up as in Exercise 4.2 part C.

4.4 [15] <4.6> When designing memory systems it becomes useful to know the frequency of memory reads versus writes and also accesses for instructions versus data. Using the average instruction-mix information for DLX in Appendix C, find

- the percentage of all memory accesses that are for data
- the percentage of data accesses that are reads
- the percentage of all memory accesses that are reads

Ignore the size of a datum when counting accesses.

4.5 [15] <4.3,4.6> Due to the lack of a PC-relative branch, a branch on a 360 often requires two instructions. This has been a major criticism of the architecture. Let's figure out what this omission costs, assuming that an extra instruction is always needed for a conditional branch on the 360, but that the extra instruction would not be necessary with PC-relative branches. Using the average data from Figure 4.28 (page 175) for branches, determine how many more instructions the standard 360 executes than a 360 with PC-relative branches. (Remember that the only branches are BC and BCR.)

4.6 [15] <4.2,4.6> We are interested in adding an instruction to the VAX architecture that compares an operand to zero and branches. Assume that

- only instructions that set the condition code for a conditional branch could be eliminated,
- 80% of the conditional branches require an instruction whose only purpose is to set the condition, and
- 90% of all branches that have an instruction that just sets the condition (i.e., the just-mentioned 80%) are based on a compare against 0.

Using the average VAX data from Figure 4.24 (page 172) what percentage more instructions would a standard VAX execute compared to the VAX with the compare-and-branch instruction added?

4.7 [18] <4.5,4.6> Compute the effective CPI for DLX. Suppose we have made the following measurements of average CPI for instructions:

All R-R instructions	1 clock cycle
Loads/stores	1.4 clock cycles
Conditional branches	
	taken 2.0 clock cycles
	not taken 1.5 clock cycles
Jumps	1.2 clock cycles

Assume that 60% of the conditional branches are taken. Average the instruction frequencies of GCC and TeX to obtain the instruction mix.

4.8 [15] <4.2,4.5> Rather than have immediates supported for many instruction types, some architectures, such as the 360, collect immediates in memory (in a literal pool) and access them from there. Suppose the VAX didn't have immediate-mode addressing, but instead put immediates in memory and accessed them using displacement-mode addressing. What would be the increase in the frequency that displacement mode was used? Use the average of the measurements in Figures 4.22 and 4.23 for this problem.

4.9 [20/10] <4.5,4.6> Consider adding a new index addressing mode to DLX. The addressing mode adds two registers and an 11-bit signed offset to get the effective address.

Our compiler will be changed so that code sequences of the form

```
ADD R1, R1, R2
LW  Rd, O(R1)    (or store)
```

will be replaced with a load (or store) using the new addressing mode. Use the overall average instruction frequencies in evaluating this addition.

- [20] Assume that the addressing mode can be used for 10% of the displacement loads and stores (accounting for both the frequency of this type of address calculation and the shorter offset). What is the ratio of instruction count on the enhanced DLX compared to the original DLX?
- [10] If the new addressing mode lengthens the clock cycle by 5%, which machine will be faster and by how much?

4.10 [12] <4.2,4.5,D> Assume the average number of instructions involved in a call and return on DLX is 8. The average frequency of a JAL instruction in the benchmarks is 1%. If all instructions on DLX take the same number of cycles, how does the percentage of cycles in calls and returns on DLX compare to the percentage of cycles in CALLS and RET on the VAX?

4.11 [22/22] <4.2,4.3,4.6,D> Some people believe that the most frequent instructions are also the simplest, while others have pointed out that the most time-consuming instructions are often not the most frequent.

- a. [22] Using the data in Figure D.1, find the CPI of the five most time-consuming instructions on the VAX that have an average execution frequency of at least 2%. Assume the overall VAX CPI is 10.
- b. [22] Find the CPI for the five most time-consuming instructions on the 360 that have at least a 3% average frequency, using the data in Figure D.2. Assume the overall 360 CPI is 4.

4.12 [20/20/10] <4.4, 4.6,D> You have been hired to try to convert the 8086 architecture to be more register-register oriented. To do this, you will need more registers, and hence more encoding space, since the encodings are already tight. Assume that you have determined that eliminating the PUSH and POP instructions can yield the encoding space needed. Suppose that increasing the number of registers reduces the frequency of each of the memory-referencing instructions (PUSH, POP, LES, and MOV) by 25%, but that each remaining PUSH or POP instruction must be replaced by a two-instruction sequence. Use the average data from Figures 4.30–4.32 (pages 177–178), the average CPI of 14.1, and Figure D.5 to answer the following questions about this new machine—the RR8086—versus the 8086.

- a. [20] Which machine executes more instructions and by how much?
- b. [20] Using the information in Appendix D, determine which machine has a higher CPI and by how much?
- c. [10] Assuming the clock rates are identical, which machine is faster and by how much?

4.13 [25/15] <4.2–4.5> Find a C compiler and compile the code shown in Exercise 4.2 for a load/store machine or one of the machines covered in this chapter. Compile the code both optimized and unoptimized.

- a. [25] Find the instruction count, dynamic instruction bytes fetched, and data accesses done for both the optimized and unoptimized versions.
- b. [15] Try to improve the code by hand, and compute the same measures as in Part a for your hand-optimized version.

4.14 [30] <4.6> If you have access to a VAX, compile the code for Spice and try to determine why it makes much smaller use of immediates than programs like GCC and TeX (see Figure 4.22 on page 169).

4.15 [30] <4.6> If you have access to an 8086-based machine, compile some programs and look at the frequency of MOV instructions. How does it correspond to the frequency in Figure 4.32 (page 178). By examining the code, can you find some reasons why the frequency of MOVs is so high?

4.16 [30/30] <4.6, 4.7> Small synthetic benchmarks can be very misleading when used for measuring instruction mixes. This is particularly true when these benchmarks are

optimized. In these exercises we want to explore these differences. These programming exercises can be done with a VAX, any load/store machine, or using the DLX compiler and simulator.

- a. [30] Compile Whetstone with optimization for a VAX, or a load/store machine similar to DLX (e.g., a DECstation or a SPARCstation), or the DLX simulator. Compute the instruction mix for the top twenty instructions. How do the optimized and unoptimized mixes compare? How does the optimized mix compare to the mix for Spice on the same or a similar machine?
- b. [30] Compile Dhrystone with optimization for a VAX, or a load/store machine similar to DLX (e.g., a DECstation or a SPARCstation), or the DLX simulator. Compute the instruction mix for the top twenty instructions. How do the optimized and unoptimized mixes compare? How does the optimized mix compare to the mix for TeX on the same or a similar machine?

4.17 [30] <4.6> Many computer manufacturers now include tools or simulators that allow you to measure the instruction set usage of a user program. Among the methods in use are machine simulation, hardware-supported trapping, and a compiler technique that instruments the object-code module by inserting counters. Find a processor available to you that includes such a tool. Use it to measure the instruction set mix for one of TeX, GCC, or Spice. Compare the results to those shown in this chapter.

4.18 [30] <4.5,4.6> DLX has only three operand formats for its register-register operations. Many operations might use the same destination register as one of the sources. We could introduce a new instruction format into DLX called R_2 that has only two operands and is a total of 24 bits in length. By using this instruction type whenever an operation had only two different register operands, we could reduce the instruction bandwidth required for a program. Modify the DLX simulator to count the frequency of register-register operations with only two different register operands. Using the benchmarks that come with the simulator, determine how much more instruction bandwidth DLX requires than DLX with the R_2 format.

4.19 [35] <D> Devise a method to measure the CPI of a machine—preferably one of the machines discussed in this chapter or a relative of DLX. Using the instruction-mix data, choose the top ten instructions and measure their CPI. How does the frequency ranking compare to the time taken? How do your measurements compare to the numbers shown in Appendix D? Try to explain any differences in both time-versus-frequency ranking and any differences between your measures and those in Appendix D.

4.20 [35] <4.5,4.6> What are the benefits of more powerful addressing modes? Assume that three VAX addressing modes—autoincrement, displacement deferred, and scaled—were added to DLX. Change the C compiler to incorporate the use of these modes. Measure the change in instruction count with these new modes for several benchmark programs. Compare the instruction mixes with those for standard DLX. How do the usage patterns compare to those for the VAX shown in Figure 4.23 (page 170)?

4.21 [35/35/30] <4.5,4.6> How much does the flexibility of memory–memory instructions reduce instruction count compared to a load/store machine? This programming assignment will help you find out.

- a. [35] Assume DLX has an instruction format that allows one of the source operands to be in memory. Modify the C code generator for DLX so that it uses this new instruction type. Use several C programs to measure the effectiveness of your change. How many more instructions does DLX require versus this new machine that appears to be closer to the 360? How often is the register–memory format used? How do the instruction mixes differ from those in Section 4.6?
- b. [35] Assume that DLX has instruction formats that allow any operand (or all three) to be memory references. Modify the C code generator for DLX so that it uses these new instruction formats. Use several programs to measure the usage of these instructions. How many more instructions does DLX require versus this new machine that appears to be closer to the VAX? How do the instruction mixes differ from those in Section 4.6? How many memory operands does the average instruction have?
- c. [30] Design an instruction format for the machines described in Parts a and b; compare the dynamic instruction bandwidth required for these two machines versus DLX.

4.22 [40] <4.6> Some manufacturers have not yet seen the value of measuring instruction set mixes. Maybe you can help them. Pick a machine for which such a tool is not widely available. Construct one for that machine. If the machine has a single-step mode—as in the VAX or 8086—you can use it to create your tool. Otherwise, an object code translation, as used in the MIPS compiler system [Chow 1986] might be more appropriate. If you measure the activity of a machine using the benchmarks in this text (GCC, Spice, and TeX), and are willing to share the results, please contact the publisher.

4.23 [25] <E> How much do the instruction set variations among the RISC machines discussed in Appendix E affect performance. Choose at least three small programs (e.g., a sort), and code these programs in DLX and two other assembly languages. What is the resulting difference in instruction count?

4.24 [40] <E> Choose one of the machines discussed in Appendix E. Modify the DLX code generator and DLX simulator to generate code for and simulate the machine you chose. Using as many benchmarks as practical, measure the instruction count differences seen between DLX and the machine you chose.

In analyzing the functions of the contemplated device, certain classificatory distinctions suggest themselves immediately ... First: Since the device is primarily a computer it will have to perform the elementary operations of arithmetic most frequently ... a central arithmetic part of the device will probably have to exist... Second: The logical control of the device, that is the proper sequencing of its operations, can be most efficiently carried out by a central control organ.

John von Neumann, *First Draft of a Report on the EDVAC* (1945)

5.1	Introduction	199
5.2	Processor Datapath	201
5.3	Basic Steps of Execution	202
5.4	Hardwired Control	204
5.5	Microprogrammed Control	208
5.6	Interrupts and Other Entanglements	214
5.7	Putting It All Together: Control for DLX	220
5.8	Fallacies and Pitfalls	238
5.9	Concluding Remarks	240
5.10	Historical Perspective and References	241
	Exercises	244

5

Basic Processor Implementation Techniques

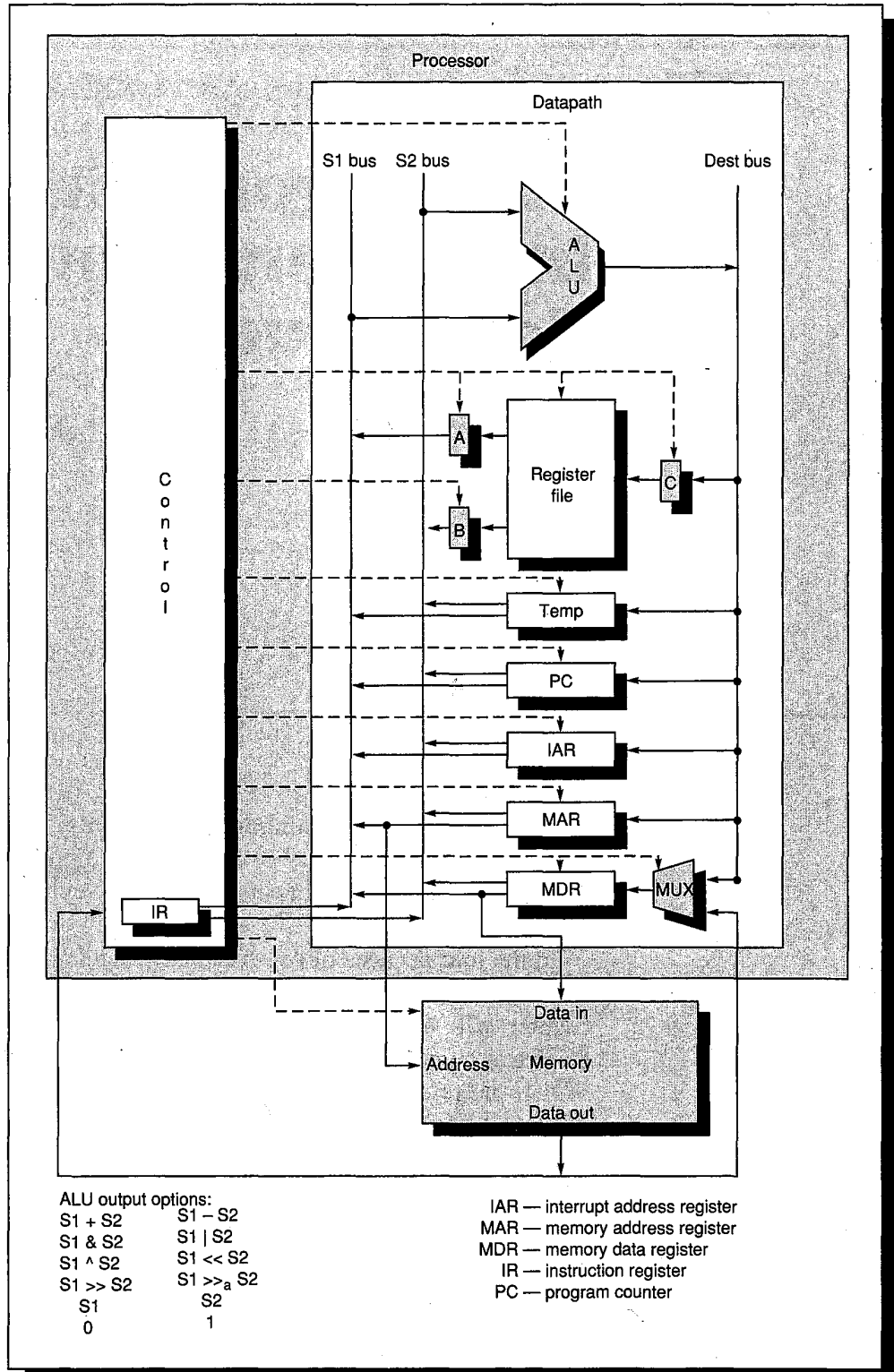
5.1

Introduction

Architecture shapes a building, but carpentry determines the quality of its construction. The carpentry of computing is implementation, which sets two of three performance components: CPI (clock cycles per instruction) and clock cycle time.

In the four decades of constructing computers, much has been learned about implementation—certainly more than can fit in one chapter. Our goal in this chapter will be to lay the foundations of processor implementation, with emphasis on control and interrupts. (Floating point is ignored in this chapter; readers are referred to Appendix A.) While some material is simple, Chapters 6 and 7 build on this foundation and show the road to faster computers. (If this is a review, go quickly over Sections 5.1 to 5.3 and then take a look at the examples in Section 5.7, which compare performance of hardwired versus microprogrammed control for DLX.)

The computer was divided into basic components by von Neumann, and these components still remain today: The CPU, or the *processor*, is the core of the computer and contains everything except the memory, input, and output. The processor is further divided into computation and control.



5.2 Processor Datapath

Today the “arithmetic” organ of von Neumann is called the *datapath*. It consists of execution units, such as arithmetic logic units (ALUs) or shifters, the registers, and the communication paths between them, as Figure 5.1 illustrates. From the programmer’s perspective, the datapath contains most of the *state* of the processor—the information that must be saved for a program to be suspended and then restored for execution to continue. In addition to the user-visible general-purpose registers, the state includes the program counter (PC), the interrupt address register (IAR), and the program status register; the latter contains all the status flags for a machine, such as interrupt enable, condition codes, and so forth.

Because an implementation is created for a specific hardware technology, it is the implementation that sets the clock cycle time. The clock cycle time in turn is determined by the slowest circuits that operate during a clock cycle period—within the processor, the datapath frequently has that honor. The datapath will also dominate the cost of the processor, typically requiring half the transistors and half the processor area. While it does all the computation, affects performance, and dominates cost, the datapath is the simplest portion of the processor to design.

Some have held the large quantity of papers on ALU designs and fast-carry schemes responsible for the loss of our rain forests, with papers on circuit designs for registers with multiple read and write ports only slightly less culpable. While this is surely an exaggeration, there are numerous options. (See Appendix A, Section A.8, for a few carry schemes.) Given the resources available and the desired goals of cost and performance, it is the designer’s job to select the best style of ALU, the proper number of ports in the register file, and then march onward.

FIGURE 5.1 (See adjoining page.) A typical processor, divided into control and datapath, plus memory. The paths for control are in dashed lines and the paths for data transfer are in solid lines. The processor uses three buses: S1, S2, and Dest. The fundamental operation of the datapath is reading operands from the register file, operating on them in the ALU, and then storing the result back. Since the register file does not need to be read and written every clock cycle, most designers follow the advice of making the frequent case fast by breaking this sequence into multiple clock cycles and making the clock cycle shorter. Thus, in this datapath there are latches on the two outputs of the register file (called A and B) and a latch on the input (C). The register file contains the 32 general-purpose registers of DLX. (Register 0 of the register file always has the value 0, matching the definition of register 0 in the DLX instruction set.) The program counter (PC) and interrupt address register (IAR) are also part of the state of the machine. There are also registers, not part of the state, used in the execution of instructions: memory address register (MAR), memory data register (MDR), instruction register (IR), and temporary register (Temp). The Temp register is a scratch register that is available for temporary storage for control to perform some DLX instructions. Note that the only path from the S1 and S2 buses to the Dest bus is through the ALU.

5.3 Basic Steps of Execution

Before discussing control, let's first review the steps of instruction execution. For the DLX instruction set (excluding floating point), all instructions can be broken into five basic steps: fetch, decode, execute, memory-access, and write-result. Each step may take one or several clock cycles in the processor shown in Figure 5.1 (page 200). Here are the five steps (see the page facing the inside back cover for a review of the register transfer language notation):

1. Instruction fetch step:

$$\text{MAR} \leftarrow \text{PC}; \text{IR} \leftarrow \text{M}[\text{MAR}]$$

Operation: Send out the PC and fetch the instruction from memory into the instruction register. PC is transferred to MAR because it has a connection to the memory address in Figure 5.1, but PC doesn't.

2. Instruction decode/register fetch step:

$$\text{A} \leftarrow \text{Rs1}; \text{B} \leftarrow \text{Rs2}; \text{PC} \leftarrow \text{PC} + 4$$

Operation: Decode the instruction and access the register file to read the registers. Also, increment the PC to point to the next instruction.

Decoding can be done in parallel with reading registers, which means that two registers' values are sent to the A and B latches **before** the instruction is decoded. How this is possible can be seen by glancing at the DLX instruction format (Figure 4.19 on page 166), which shows that the *source registers* are always at the same location in an instruction. Thus, registers can be read because the register specifiers are unambiguous. (This technique is known as *fixed-field decoding*.) Since the immediate portion of an instruction is also identical in every DLX format, the sign-extended immediate is also calculated during this step in case it is needed in the next step.

3. Execution/effective address step:

The ALU is operating on the operands prepared in the prior step, performing one of three functions depending on the DLX instruction type.

Memory reference:

$$\text{MAR} \leftarrow \text{A} + (\text{IR}_{16})^{16}\#\#\text{IR}_{16..31}; \text{MDR} \leftarrow \text{Rd}$$

Operation: The ALU is adding the operands to form the effective address, and the MDR is loaded for a store.

ALU instruction:

$$\text{ALUoutput} \leftarrow A \text{ op } (B \text{ or } (\text{IR}_{16})^{16} \# \# \text{IR}_{16..31})$$

Operation: The ALU is performing the operation specified by the opcode on the value in A (Rs1) and on the value in B or the sign-extended immediate.

Branch/Jump:

$$\text{ALUoutput} \leftarrow \text{PC} + (\text{IR}_{16})^{16} \# \# \text{IR}_{16..31}; \text{ cond} \leftarrow (A \text{ op } 0)$$

Operation: The ALU is adding the PC to the sign-extended immediate value (16-bit for branch and 26-bit for jump) to compute the address of the branch target. For conditional branches, a register, which has been read in the prior step, is checked to decide if this address should be inserted into the PC. The comparison operation *op* is the relational operator determined by the opcode; for example, *op* is “==” for the instruction BEQZ.

The load/store architecture of DLX means that effective address and execution steps can be combined into a single step, since no instruction needs to both calculate an address and perform an operation on the data. Other integer instructions not included above are JAL and TRAP. These are similar to jumps, except JAL stores the return address in R31 and TRAP stores it in IAR.

4. Memory access/branch completion step: The only DLX instructions active in this step are loads, stores, branches, and jumps.

Memory reference:

$$\text{MDR} \leftarrow \text{M}[\text{MAR}] \text{ or } \text{M}[\text{MAR}] \leftarrow \text{MDR}$$

Operation: Access memory if needed. If instruction is a load, data returns from memory; if it is a store, then the data writes into memory. In either case the address used is the one computed during the prior step.

Branch:

$$\text{if (cond) PC} \leftarrow \text{ALUoutput (branch)}$$

Operation: If the instruction branches, the PC is replaced with the branch destination address. For jumps the condition is always true.

5. Write-back step:

$$\text{Rd} \leftarrow \text{ALUoutput or MDR}$$

Operation: Write the result into the register file, whether coming from the memory system or from the ALU.

Now that we have had an overview of the work that must be performed to execute an instruction, we are ready to look at the two main techniques for implementing control.

5.4 Hardwired Control

If the datapath design is simple, then some part of processor design must be difficult, and that part is control. Specifying control is on the critical path of any computer project; and it is where most errors are found when a new computer is debugged. Control can be simplified—the easiest way is to simplify the instruction set—but that is the subject of Chapters 3 and 4.

Given an instruction set description, such as the description of DLX in Chapter 4, and a datapath design, such as Figure 5.1 (page 200), the next step is defining the control unit. The control unit tells the datapath what to do every clock cycle during the execution of instructions. This is typically specified by a *finite-state diagram*. Every state corresponds to one clock cycle, and the operations to be performed during the clock cycle are written within the state. Each instruction takes several clock cycles to complete; Chapter 6 shows how to overlap execution to reduce the clock cycles per instruction to as low as one.

Figure 5.2 shows a portion of a finite-state diagram for the first two steps of instruction execution in DLX. The first step is spread over all three states: The memory-address register is loaded from PC during the first state, the instruction register is loaded from memory during the second state, and the PC is incremented in the third state. This third state also performs step 2, loading the two register operands, Rs1 and Rs2, into the A and B registers for use in the later states. In Section 5.7 the full finite-state diagram for DLX is shown.

Turning a state diagram into hardware is the next step. The alternatives for doing this depend on the implementation technology. One way to bound the complexity of control is by the product

$$\text{States} * \text{Control inputs} * \text{Control outputs}$$

where

States = the number of states in the finite-state machine controller;

Control inputs = the number of signals examined by the control unit;

Control outputs = the number of control outputs generated for the hardware, including bits to specify the next state.

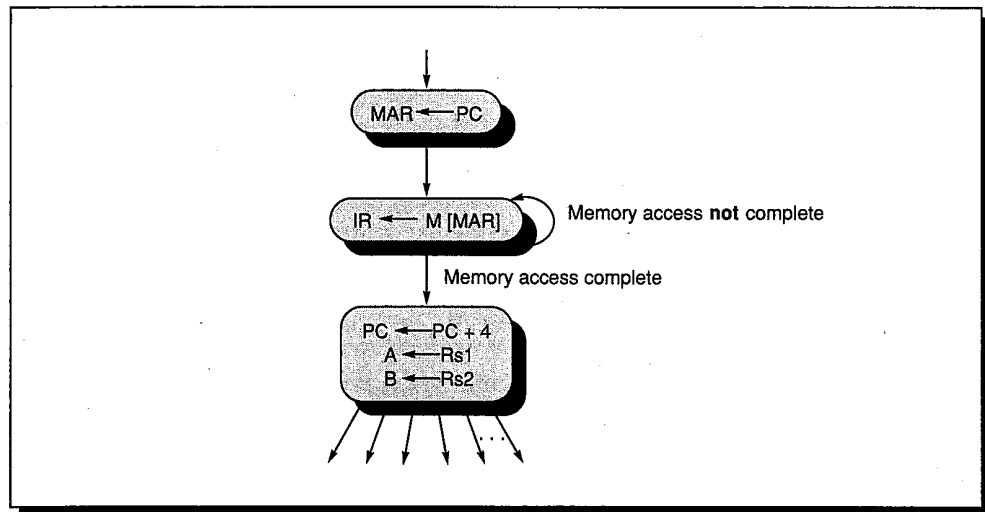


FIGURE 5.2 The top level of the DLX finite-state diagram. The first two steps of instruction execution, instruction fetch and instruction decode/register fetch, are shown. The second state repeats until the instruction is fetched from memory. The last three steps of instruction execution—execution/effective address, memory access, and write back—are found in Section 5.7.

Figure 5.3 shows an organization for control of DLX. Let's say the DLX finite-state diagram contains 50 states, requiring 6 bits to represent the state. Thus, the control inputs must include these 6 bits, some number of bits (say 3) to select conditions from the datapath and memory interface unit, plus instruction bits. Register specifiers and immediates are sent directly to the hardware, so there is no need to send all 32 bits of DLX instructions as control inputs. The DLX opcode is 6 bits, and only 6 bits of the extended opcode (the “func” field) are used, making a total of 12 instruction bits for control inputs. Given those inputs, control can be specified as a big table. Each row of the table contains the values of the control lines to perform the operations required by that state and supplies the next state number. Let's assume there are 40 control lines.

Reducing Hardware Costs of Hardwired Control

The straightforward implementation of a table is with a *read only memory* (ROM). In this example, 2^{21} words, each 40 bits wide (10 MB of ROM!), would be required. It will be a long time before we can afford this much hardware for control. Fortunately, little of this table has unique information, so its size can be reduced by keeping only the rows with unique information—at the cost of more complicated address decoding. Such a hardware construct is called a *programmed logic array* (PLA). This essentially reduces the hardware from 2^{21} words to 50 words while increasing address decoding logic. Computer-aided design programs can reduce the hardware requirements even further by

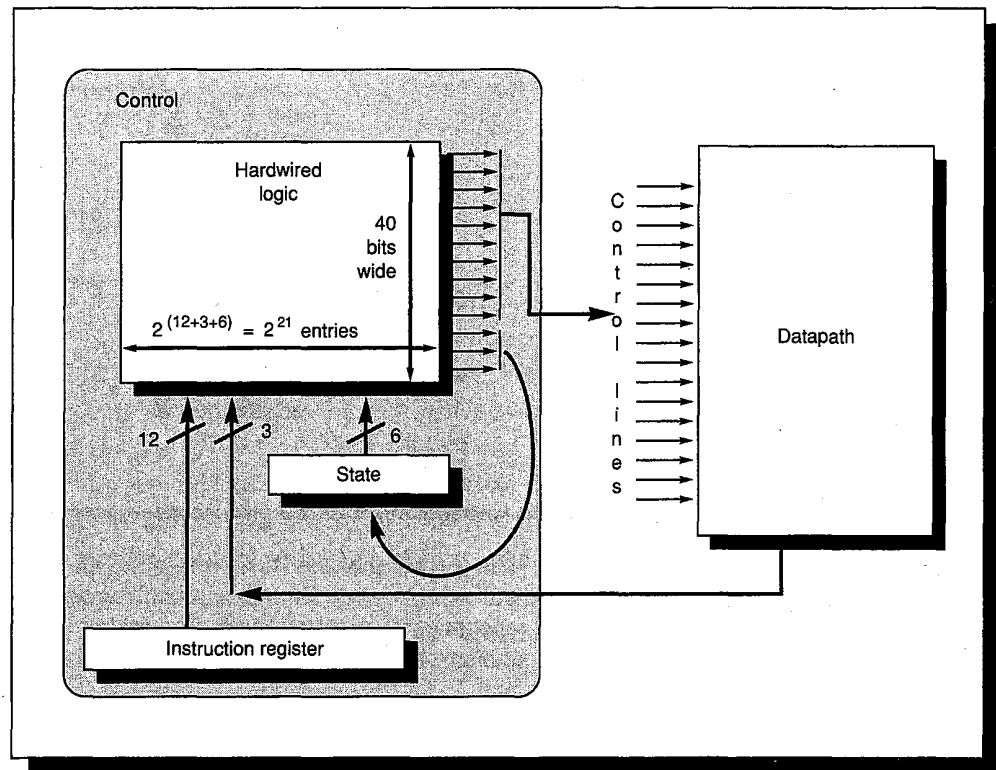


FIGURE 5.3 Control specified as a table for a simple instruction set. The control inputs consist of 6 input lines for the 50 states ($\log_2 50=5.6$), 3 inputs from the datapath, and 12 instruction bits (the 6-bit opcode plus 6 bits of the extended opcode). The number of control lines is assumed to be 40.

minimizing the number of “minterms,” which is essentially the number of unique rows. In real machines, even a single PLA is sometimes prohibitive because its size grows as the product of the unique rows times the sum of the inputs and outputs. In such a case, a large table is factored into several smaller PLAs, whose outputs are multiplexed to choose the correct control.

Oddly enough, the numbering of the states in the finite-state diagram can make a difference in the size of the PLA. The idea here is to try to assign similar state numbers to states that perform similar operations. Differentiating the bit patterns that represent the state number by only one bit—say 010010 and 010011—make the inputs close for the same output. There are also computer-aided design programs to help with this *state-assignment problem*.

Since the instruction bits are also inputs to the control PLA, they can affect the complexity of the PLA just as numbering of the states does. Thus, care should be taken when selecting opcodes since it may affect the cost of control.

Readers interested in taking this design further are referred to the many excellent texts on logic design.

Performance of Hardwired Control

When designing the detailed control for a machine, we want to minimize the average CPI, the clock cycle, the amount of hardware to specify control, and the time to develop a correct controller. Minimizing CPI means reducing the average number of states along the path of execution of an instruction, since each clock cycle corresponds to a state. This is typically done by making changes to the datapath to combine or eliminate states.

Example

Let's change the hardware so that the PC can be used directly to address memory without going through MAR first. How should the state diagram be changed to take advantage of this improvement, and what would be the change in performance?

Answer

From Figure 5.2 (page 205) we see that the first state copies PC into MAR. This proposed hardware change makes that state unnecessary, and Figure 5.4 shows the appropriately modified state diagram. This change saves one clock cycle from every instruction. Suppose the average number of CPI was originally 7. Provided there was no impact on clock cycle time, this change would make the machine 17% faster.

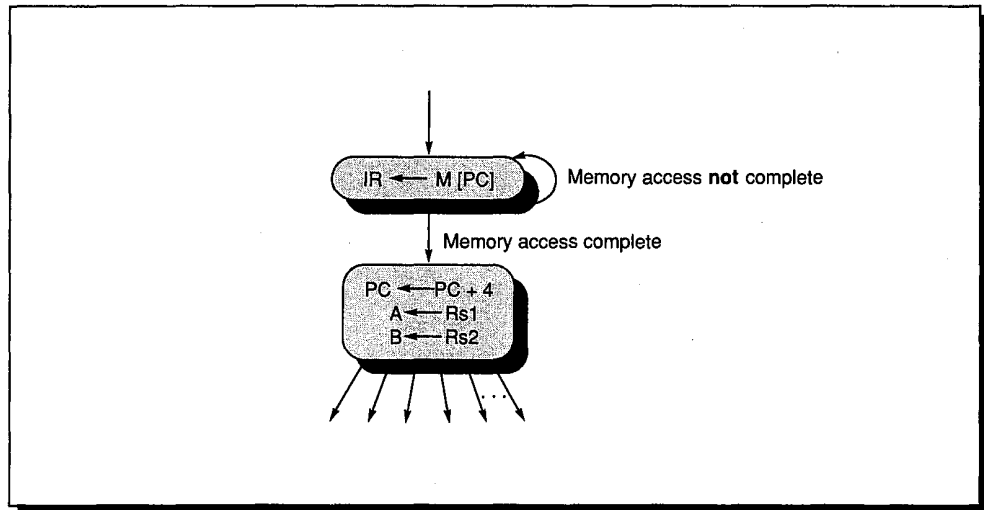


FIGURE 5.4 Figure 5.2 modified to remove the loading of MAR from PC in the first state and to use the PC value directly to address memory.

5.5 Microprogrammed Control

After constructing the first full-scale, operational, stored-program computer in 1949, Maurice Wilkes reflected on the process. I/O was easy—teletypewriters could just be purchased directly from the telegraph company. Memory and the datapath were highly repetitive, and that made things simpler. But control was neither easy nor repetitive, so Wilkes set out to discover a better way to design control. His solution was to turn the control unit into a miniature computer by having a table to specify control of the datapath and a second table to determine control flow at the micro level. Wilkes called his invention *microprogramming* and attached the prefix “micro” to traditional terms used at the control level: microinstruction, microcode, microprogram, and so on. (To avoid confusion the prefix “macro” is sometimes used to describe the higher level, e.g., macroinstruction and macroprogram.) Microinstructions specify all the control signals for the datapath, plus the ability to conditionally decide which micro-

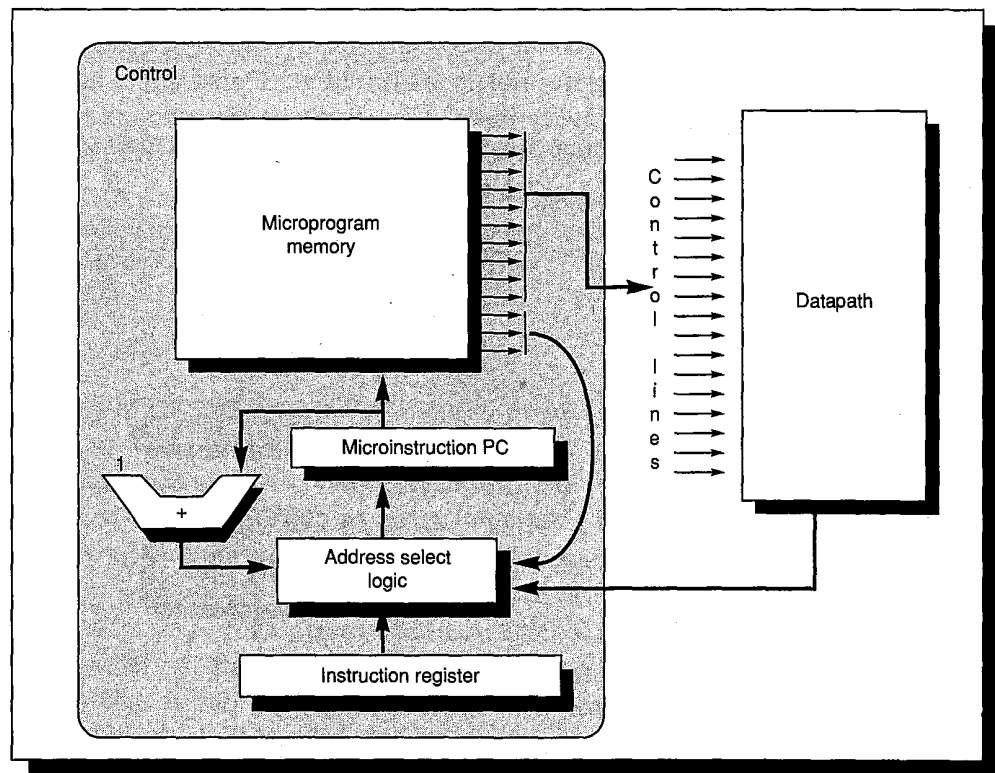


FIGURE 5.5 A basic microcoded engine. Unlike Figure 5.3 (page 206), there is an incrementer and special logic to select the next microinstruction. There are two approaches to specifying the next microinstruction: use a microinstruction program counter, as shown above, or include a next microinstruction address in every microinstruction. Microprogram memory is sometimes called ROM because most early machines use ROM for control stores.

instruction should be executed next. As the name “microprogramming” suggests, once the datapath and memory for the microinstructions are designed, control becomes essentially a programming task; that is, the task of writing an interpreter for the instruction set. The invention of microprogramming enabled the instruction set to be changed by altering the contents of control store without touching the hardware. As we will see in Section 5.10, this ability played an important role in the IBM 360 family—one that was a surprise to its designers.

Figure 5.5 shows an organization for a simple microprogrammed control. The structure of a microprogram is very similar to the state diagram, with a microinstruction for each state in the diagram.

ABCs of Microprogramming

While it doesn’t matter to the hardware how the control lines are grouped within a microinstruction, control lines performing related functions are traditionally placed next to each other for ease of understanding. Groups of related control lines are called *fields* and are given names in a microinstruction format. Figure 5.6 shows a microinstruction format with eight fields, each named to reflect its function. Microprogramming can be thought of as supplying the proper bit pattern in each field, much like assembly language programming of “macroinstructions.”

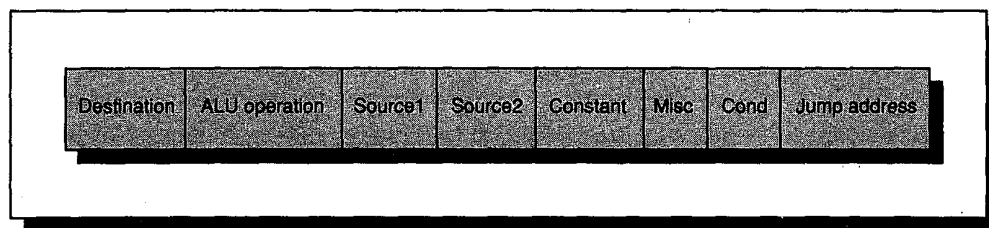


FIGURE 5.6 Example microinstruction with eight fields (used for DLX in Section 5.7).

A program counter can be used to supply the next microinstruction, as shown in Figure 5.5, but some computers dedicate a field in every microinstruction to the address of the next instruction. Some even provide multiple next-address fields to handle conditional branches.

While conditional branches could be used to decode an instruction by testing the opcode one bit at a time, this tedious approach is too slow in practice. The simplest fast instruction decoding scheme is to jam the macroinstruction opcode into the middle of the address of the next microinstruction, similar to an indexed jump instruction in assembly language. A more refined approach is to use the opcode to index a table containing microinstruction addresses that supply the next address, similar to a jump table in assembly code.

The microprogram memory, or *control store*, is the most visible and easily measured hardware in microprogrammed control; hence, it is the focus of techniques to reduce hardware costs. Techniques to trim control-store size include reducing the number of microinstructions, reducing the width of each microinstruction, or both. Just as cost is traditionally measured by control-store size, performance is traditionally measured by CPI. The wise microprogrammer knows the frequency of macroinstructions by using statistics like those in Chapter 4, and hence knows where and how time is best spent—instructions demanding the largest part of execution time are optimized for speed, and the others are optimized for space.

In four decades of microprogramming history there have been a wide variety of terms and techniques for microprogramming. In fact, a workshop has met annually on this subject since 1968. Before looking at a few examples, let us remember that control techniques—whether hardwired or microcoded—are judged by their impact on hardware cost, clock cycle time, CPI, and development time. In the next two sections we will examine how hardware costs can be lowered by reducing control-store size. First we look at two techniques to reduce the width of microinstructions, then one technique to reduce the number of microinstructions.

Reducing Hardware Costs by Encoding Control Lines

The ideal approach to reducing control store is to first write the complete microprogram in a symbolic notation and then measure how control lines are set in each microinstruction. By taking measurements we are able to recognize control bits that can be encoded into a smaller field. If no more than one of, say, 8 lines is set simultaneously in the same microinstruction, then they can be encoded into a 3-bit field ($\log_2 8 = 3$). This change saves 5 bits in every microinstruction and does not hurt CPI, though it does mean the extra hardware cost of a 3-to-8 decoder needed to generate the original 8 control lines. Nevertheless, shaving 5 bits off control-store width will usually overcome the cost of the decoder.

This technique of reducing field width is called *encoding*. To further save space, control lines may be encoded together if they are only occasionally set in the same microinstruction; two microinstructions instead of one are then required when both must be set. As long as this doesn't happen in critical routines, the narrower microinstruction may justify a few extra words of control store.

There are dangers to encoding. For example, if an encoded control line is on the critical timing path, or if the hardware it controls is on the critical path, then the clock cycle time will suffer. A more subtle danger is that a later revision of the microcode might encounter situations where control lines would be set in the same microinstruction, either hurting performance or requiring changes to the hardware that could lengthen the development cycle.

Example

Assume we want to encode the three fields that specify a register on a bus—Destination, Source1, and Source2—in the DLX microinstruction format in Figure 5.6. How many bits of control store can be saved versus unencoded fields?

Answer

Figure 5.7 lists the registers for each source and destination of the datapath in Figure 5.1 (page 200). Note that the destination field must be able to specify that nothing is modified. Without encoding, the 3 fields require $7 + 9 + 9$, or 25 bits. Since $\log_2 7 \approx 2.8$ and $\log_2 9 \approx 3.2$, the encoded fields require $3 + 4 + 4$, or 11 bits. Thus, encoding these 3 fields saves 14 bits per microinstruction.

Number	Destination	Source1/Source2
0	(None)	A/B
1	C	Temp
2	Temp	PC
3	PC	IAR
4	IAR	MAR
5	MAR	MDR
6	MDR	IR (16-bit imm)
7	---	IR (26-bit imm)
8	---	Constant

FIGURE 5.7 The sources and destinations specified in the three fields of Figure 5.6 from the datapath description in Figure 5.1. A and B are not separate entries because A can only transfer on the S1 bus and B can only transfer on the S2 bus (see Figure 5.1 on pages 200–201). The last entry in the third column, Constant, is used by control to specify a constant needed in an ALU operation (e.g., 4). See Section 5.7 for its use.

Reducing Hardware Costs with Multiple Microinstruction Formats

Microinstructions can be made narrower still if they are broken into different formats and given an opcode or *format field* to distinguish them. The format field gives all the unspecified control lines their default values, so as not to change anything else in the machine, and is similar to the opcode of a macroinstruction.

Reducing hardware costs by using format fields has its own performance cost—namely, executing more microinstructions. Generally, a microprogram using a single microinstruction format can specify any combination of operations in a datapath and will take fewer clock cycles than a microprogram made up of restricted microinstructions. Narrower machines are cheaper because memory chips are also narrow and tall: It takes many fewer chips for a 16K word by 24-bit memory than for a 4K word by 96-bit memory. (When control memory is on the processor chip, this hardware advantage is no longer true.)

This narrow but tall approach is often called *vertical microcode*, while the wide but short approach is called *horizontal microcode*. It should be noted that the terms “vertical microcode” and “horizontal microcode” have no universal definition—the designers of the 8086 considered its 21-bit microinstruction to be more horizontal than other single-chip computers of the time. The related terms *maximally encoded* and *minimally encoded* lead to less confusion.

Figure 5.8 plots control-store size against microinstruction width for three families of computers. Notice that for each family the total size is similar, even though the width varies by a factor of 6. As a rule, minimally encoded control stores use more bits, and the narrow but tall aspect of memory chips means that maximally encoded control stores naturally have more entries. Sometimes designers of minimally encoded machines don't have the option of shorter RAM chips, causing wide microinstruction machines to end up with many words of control store. Since the hardware costs are not lower if microcode doesn't use up all the space in control store, machines in this class can end up with much larger control stores than expected from other implementations. The ECL RAMs available to build the VAX 8800, for example, led to 2000 K bits of control store.

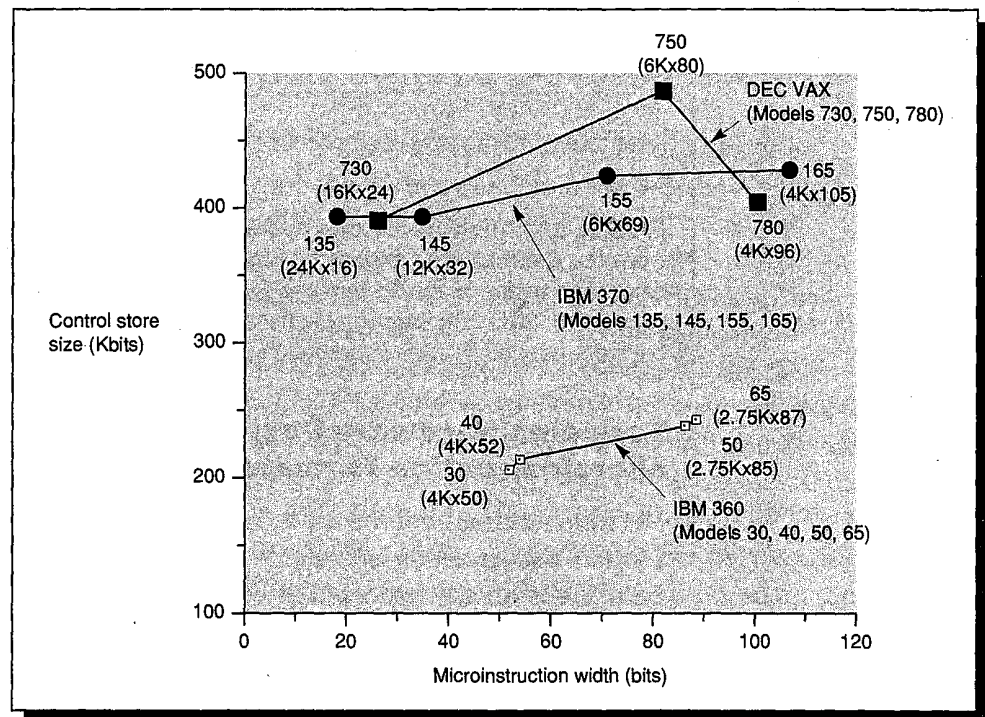


FIGURE 5.8 Size of control store versus width of microinstructions for 11 computer models. Each point is identified by the length and width of control store (not including parity). Models selected from each family are ones that shipped about the same time: IBM 360 models 30, 40, 50, and 65 all shipped in 1965; IBM 370 models 145, 155, and 165 shipped in 1971, with the 135 following in the next year; and the VAX model 780 was shipped in 1978, followed by the 750 in 1980 and the 730 in 1982. The development of the VAX designs all overlapped one another inside DEC.

Reducing Hardware Costs by Adding Hardwired Control to Share Microcode

The other approach to reducing control store is to reduce the number of microinstructions rather than their width. Microsubroutines provide one approach, as well as routines with common “tail” sequences sharing code by jumps.

More sharing can be done with hardwired control assistance. For example, many microarchitectures allow bits of the instruction register to specify the correct register. Another common assist is to use portions of the instruction register to specify the ALU operation. Each of these assists is under microprogrammed control and is invoked with a special value in the appropriate field. The 8086 uses both techniques, giving one 4-line routine responsibility for 32 opcodes. The drawback of adding hardwired control is it may stretch the development cycle because it no longer involves programming, but requires hardware layout for designing and debugging.

This section and the previous two give techniques for reducing cost. The following sections present three techniques for improving performance.

Reducing CPI with Special Case Microcode

As we have noted, the wise microprogrammer knows when to save space and when to spend it. An instance of this is dedicating extra microcode for frequent instructions, thereby reducing CPI. For example, the VAX 8800 uses its large control store for many versions of the CALLS instruction, optimized for register saving depending upon the value in the register-save mask. Candidates for special case microcode can be uncovered by instruction mix measurements, such as those found in Chapter 4 or in Appendix B, or by counting the frequency of use of each microinstruction in an existing implementation (see Emer and Clark [1984]).

Reducing CPI by Adding Hardwired Control

Adding hardwired control can reduce costs as well as improve performance. For example, VAX operands can be in memory or registers, but later machines reduce CPI by having special code for register–register or register–memory moves and adds: `ADDL2 Rn, 10 (Rm)` takes five or more cycles on the 780, but as few as one on the 8600. Another example is in the memory interface, where the straightforward solution is for microcode to continuously test and branch until memory is ready. Because of the delay between the time a condition becomes true and the time the next microinstruction is read, this approach can add one extra clock to each memory access. The importance of the memory interface is underlined by the 780 and 8800 statistics—20% of the 780 clock cycles and 23% of the 8800 are waiting for memory to be ready, these are called

stalls. A *stall* is where an instruction must pause one or more clock cycles waiting for some resource to be available. In this chapter stalls occur only when waiting for memory; in the next chapter we'll see other reasons for stalls.

Many machines approach this problem by having the hardware stall a microinstruction that tries to access the memory-data register before the memory operation is completed. (This can be accomplished by freezing the microinstruction address so that the same microinstruction is executed until the condition is met.) The instant the memory reference is ready, the microinstruction that needs the data is allowed to complete, avoiding the extra clock delay to access control memory.

Reducing CPI by Parallelism

Sometimes CPI can be reduced with more operations per microinstruction. This technique, which usually requires a wider microinstruction, increases parallelism with more datapath operations. It is another characteristic of machines labeled horizontal. Examples of this performance gain can be seen in the fact that the fastest models of each family in Figure 5.8 also have the widest microinstructions. Making the microinstruction wider does not guarantee increased performance, however. An example where the potential gain was not realized is found in a microprocessor very similar to the 8086, except that another bus was added to the datapath, requiring six more bits in its microinstruction. This could have reduced the execution phase from three clock cycles to two for many popular 8086 instructions. Unfortunately, these popular macroinstructions were grouped with macroinstructions that couldn't take advantage of this optimization, so they all had to run at the slower rate.

5.6

Interrupts and Other Entanglements

Control is the hard part of processor design, and the hard part of control is *interrupts*—events other than branches that change the normal flow of instruction execution. Detecting interrupt conditions within an instruction can often be on the critical timing path of a machine, possibly affecting the clock cycle time, and thus performance. Without proper attention to interrupts during design, adding interrupts to a complicated implementation can even foul up the works so as to make the design impracticable.

Invented to detect arithmetic errors and signal real-time events, interrupts have been handed a flock of difficult duties. Here are 11 examples:

- I/O device request

- Invoking an operating system service from a user program

- Tracing instruction execution

Breakpoint (programmer-requested interrupt)
 Arithmetic overflow or underflow
 Page fault (not in main memory)
 Misaligned memory accesses (if alignment is required)
 Memory-protection violation
 Using an undefined instruction
 Hardware malfunctions
 Power failure

	IBM 360	VAX	Motorola 680x0	Intel 80x86
I/O device request	Input/output interruption	Device interrupt	Exception (Level 0...7 autovector)	Vectored interrupt
Invoking the operating system service from a user program	Supervisor call interruption	Exception (change mode supervisor trap)	Exception (unimplemented instruction)—on Macintosh	Interrupt (INT instruction)
Tracing instruction execution	NA	Exception (trace fault)	Exception (trace)	Interrupt (single-step trap)
Breakpoint	NA	Exception (breakpoint fault)	Exception (illegal instruction or breakpoint)	Interrupt (breakpoint trap)
Arithmetic overflow or underflow	Program interruption (overflow or underflow exception)	Exception (integer overflow trap or floating underflow fault)	Exception (floating-point coprocessor errors)	Interrupt (overflow trap or math unit exception)
Page fault (not in main memory)	NA (only in 370)	Exception (translation not valid fault)	Exception (memory-management unit errors)	Interrupt (page fault)
Misaligned memory accesses	Program interruption (specification exception)	NA	Exception (address error)	NA
Memory protection violations	Program interruption (protection exception)	Exception (access control violation fault)	Exception (bus error)	Interrupt (protection exception)
Using undefined instructions	Program interruption (operation exception)	Exception (opcode privileged/reserved fault)	Exception (illegal instruction or breakpoint/unimplemented instruction)	Interrupt (invalid opcode)
Hardware malfunctions	Machine-check interruption	Exception (machine-check abort)	Exception (bus error)	NA
Power failure	Machine-check interruption	Urgent interrupt	NA	Nonmaskable interrupt

FIGURE 5.9 Names of 11 interrupt classes on four computers. Every event on the IBM 360 and 80x86 is called an *interrupt*, while every event on the 680x0 is called an *exception*. VAX divides events into *interrupts* or *exceptions*. Adjectives *device*, *software*, and *urgent* are used with VAX interrupts, while VAX exceptions are subdivided into *faults*, *traps*, and *aborts*.

The enlarged responsibility of interrupts has led to the confusing situation of each computer vendor inventing a different term for the same event, as Figure 5.9 on page 215 illustrates. Intel and IBM still call such events *interrupts*, but Motorola calls them *exceptions*; and, depending on the circumstances, DEC calls them *exceptions*, *faults*, *aborts*, *traps*, or *interrupts*. To give some idea of how often interrupts occur, Figure 5.10 shows the frequency on the VAX 8800.

Event	Time between events
I/O interrupt	2.7 ms
Interval timer interrupt	10.0 ms
Software interrupt	1.5 ms
Any interrupt	0.9 ms
Any hardware interrupt	2.1 ms

FIGURE 5.10 Frequency of different interrupts on the VAX 8800 running a multiuser workload on the VMS timesharing system. Real-time operating systems used in embedded controllers may have a higher interrupt rate than a general-purpose timesharing system. (Collected by Clark, Bannon, and Keller [1988].)

Clearly, there is no consistent convention for naming these events. Rather than imposing one, then, let's review the reasons for the different names. The events can be characterized on five independent axes:

1. Synchronous versus asynchronous. If the event occurs at the same place every time the program is executed with the same data and memory allocation, the event is synchronous. With the exception of hardware malfunctions, asynchronous events are caused by devices external to the processor and memory.
2. User request versus coerced. If the user task directly asks for it, it is a user-request event.
3. User maskable versus user nonmaskable. If it can be masked or disabled by a user task, the event is user maskable.
4. Within versus between instructions. This classification depends on whether the event prevents instruction completion by occurring in the middle of execution—no matter how short—or whether it is recognized between instructions.
5. Resume versus terminate. If the program's execution stops after the interrupt, it is a terminating event.

The difficult task is implementing interrupts occurring within instructions where the instruction must be resumed. Another program must be invoked to collect the state of the program, correct the cause of an interrupt, and then restore the state of the program before an instruction can be tried again.

Figure 5.11 classifies the examples from Figure 5.9 according to these five categories.

	Synchronous vs. asynchronous	User request vs. coerced	User maskable vs. nonmaskable	Within vs. between instructions	Resume vs. terminate
I/O device request	Asynchronous	Coerced	Nonmaskable	Between	Resume
Invoking operating system service	Synchronous	User request	Nonmaskable	Between	Resume
Tracing instruction execution	Synchronous	User request	User maskable	Between	Resume
Breakpoint	Synchronous	User request	User maskable	Between	Resume
Integer arithmetic overflow	Synchronous	Coerced	User maskable	Within	Terminate
Floating-point arithmetic overflow or underflow	Synchronous	Coerced	User maskable	Within	Resume
Page fault	Synchronous	Coerced	Nonmaskable	Within	Resume
Misaligned memory accesses	Synchronous	Coerced	User maskable	Within	Terminate
Memory-protection violations	Synchronous	Coerced	Nonmaskable	Within	Terminate
Using undefined instructions	Synchronous	Coerced	Nonmaskable	Within	Terminate
Hardware malfunctions	Asynchronous	Coerced	Nonmaskable	Within	Terminate
Power failure	Asynchronous	Coerced	Nonmaskable	Within	Terminate

FIGURE 5.11 The events of Figure 5.9 classified using five categories.

How Control Checks for Interrupts

Integrating interrupts with control means modifying the finite-state diagram to check for interrupts. Interrupts that occur between instructions are checked either at the beginning of the finite-state diagram—before an instruction is decoded—or at the end—after the execution of an instruction is completed. Interrupts that can occur within an instruction are generally detected in the state that causes the action or in a state that follows it. For example, Figure 5.12 shows Figure 5.4 (page 207) modified to check for interrupts.

We assume DLX transfers the return address into a new programmer-visible register, the interrupt return-address register. Control then loads PC with the address of the interrupt routine for that interrupt.

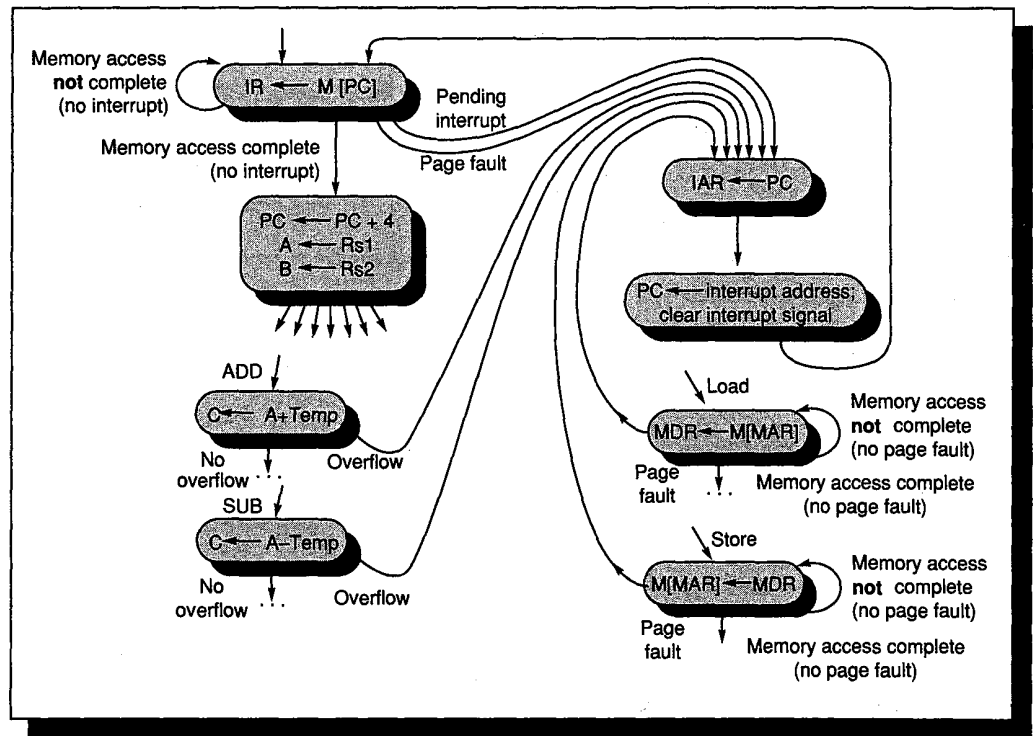


FIGURE 5.12 The top-level view of the DLX finite-state diagram (Figure 5.4 on page 207) modified to check for interrupts. Either a between interrupt or an instruction page fault invokes the control that saves the PC and then loads it with the address of the appropriate interrupt routine. The lower portion of the figure shows interrupts resulting in page faults of data accesses or arithmetic overflow.

What's Hard About Interrupts

The conflicting terminology is confusing, but that is not what makes the hard part of control hard. Even though interrupts are rare, the hardware must be designed so that the full state of the machine can be saved, including an indication of the offending event, and the PC of the instruction to be executed after the interrupt is serviced. This difficulty is exacerbated by events occurring during the middle of execution, for many instructions also require the hardware to restore the machine to the state just before the event occurred—the beginning of the instruction. This last requirement is so difficult that computers are awarded the title *restartable* if they pass that test. That supercomputers and many early microprocessors do not earn that badge of honor illustrates both the difficulty of interrupts and the potential cost in hardware complexity and execution speed.

No engineers deserve more admiration than those who built the first VAX, DEC's first restartable minicomputer. The variable-length instructions mean the computer can fetch 50 bytes of one instruction before discovering that the next

byte of the instruction is not in main memory—a situation that requires the saved PC to point 50 bytes earlier. Imagine the difficulties of restarting an instruction with six operands, each of which could be misaligned and thus be partially in memory and partially on disk!

The instructions that are hardest to restart are those that modify some of the machine state before it is known whether interrupts can occur. The VAX autoincrement and autodecrement addressing modes would naturally modify registers during the addressing phase of execution rather than at the writeback phase, and so would be vulnerable to this difficulty. To avoid this problem, recent VAXes keep a history queue of the register specifiers and the operations on the registers, so that the operations can be reversed on an interrupt. Another approach, used on the earlier VAXes, is to record the specifiers and the original values of the registers, restoring the original values on interrupt. (The primary difference is that it only takes a few bits to record how the address was changed due to autoincrement or autodecrement versus the full 32-bit register value.)

It is not just addressing modes that make the VAX difficult to restart; long-running instructions mean that interrupts must be checked in the middle of execution to prevent long interrupt latency. `MOV3`, for example, copies up to 2^{16} bytes and can take tens of milliseconds to finish—far too long to wait for an urgent event. On the other hand, even if there were a way to undo copying in the middle of execution so that `MOV3` could be restarted, interrupts would occur so frequently, relative to this long-running instruction (see Figure 5.10 on page 216), that `MOV3` would be restarted repeatedly under those conditions. Such wasted effort from incomplete copies would render `MOV3` worse than useless.

DEC divided the problem to conquer it. First, the operands—source address, length, and destination address—are fetched from memory and placed into general-purpose registers R1, R2, and R3. If an interrupt occurs during this first phase, these registers are restored, and the `MOV3` is restarted from scratch. After this first phase, every time a byte is copied, the length (R2) is decremented and addresses (R1 and R3) are incremented. If an interrupt occurs during this second phase, `MOV3` sets the *first part done* (FPD) bit in the program status word. When the interrupt is serviced and the instruction is reexecuted, it first checks the FPD bit to see if the operands have already been placed in registers. If so, the VAX doesn't fetch the address and length operands, but just continues with the current values in the registers, since that is all that remains to be copied. This permits more rapid response to interrupts while allowing long-running instructions to make progress between interrupts.

IBM had a similar problem. The 360 included the `MVC` instruction, which copies up to 256 bytes of data. For the early machines without virtual memory, the machine simply waited until the instruction was completed before servicing interrupts. With the inclusion of virtual memory in the 370, the problem could no longer be ignored. Control first tries to access all possible pages, forcing all possible virtual memory miss interrupts to occur before moving any data. If any interrupts occur in this phase, the instruction is restarted. Control then ignores interrupts until the instruction is complete. To allow longer copies, the 370

includes MVCL, which can move up to 2^{24} bytes. The operands are in registers and are updated as a part of execution—like the VAX, except that there is no need for FPD since the operands are always in registers. (Or, to speak historically, the VAX solution is like the IBM 370, which came first.)

5.7

Putting It All Together: Control for DLX

The control for DLX is presented here to tie together the ideas from the previous three sections. We begin with a finite-state diagram to represent hardwired control and end with microprogrammed control. Both versions of DLX control are used to demonstrate tradeoffs to reduce cost or to improve performance. Because the figures are already too large, the checking for data page faults or arithmetic overflow shown in Figure 5.12 (page 218) is not included in this section. (Exercise 5.12 adds them.)

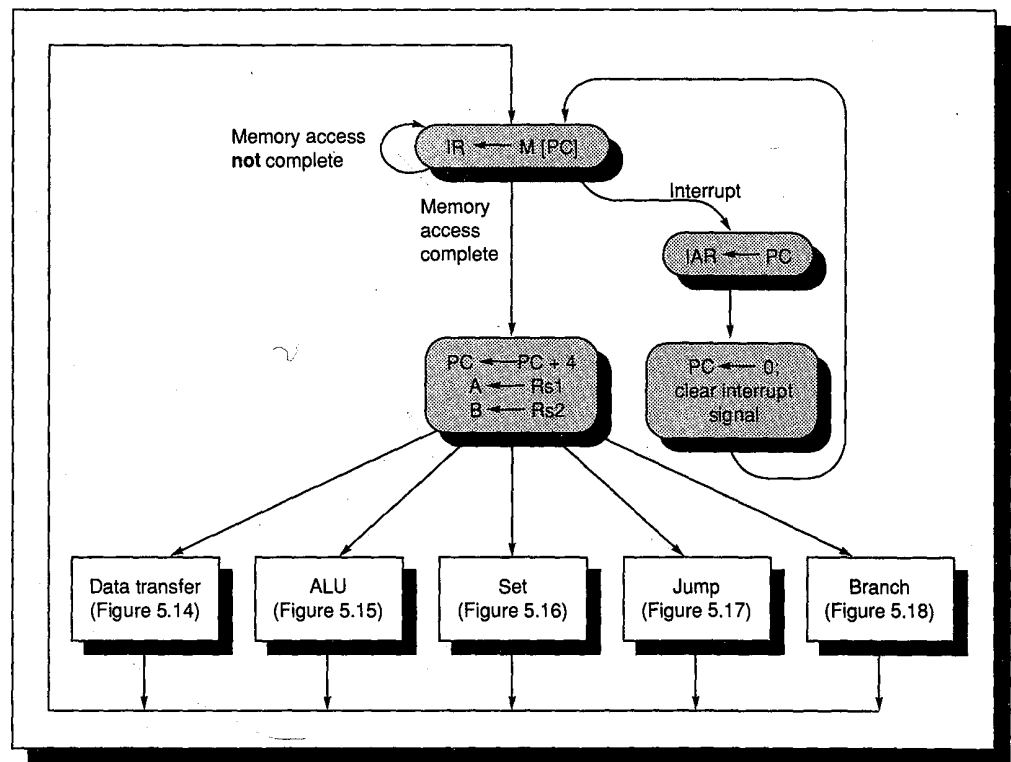


FIGURE 5.13 The top-level view of the DLX finite-state diagram for the non-floating-point instructions. The first two steps of instruction execution—instruction fetch and instruction decode/register fetch—are shown. The first state repeats until the instruction is fetched from memory or an interrupt is detected. If an interrupt is detected, the PC is saved in IAR and PC is set to the address of the interrupt routine. The last three steps of instruction execution—execution/effective address, memory access, and write back—are shown in Figures 5.14 to 5.18 on pages 221–224.

Rather than trying to draw the DLX finite-state machine in a single figure showing all 52 states, Figure 5.13 (see page 220) shows just the top level, containing 4 states plus references to the rest of the states detailed in Figures 5.14 (below) through 5.18 (page 224). Unlike Figure 5.2 (page 205), Figure 5.13 takes advantage of the change to the datapath allowing PC to address memory directly without going through MAR (Figure 5.4 on page 207).

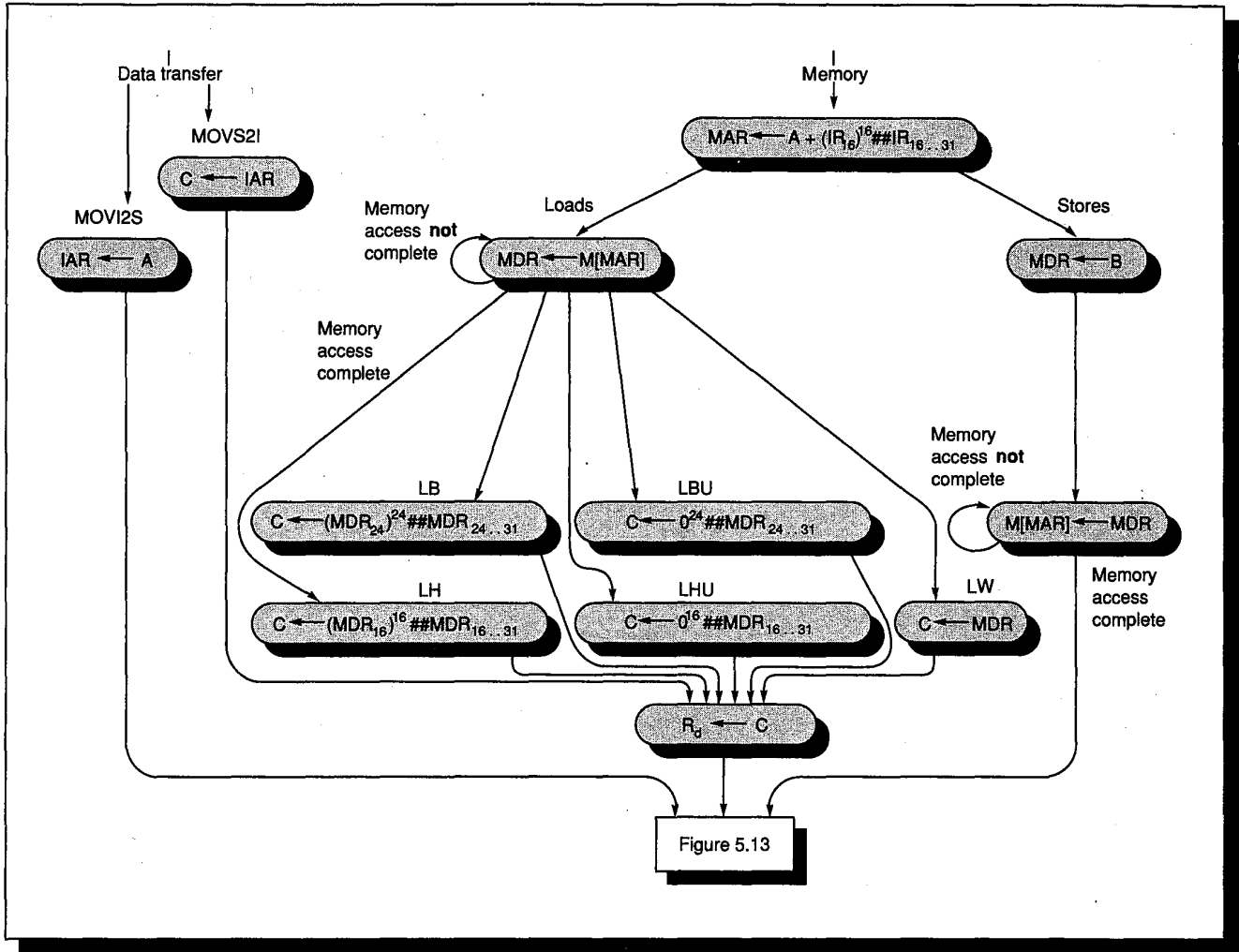


FIGURE 5.14 The effective address calculation, memory-access, and write-back states for the memory-access and data-transfer instructions of DLX. For loads, the second state repeats until the data is fetched from memory. The final state of stores repeats until the write is complete. While the operation of all five loads is shown in the states of this figure, the proper operation of writes depends on the memory system writing bytes and halfwords, without disturbing the rest of the word in memory, and correctly aligning the bytes and halfwords (see Figure 3.10, page 97) over the proper bytes of memory. On completion of execution control transfers to Figure 5.13, found on page 220.

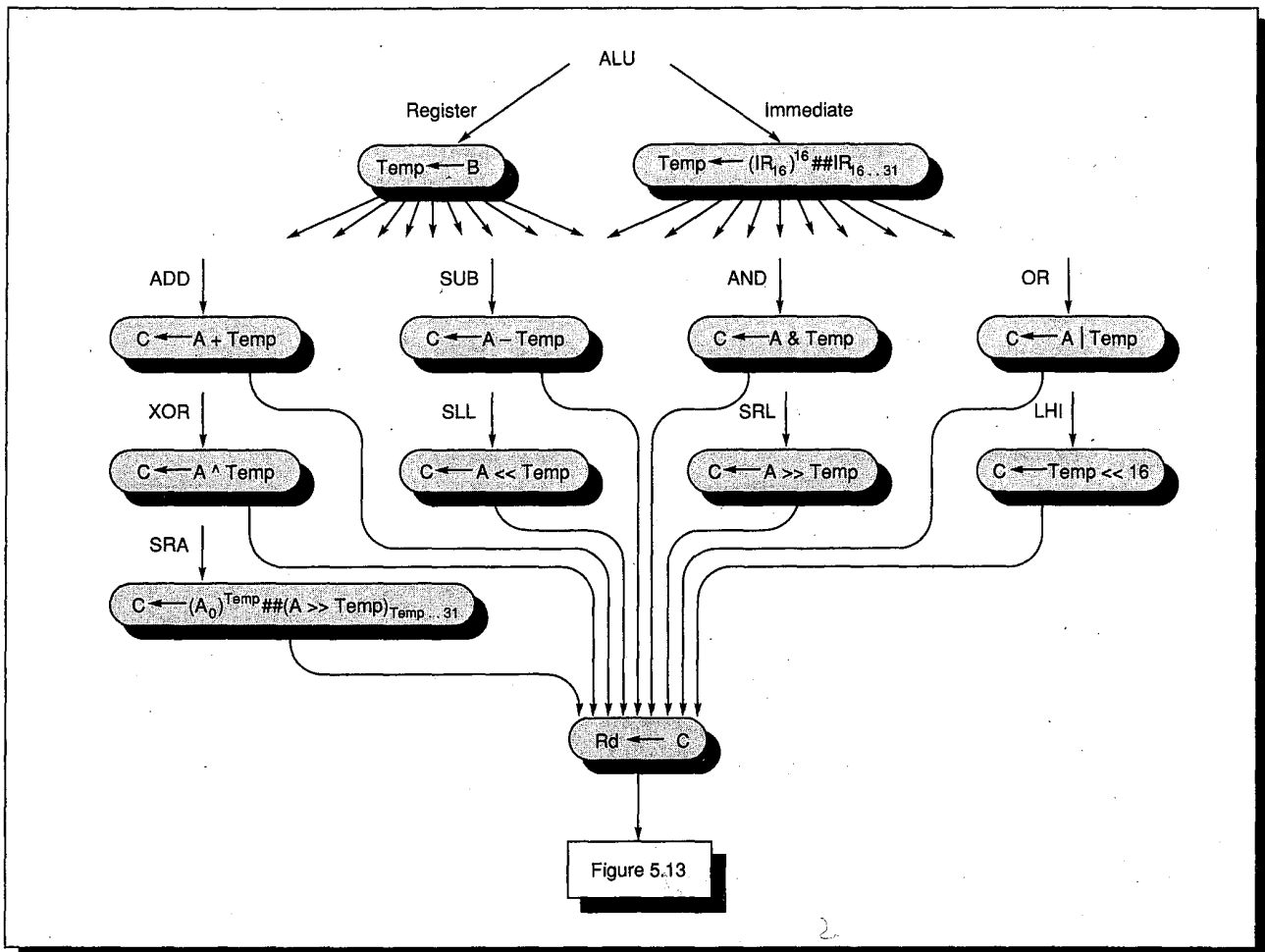


FIGURE 5.15 The execution and write-back states for the ALU instructions of DLX. After putting a register or the sign-extended 16-bit immediate into Temp, 1 of the 9 instructions is executed, and the result (C) is written back into the register file. Only SRA and LHI may not be self-explanatory: The SRA instruction shifts right while it sign extends the operand and LHI loads the upper 16 bits of the register while zeroing the lower 16 bits. (The C operators << and >> shift left and right, respectively; they fill with zeros unless bits are concatenated explicitly using ##, e.g., sign extension). As mentioned above, the check for overflow in ADD and SUB is not included to simplify the figure. On completion of execution control transfers to Figure 5.13 (page 220).

FIGURE 5.16 (See adjoining page.) The execution and write-back states for the Set instructions of DLX. After putting a register or the sign-extended 16-bit immediate into Temp, 1 of the 6 instructions compares A to Temp and then sets C to 1 or 0, depending on whether the condition is true or false. C is then written back into the register file, and then execution control transfers to Figure 5.13 (page 220). The dashed lines in this figure and Figure 5.18 are used to make it easier to follow intersecting lines.

FIGURE 5.17 (See adjoining page.) The execution and write-back states for the jump instructions of DLX. With jump and link instructions, the return address is first placed in C before the new value is loaded into PC. Trap saves it in IAR. Note that the immediate in these instructions is 10 bits longer than the 16-bit immediate in all other instructions. Jump and link instructions conclude by writing the return address back into R31. On completion of execution, control transfers to Figure 5.13 (page 220).

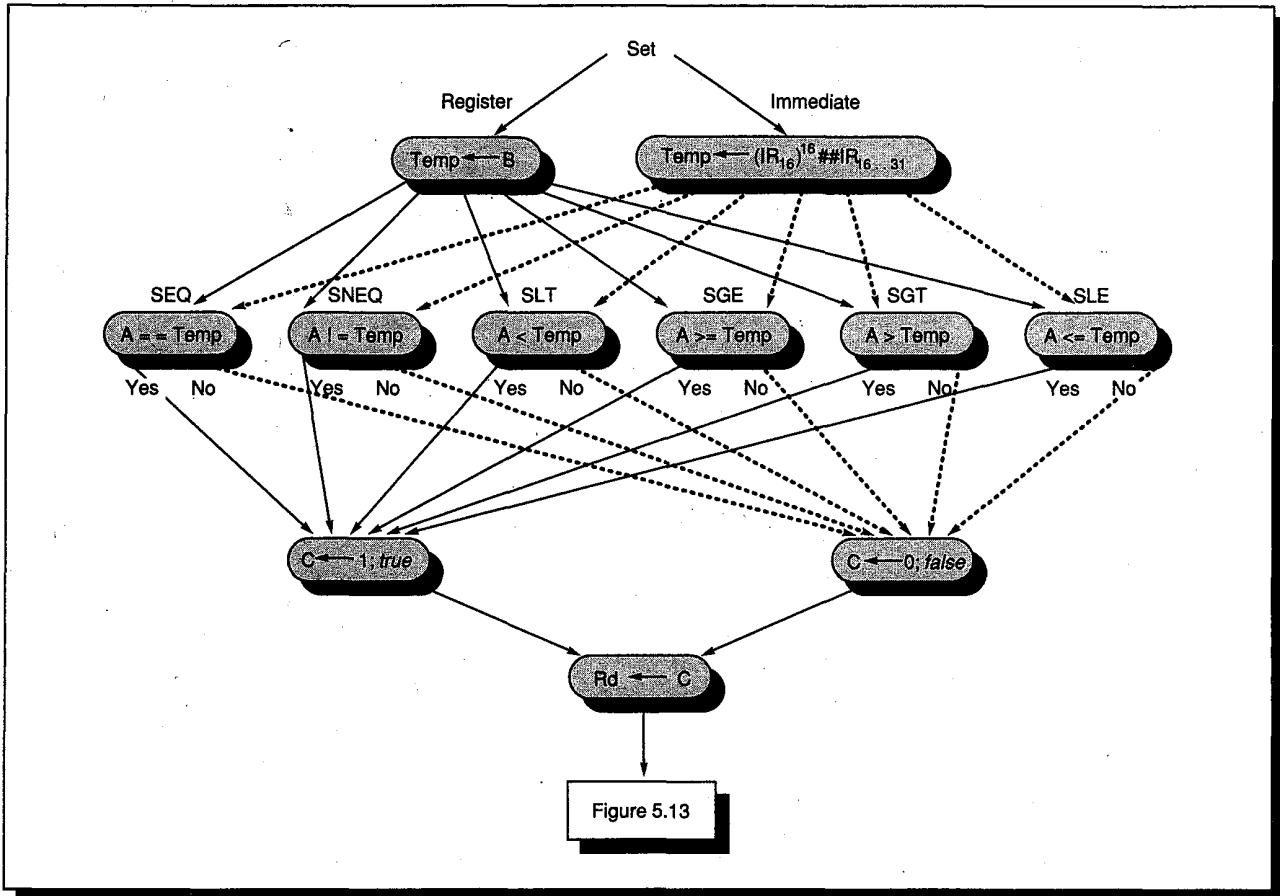


Figure 5.13

FIGURE 5.16

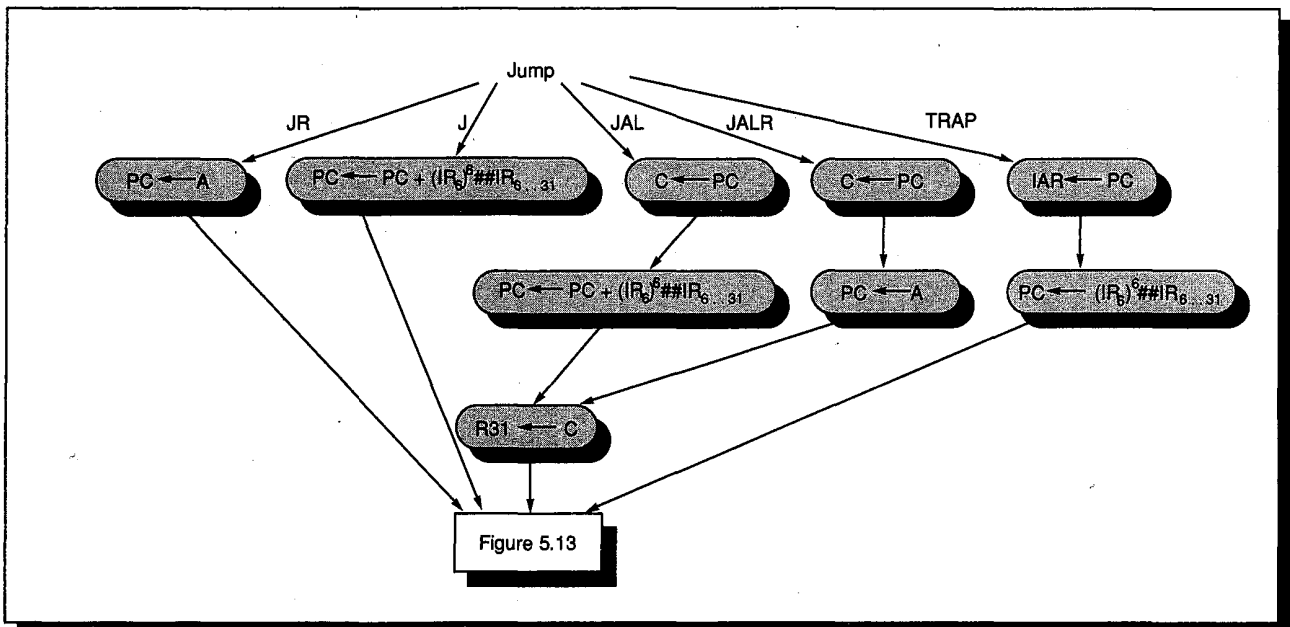


Figure 5.13

FIGURE 5.17

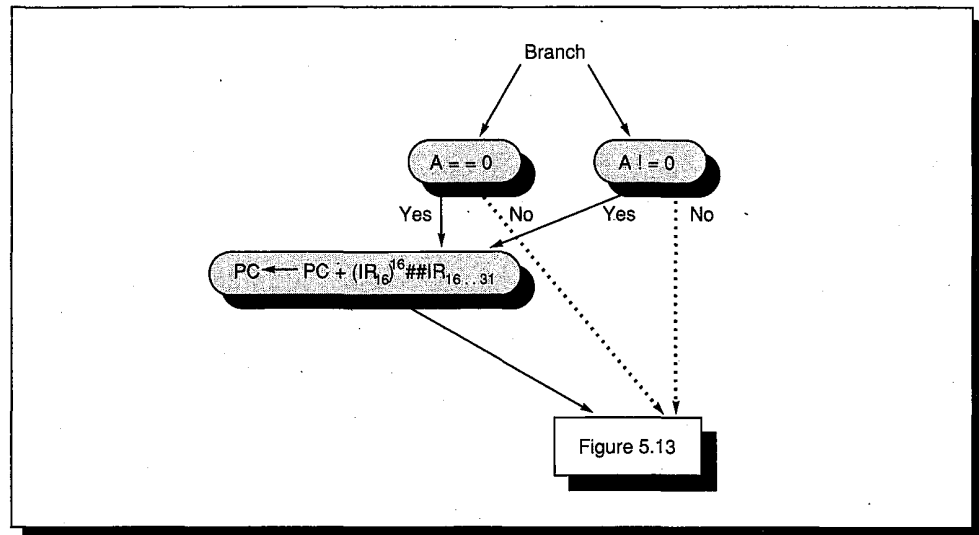


FIGURE 5.18 The execution states for the branch instructions of DLX. The PC is loaded with the sum of the PC and the immediate only if the condition is true. On completion of execution, control transfers to Figure 5.13, found on page 220.

Performance of Hardwired Control for DLX

As stated in Section 5.4, the goal for control designs is to minimize CPI, clock cycle time, amount of control hardware, and development time. CPI is just the average number of states along the execution path of an instruction.

Example

Let's assume that hardwired control directly implements the finite-state diagram in Figures 5.13 to 5.18. What is the CPI for DLX running GCC?

Answer

The number of clock cycles to execute each DLX instruction is determined by simply counting the states of an instruction. Starting at the top, every instruction spends at least two clock cycles in the states in Figure 5.13 (ignoring interrupts). The actual number depends on the average number of times the state accessing memory must repeat because memory is not ready. (These wasted clock cycles are usually called *memory stall cycles* or *wait states*.) In cache-based machines this value is typically 0 (i.e., no repetitions since cache access is 1 cycle) when the data is found in the cache, and 10 or higher when it is not.

The time for the remaining portion of instruction execution comes from the additional figures. Besides two cycles for fetch and decode, loads take four more cycles plus clock cycles waiting for the data access, while stores take just three more clock cycles plus wait states. Three extra clock cycles are also needed by ALU instructions, and set instructions take four. Figure 5.17 shows that jumps take just one extra clock cycle with jump and links taking three. Branches depend on the result: Taken branches use two more clock cycles while

DLX instructions	Minimum clock cycles	Memory accesses	Total clock cycles
Loads	6	2	8
Stores	5	2	7
ALU	5	1	6
Set	6	1	7
Jumps	3	1	4
Jump and links	5	1	6
Branch (taken)	4	1	5
Branch (not taken)	3	1	4

FIGURE 5.19 Clock cycles per instruction for DLX categories using the state diagram in Figures 5.13 through 5.18. Determining the total clock cycles per category requires multiplying the number of memory accesses—including instruction fetches—times the average number of wait states, and adding this product to the minimum number of clock cycles. We assume an average of 1 clock cycle per memory access. For example, loads take eight clock cycles if the average number of wait states is one.

untaken branches need just one. Adding these times to the first portion of instruction execution yields the clock cycles per DLX instruction class shown in Figure 5.19.

From Chapter 2, one way to calculate CPI is

$$\text{CPI} = \sum_{i=1}^n \left(\text{CPI}_i * \frac{I_i}{\text{Instruction count}} \right)$$

Using the DLX instruction mix from Figure C.4 in Appendix C for GCC (normalized to 100%), the percentage of taken branches from Figure 3.22 (page 107), and one for the average number of wait states per memory access, the DLX CPI for this datapath and state diagram is calculated:

Loads	8 * 21%	=	1.68
Stores	7 * 12%	=	0.84
ALU	6 * 37%	=	2.22
Set	7 * 6%	=	0.42
Jumps	4 * 2%	=	0.08
Jump and links	6 * 0%	=	0.00
Branch (taken)	5 * 12%	=	0.60
Branch (not taken)	4 * 11%	=	0.44
	Total CPI:		6.28

Thus, the DLX CPI for GCC is about 6.3.

Improving DLX Performance When Control Is Hardwired

As mentioned above, performance is improved by reducing the number of states an instruction must pass through during execution. Sometimes, performance can be improved by removing intermediate calculations that select one of several options, either by adding hardware that uses information in the opcode to later select the appropriate option, or by simply increasing the number of states.

Example

Let's look at improving the performance of ALU instructions by removing the top two states in Figure 5.15 on page 222, which load either a register or an immediate into Temp. One approach uses a new hardware option. Let's call it "X" (see Figure 5.20). The X option selects either the B register or the 16-bit immediate, depending on the opcode in IR. A second approach is simply to increase the number of execution states so that there are separate states for ALU instructions using immediate versus ALU instructions using registers.

Answer

For each option, what would be the change in performance, and how should the state diagram be changed? Also, how many states are needed in each option?

Either change reduces ALU execution time from five to four clock cycles plus wait states. From Figure C.4, ALU operations are about 37% of the instructions for GCC, lowering CPI from 6.3 to 5.9, and making the machine about 7% faster. Figure 5.20 shows Figure 5.15 modified to use the X option instead of the two states that load Temp, while Figure 5.21 simply has many more states to achieve the same result. The total number of states are 50 and 58, respectively.

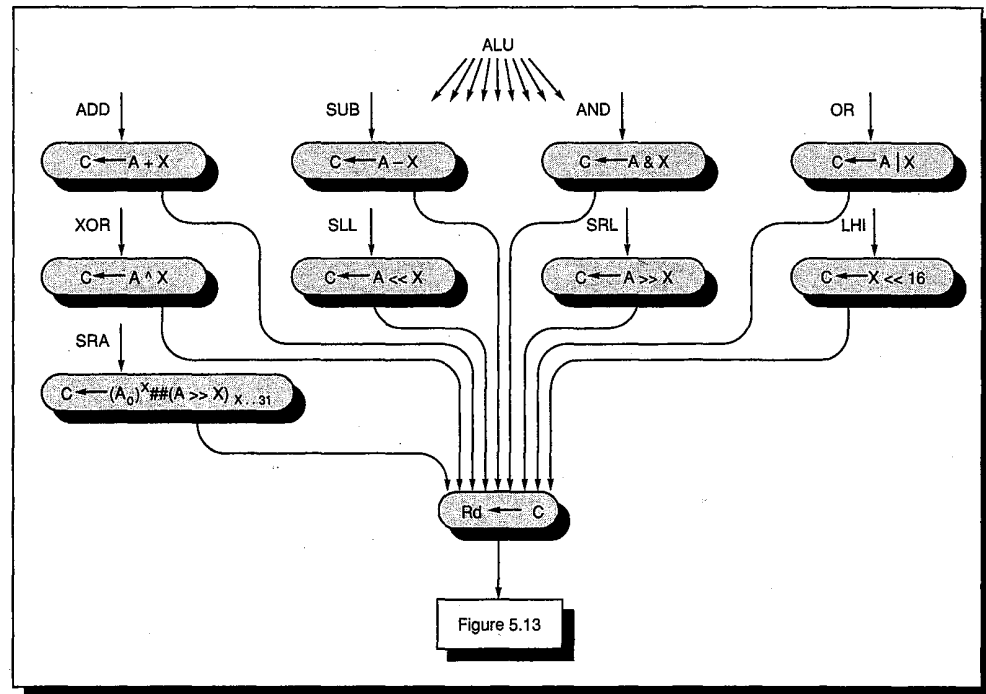


FIGURE 5.20 Figure 5.15 modified to remove the two states loading Temp. The states use the new X option to mean that either B or $(IR_{16})^{16}##IR_{16..31}$ is the operand, depending on the DLX opcode.

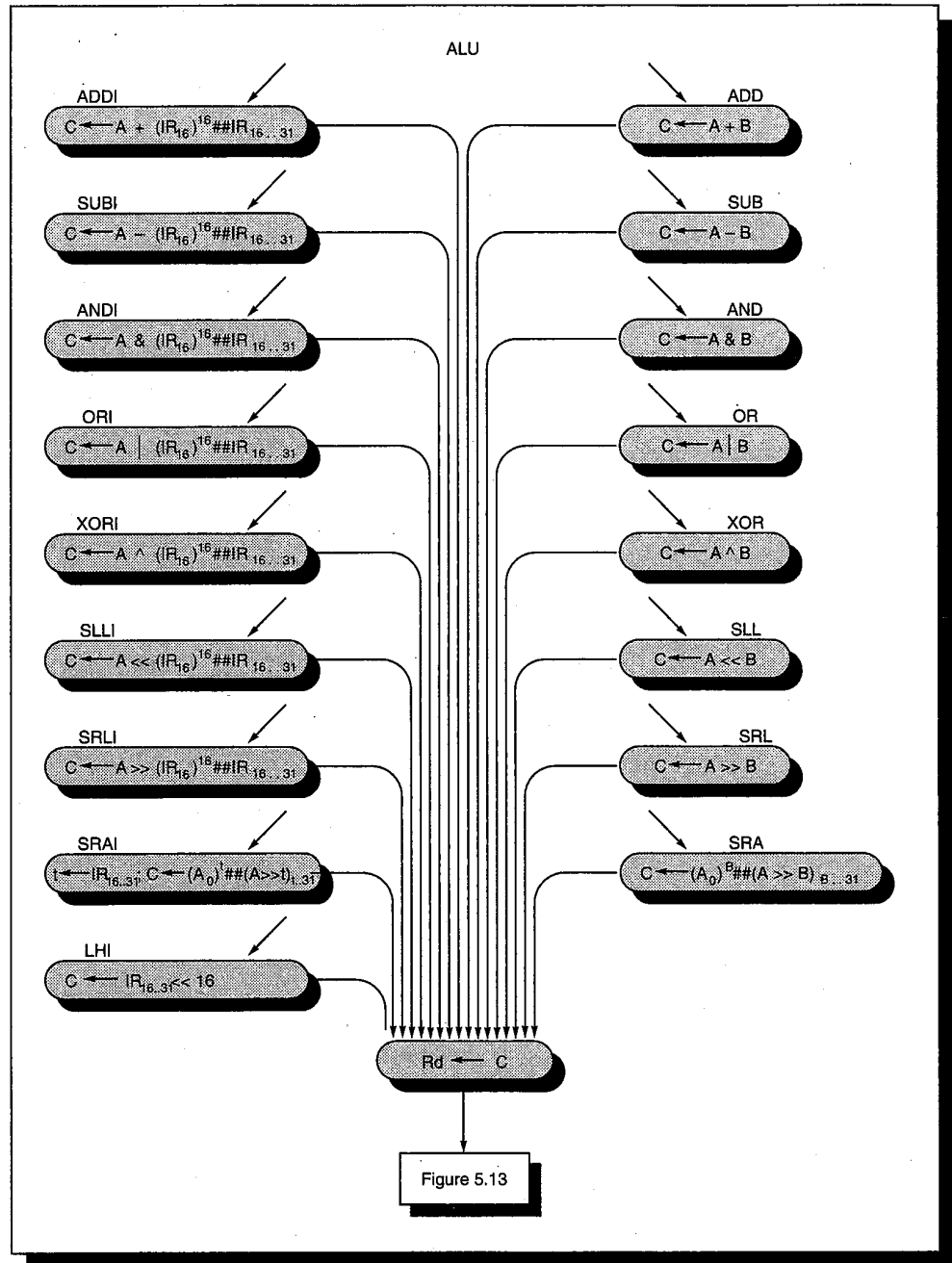


FIGURE 5.21 Figure 5.15 modified to remove the two states loading Temp. Unlike Figure 5.20, this requires no new hardware options in the datapath, but simply more control states.

Control can affect the clock cycle time, either because control itself takes longer than the corresponding operations in the datapath, or because the datapath operations selected by control lengthens the worst-case clock cycle time.

Example

Assume a machine with a 10-ns clock cycle (100-MHz clock rate). Suppose that on closer inspection the designer discovered that all states could be executed in 9 ns, except states that use the shifter. Would it be wise to split those states, taking two 9-ns clock cycles for shift states and one 9-ns clock for everything else?

Answer

Assuming the improvement in the previous example, the average instruction execution time for the 100-MHz machine is 5.9×10 ns or 59 ns. The shifter is only used in the states of four instructions: SLL, SRL, SRA, and LHI (see Figure 5.20). In fact, each of these instructions takes 5 clock cycles (including one wait state for memory access), and only one of the five original clock cycles need be split into two new clock cycles. Thus, the average execution time of these instructions changes from 5×10 ns, or 50 ns, to 6×9 ns, or 54 ns. From Figure C.4 these 4 instructions are about 11% of the instructions executed for GCC (after normalization), making the average instruction execution time $89\% \times (5.9 \times 9 \text{ ns}) + 11\% \times 54 \text{ ns}$ or 53 ns. Thus, splitting the shift state results in a machine that is about 10% faster—a wise decision. (See Exercise 5.8 for a more sophisticated version of this tradeoff.)

Hardwired control is completed by listing the control signals activated in each state, assigning numbers to the states, and finally generating the PLA. Now let's implement control using microcode in a ROM.

Microcoded Control for DLX

A custom format such as this is a slave to the architecture of the hardware and instruction set which it serves. The format must strike a proper compromise between ROM size, ROM-output decoding circuitry size, and machine execution rate.

Jim McKeiv et al. [1977]

Before microprogramming can commence, the microinstruction set must be determined. The first step is to list the possible entries for each field of the DLX microinstruction format from Figure 5.6 on page 209. Figure 5.7 on page 211 lists them for the Destination, Source1, and Source2 fields. Figure 5.22 below shows the values for the remaining fields.

Sequencing of microinstructions requires further explanation. The microprogrammed control includes a microprogram counter to specify the address of the next microinstruction if a branch is not taken, as in Figure 5.5 on page 208. In addition to the branches using the Jump address field, three tables are used to decode the DLX macroinstructions. These tables are indexed with the opcodes of the DLX instruction, and supply a microprogram address depending on the value in the opcode. Their use will become clear as we examine the DLX microprogram.

Value	ALU	Misc	Cond
0	ADD +	Instr Read $IR \leftarrow M[PC]$	--- <i>Go to next sequential microinstruction</i>
1	SUB -	Data Read $MDR \leftarrow M[MAR]$	Uncond <i>Always jump</i>
2	RSUB $-_r$ (reverse sub)	Write $M[MAR] \leftarrow MDR$	Int? <i>Pending (between instruction) interrupt?</i>
3	AND &	AB \leftarrow RF <i>Load A&B from Reg. File</i>	Mem? <i>Memory access not complete?</i>
4	OR /	Rd \leftarrow C <i>Write Rd</i>	Zero? <i>Is the ALU output zero?</i>
5	XOR ^	R31 \leftarrow C <i>Write R31 (for call)</i>	Negative? <i>Is the ALU output less than zero?</i>
6	SLL <<		Load? <i>Is the macroinstruction a DLX load?</i>
7	SRL >>		Decode1 <i>Address table 1 determines next microinstruction (uses main opcode)</i>
8	SRA >> _a		Decode2 <i>Address table 2 determines next microinstruction (uses "func" opcode)</i>
9	Pass S1 S1		Decode3 <i>Address table 3 determines next microinstruction (uses main opcode)</i>
10	Pass S2 S2		

FIGURE 5.22 The options for three fields of the DLX microinstruction format in Figure 5.6 on page 209. The possible names are shown on the left of the field name, with an explanation of each field to the right. The real microinstruction would contain a bit pattern corresponding to the number in the first column. Combined with Figure 5.7 (page 211), all the fields are defined except the Constant and Jump address fields, which contain numbers supplied by the microprogrammer. >>_a is an abbreviation for shift right arithmetic and $-_r$ means reverse subtract ($B -_r A = A - B$).

Following the lead of the state diagram, the DLX microprogram is divided into Figures 5.23, 5.25, 5.27, 5.28, and 5.29, with each section of microcode corresponding to one of Figures 5.13 to 5.18 (pages 220–224). The first state in Figure 5.13 becomes the first two microinstructions in Figure 5.23. The first microinstruction (address 0) branches to microinstruction 3 if there is an interrupt pending. Microinstruction 1 fetches an instruction from memory, branching back to itself as long as the memory access is not complete. Microinstruction 2 increments the PC by 4, loads A and B, and then does the first-level decoding. The address of the next microinstruction then depends on which macroinstruction is in the instruction register. The microinstruction addresses for this first-level macroinstruction decode are specified in Figure 5.24. (In reality, the table shown in this figure is specified after the microprogram is written, as both the number of entries and the corresponding locations aren't known until then.)

Loc	Label	Dest	ALU	S1	S2	C	Misc	Cond	Jump label	Comment
0	Ifetch:							Interrupt?	Intrpt	Check interrupt
1	Iloop:						Instr Read	Mem?	Iloop	$IR \leftarrow M[PC]$; wait for memory
2		PC	ADD	PC	Constant	4	AB←RF	Decode1		
3	Intrpt:	IAR	Pass S1	PC						Interrupt
4		PC	Pass S2		Constant	0		Uncond	Ifetch	$PC \leftarrow 0$ & go fetch next instruction

FIGURE 5.23 The first section of the DLX microprogram, corresponding to the states in Figure 5.13 (page 220). The first column contains the absolute address of the microinstruction, followed by a label. The rest of the fields contain values from Figures 5.7 (page 211) and 5.22 for the microinstruction format in Figure 5.6 (page 209). As an example, microinstruction 2 corresponds to the second state of Figure 5.13. It sends the output from the ALU into PC, tells the ALU to add, puts PC onto the Source1 bus, and a constant from the microinstruction (whose value is 4) onto the Source2 bus. In addition, A and B are loaded from the register file according to the specifiers in IR. Finally, the address of the next microinstruction to be executed comes from decode table 1 (Figure 5.24), which depends on the opcode in the instruction register (IR).

Opcodes (symbolically specified)	Absolute address	Label	Figure
Memory	5	Mem:	5.25
Move to special	20	MovI2S:	5.25
Move from special	21	MovS2I:	5.25
S2 = B	23	Reg:	5.27
S2 = Immediate	24	Imm:	5.27
Branch equal zero	50	Beq:	5.29
Branch not equal zero	52	Bne:	5.29
Jump	54	Jump:	5.29
Jump register	55	JReg:	5.29
Jump and link	56	JAL:	5.29
Jump and link register	58	JALR:	5.29
Trap	60	Trap:	5.29

FIGURE 5.24 Opcodes and corresponding addresses for decode table 1. The opcodes are shown symbolically on the left, followed by the addresses with the absolute microinstruction address, a label, and the figure where the microcode can be found. If this table were implemented with a ROM it would contain 64 entries corresponding to the 6-bit opcode of DLX. As this would clearly result in many redundant or unspecified entries, a PLA could be used to minimize hardware.

Figure 5.25 contains the DLX load and store instructions. Microinstruction 5 calculates the effective address, and branches to microinstruction 9 if the

macroinstruction in the IR is a load. If not, microinstruction 6 loads MDR with the value to be stored, and microinstruction 7 jumps to itself until the memory is finished writing the data. Microinstruction 8 then jumps back to microinstruction 0 (Figure 5.23) to begin the execution cycle all over again. If the macroinstruction was a load, microinstruction 9 loops until the data has been read. Microinstruction 10 then uses decode table 2 (specified in Figure 5.26) to specify the address of the next microinstruction. Unlike the first decode table, this table is used by other microinstructions. (There is no conflict in multiple uses since the opcodes for each instance are different.)

Suppose the instruction were load halfword. Figure 5.26 shows that the result of decode 2 would be to jump to microinstruction 15. This microinstruction shifts the contents of MDR to the left 16 bits and stores the result in Temp. The following microinstruction shifts Temp right arithmetically 16 bits and puts the result in C. C now contains the 16 rightmost bits of MDR, with the upper 16 bits containing the extended sign. This microinstruction jumps to location 22, which writes C back into the destination register specifier in IR, and then jumps to fetch the next macroinstruction starting at location 0 (Figure 5.23).

Loc	Label	Dest	ALU	S1	S2	C	Misc	Cond	Jump label	Comment
5	Mem:	MAR	ADD	A	imm16			Load?	Load	Memory instruct.
6	Store:	MDR	Pass S2		B					Store
7	Dloop:						Data write	Mem?	Dloop	
8								Uncond	Ifetch	Fetch next instr.
9	Load:						Data read	Mem?	Load	Load MDR
10								Decode2		
11	LB:	Temp	SLL	MDR	Constant	24				Load byte; shift left to remove upper 24 bits
12		C	SRA	Temp	Constant	24		Uncond	Write1	Shift right arithmetic to sign extend
13	LBU:	Temp	SLL	MDR	Constant	24				LB unsigned
14		C	SRL	Temp	Constant	24		Uncond	Write1	Shift right logical
15	LH:	Temp	SLL	MDR	Constant	16				Load half
16		C	SRA	Temp	Constant	16		Uncond	Write1	Shift right arithmetic
17	LHU:	Temp	SLL	MDR	Constant	16				LH Unsigned
18		C	SRL	Temp	Constant	16		Uncond	Write1	Shift right logical
19	LW:	C	Pass S1	MDR				Uncond	Write1	Load word
20	MovI2S:	IAR	Pass S1	A				Uncond	Ifetch	Move to special
21	MovS2I:	C	Pass S1	IAR						Move from spec.
22	Write1:						Rd←C	Uncond	Ifetch	Write back & go fetch next instruction

FIGURE 5.25 The section of the DLX microprogram for loads and stores, corresponding to the states in Figure 5.14 (page 221). The microcode for bytes and halfwords takes an extra microinstruction to align the data (see Figure 3.10, page 97). Note that microinstruction 5 loads A from Rd, just in case the instruction is a store. The label Ifetch is for microinstruction 0 in Figure 5.23 on page 230.

Opcode	Absolute address	Label	Figure
Load byte	11	LB:	5.25
Load byte unsigned	13	LBU:	5.25
Load half	15	LH:	5.25
Load half unsigned	17	LHU:	5.25
Load word	19	LW:	5.25
ADD	25	ADD/I:	5.27
SUB	26	SUB/I:	5.27
AND	27	AND/I:	5.27
OR	28	OR/I:	5.27
XOR	29	XOR/I:	5.27
SLL	30	SLL/I:	5.27
SRL	31	SRL/I:	5.27
SRA	32	SRA/I:	5.27
LHI	33	LHI:	5.27
Set equal	35	SEQ/I:	5.28
Set not equal	37	SNE/I:	5.28
Set less than	39	SLT/I:	5.28
Set greater than or equal	41	SGE/I:	5.28
Set greater than	43	SGT/I:	5.28
Set less than or equal	45	SLE/I:	5.28

FIGURE 5.26 Opcodes and corresponding addresses for decode tables 2 and 3. The opcodes are shown symbolically on the left, followed by the absolute microinstruction address, the corresponding label, and the figure where the microcode can be found. Since the opcodes are shown symbolically, and they go to the same place in both tables, the same information can be used for specifying decode tables 2 and 3. This similarity is attributable to the immediate version and register version of the DLX instructions sharing the same microcode. If a table were implemented with a ROM, it would contain 64 entries corresponding to the 6-bit opcode of DLX. Again, the many redundant or unspecified entries suggest the use of a PLA to minimize hardware cost.

The ALU instructions are found in Figure 5.27. The first two microinstructions correspond to the states at the top of Figure 5.15 (page 222). After loading Temp with either the register or the immediate, each uses a decode table to vector to the microinstruction that executes the ALU instruction. To save microcode space, the same microinstruction is used whether the operand is a register or an immediate. One of the microinstructions between 25 and 33 is executed, storing its result in C. It then jumps to microinstruction 34, which stores C into the register specified in the IR, and in turn jumps to fetch the next macroinstruction.

Loc	Label	Dest	ALU	S1	S2	C	Misc	Cond	Jump label	Comment
23	Reg:	Temp	Pass S2		B			Decode2		<i>source2 = reg</i>
24	Imm:	Temp	Pass S2		Imm			Decode3		<i>source2 = imm.</i>
25	ADD/I:	C	ADD	A	Temp			Uncond	Write2	<i>ADD</i>
26	SUB/I:	C	SUB	A	Temp			Uncond	Write2	<i>SUB</i>
27	AND/I:	C	AND	A	Temp			Uncond	Write2	<i>AND</i>
28	OR/I:	C	OR	A	Temp			Uncond	Write2	<i>OR</i>
29	XOR/I:	C	XOR	A	Temp			Uncond	Write2	<i>XOR</i>
30	SLL/I:	C	SLL	A	Temp			Uncond	Write2	<i>SLL</i>
31	SRL/I:	C	SRL	A	Temp			Uncond	Write2	<i>SRL</i>
32	SRA/I:	C	SRA	A	Temp			Uncond	Write2	<i>SRA</i>
33	LHI:	C	SLL	Temp	Constant	16		Uncond	Write2	<i>LHI</i>
34	Write2:						Rd←C	Uncond	Ifetch	<i>Write back & go fetch next instruction</i>

FIGURE 5.27 Like the first two states in Figure 5.15 (page 222), microinstructions 23 and 24 load Temp with an operand and then vector to the appropriate microinstruction, depending on the opcode in IR. One of the nine following microinstructions is executed, leaving its result in C. C is written back into the register specified in the register destination field of DLX macroinstruction in IR in microinstruction 34.

Loc	Label	Dest	ALU	S1	S2	C	Misc	Cond	Jump label	Comment
35	SEQ/I:		SUB	A	Temp			Zero?	Set1	<i>Set equal</i>
36		C	Pass S2		Constant	0		Uncond	Write4	<i>A≠T (set to false)</i>
37	SNE/I:		SUB	A	Temp			Zero?	Set0	<i>Set not equal</i>
38		C	Pass S2		Constant	1		Uncond	Write4	<i>A≠T (set to true)</i>
39	SLT/I:		SUB	A	Temp			Negative?	Set1	<i>Set less than</i>
40		C	Pass S2		Constant	0		Uncond	Write4	<i>A≥T (set to false)</i>
41	SGE/I:		SUB	A	Temp			Negative?	Set0	<i>Set GT or equal</i>
42		C	Pass S2		Constant	1		Uncond	Write4	<i>A≥T (set to true)</i>
43	SGT/I:		RSUB	A	Temp			Negative?	Set1	<i>Set greater than</i>
44		C	Pass S2		Constant	0		Uncond	Write4	<i>T≥A (set to false)</i>
45	SLE/I:		RSUB	A	Temp			Negative?	Set0	<i>Set LT or equal</i>
46		C	Pass S2		Constant	1		Uncond	Write4	<i>T≥A (set to true)</i>
47	Set0:	C	Pass S2		Constant	0		Uncond	Write4	<i>Set to 0 = false</i>
48	Set1:	C	Pass S2		Constant	1				<i>Set to 1 = true</i>
49	Write4:						Rd←C	Uncond	Ifetch	<i>Write back & fetch next instruction</i>

FIGURE 5.28 Corresponding to Figure 5.16 (pages 222–223), this microcode performs the DLX Set instructions. As in the previous figure, to save space these same microinstructions execute either the version of set using registers or the version using immediates. The tricky microcode is found in microinstructions 43 and 45, where the subtraction Temp – A is unlike the earlier microcode. Remember that $A -_r \text{Temp} = \text{Temp} - A$ (see Figure 5.22 on page 229).

Figure 5.28 corresponds to the states in Figure 5.16 (pages 222–223), except that the top two states that load Temp are microinstructions 23 and 24 of the previous figure; the decode tables will either jump to locations 25 to 34 in Figure 5.27, or 35 to 45 in Figure 5.28, depending on the opcode. The microinstructions for Set perform relative tests by having the ALU subtract Temp from A and then test the ALU output to see if the result is zero or negative. Depending on the test result, C is set to 1 or 0 and written back in the register file before going to fetch the next macroinstruction. Tests for $A = \text{Temp}$, $A \neq \text{Temp}$, $A < \text{Temp}$, and $A \geq \text{Temp}$ are straightforward using these conditions on the ALU output $A - \text{Temp}$. $A > \text{Temp}$ and $A \leq \text{Temp}$, on the other hand, are not simple, but can be done using the negative condition with the subtraction reversed:

$$(\text{Temp} - A < 0) = (\text{Temp} < A) = (A > \text{Temp})$$

If the result is negative, then $A > \text{Temp}$, otherwise $A \leq \text{Temp}$. Voila!

Figure 5.29 contains the last of the DLX microcode and corresponds to the states found in Figures 5.17 and 5.18 (pages 222–224). Microinstruction 50, corresponding to the macroinstruction branch on equal zero, tests if A equals zero. If it does, the macroinstruction branch succeeds, and the microinstruction jumps to the microinstruction 53. This microinstruction loads the PC with the PC-relative address and then jumps to the microcode that fetches the new macroinstruction (location 0). If A does not equal zero, the macroinstruction branch fails, so that the next sequential microinstruction (51) executes, jumping to location 0 without changing the PC.

A state usually corresponds to a single microinstruction, although in a few cases above two microinstructions were needed. The jump and link instructions have the reverse case, with two states collapsing into one microinstruction. The actions in the last two states of jump and link in Figure 5.17 are found in microinstruction 57, and similarly for the jump and link register with microinstruction 59. These microinstructions load the PC with the PC-relative branch address and save C into R31.

Loc	Label	Dest	ALU	S1	S2	C	Misc	Cond	Jump label	Comment
50	Beq:		SUB	A	Constant	0		0?	Branch	<i>Instr is branch =0</i>
51								Uncond	Ifetch	<i>≠0: not taken</i>
52	Bne:		SUB	A	Constant	0		0?	Ifetch	<i>Instr is branch ≠0</i>
53	Branch:	PC	ADD	PC	imm16			Uncond	Ifetch	<i>≠0: taken</i>
54	Jump:	PC	ADD	PC	imm26			Uncond	Ifetch	<i>Jump</i>
55	JReg:	PC	Pass S1	A				Uncond	Ifetch	<i>Jump register</i>
56	JAL:	C	Pass S1	PC						<i>Jump and link</i>
57		PC	ADD	PC	imm26		R31←C	Uncond	Ifetch	<i>Jump & save PC</i>
58	JALR:	C	Pass S1	PC						<i>Jump & link reg</i>
59		PC	Pass S1	A			R31←C	Uncond	Ifetch	<i>Jump & save PC</i>
60	Trap:	IAR	Pass S1	PC						<i>Trap</i>
61		PC	Pass S2		imm26			Uncond	Ifetch	

FIGURE 5.29 The microcode for branch and jump DLX instructions, corresponding to the states in Figures 5.17 and 5.18 on pages 222–224.

Performance of Microcoded Control for DLX

Before trying to improve performance or reduce costs of control, the existing performance must be assessed. Again, the process is to count the clock cycles for each instruction, but this time there is a larger variety in performance.

All instructions execute microinstructions 0, 1, and 2 in Figure 5.23 (page 230), giving a base of 3 clocks plus wait states, depending on the repetition of microinstruction 1. The clock cycles for the rest of the categories are:

- 4 for stores, plus wait states
- 5 for load word, plus wait states
- 6 for load byte or load half (signed or unsigned), plus wait states
- 3 for ALU
- 4 for set
- 2 for branch equal zero (taken or untaken)
- 2 for branch not equal zero (taken)
- 1 for branch not equal zero (untaken)
- 1 for jumps
- 2 for jump and links

Using the instruction mix for GCC in Figure C.4, and assuming an average of 1 wait state per memory access, the CPI is 7.68. This is higher than the hardwired control CPI, because the test for interrupt takes another clock cycle at the beginning, loads and stores are slower, and branch equal zero is slower for the untaken case.

Reducing Cost and Improving Performance of DLX When Control Is Microcoded

The size of a completely unencoded version of the DLX microinstruction is calculated from the number of entries in Figures 5.7 (page 211) and 5.22 (page 229) plus the size of the Constant and Jump address fields. The largest constant in the fields is 24, which requires 5 bits, and the largest address is 61, which requires 6. Figure 5.30 shows the microinstruction fields, the unencoded widths, and the encoded widths. Encoding almost halves the size of control store.

	Dest	ALU operation	Source1	Source2	Constant	Misc	Cond	Jump address	Total
Unencoded	7	11	9	9	5	6	10	6	= 63 bits
Encoded	3	4	4	4	5	3	4	6	= 33 bits

FIGURE 5.30 Width of field in bits of unencoded and encoded microinstruction formats. Note that the Constant and Jump address fields are not encoded in this example, placing fewer restrictions on the microprogram using the encoded format.

The microinstruction can be further shrunk by introducing multiple microinstruction formats and by combining independent fields.

Example

Figure 5.31 shows an encoded version of the original DLX microinstruction format and the version with two formats: one for ALU operations and one for miscellaneous and branch operations. A bit is added to distinguish the two formats. The ALU/Jump (A/J) microinstruction performs the ALU operations specified in the microinstruction; the address of the next microinstruction is specified in the Jump address. For the Transfer/Misc/Branch (T/M/B) microinstruction, the ALU performs Pass S1, while the Misc and Cond fields specify the rest of the operations. The primary change in interpretation of the fields in the new formats is that the ALU condition being tested in the T/M/B format refers to the ALU output from the *prior* A/J microinstruction since there is no ALU operation in T/M/B format. In both formats the Constant and Jump fields are combined into a single field under the assumption they are not used at the same time. (For the A/J format, the appearance of a constant in a source field results in fetching the following microinstruction.) The new formats shrink width from the original 33 bits to 22 bits, but the actual size savings depends on the number of extra microinstructions needed because of the reduced options.

What is the increase in number of microinstructions, compared to the single format, for the microcode in Figure 5.23 (page 230)?

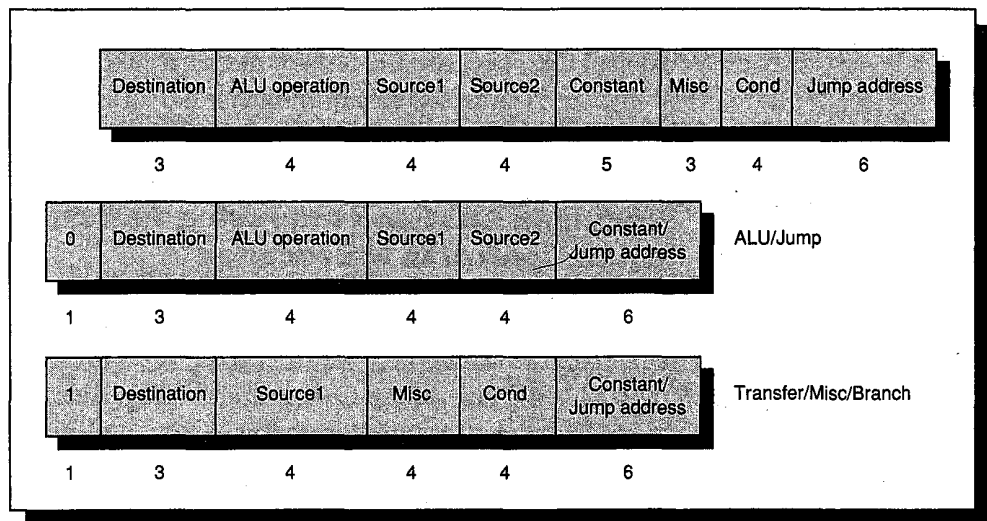


FIGURE 5.31 The original DLX microinstruction format at the top and the dual-format version below. Note that the Misc field is expanded from 3 to 4 bits in the T/M/B to make the two formats the same length.

Answer

Figure 5.32 shows the increase in the number of microinstructions over Figure 5.23 (page 230) because of the restrictions of each format. The five microinstructions in the original format expand to six in the new format. Microinstruction 2 is the only one that expands to two microinstructions for this example.

Loc	Label	Type	Dest	ALU	S1	S2	Misc	Cond	Const/ Jump	Comment
0	Ifetch:	M/T/B		---		---		Interrupt?	Intrpt	<i>Check interrupt</i>
1	Iloop:	M/T/B		---		---	Instr Read	Mem?	Iloop	<i>IR ← M[PC]; wait for memory</i>
2		A/J	PC	ADD	PC	Constant	---	---	4	<i>Increment PC</i>
3		M/T/B		---		---	AB← RF	Decode1		
4	Intrpt:	A/J	IAR	Pass	S1	PC	---	---	5	<i>Interrupt</i>
5		A/J	PC	SUB	Temp	Temp	---	---	Ifetch	<i>PC ← 0 (t minus t=0) & go fetch next instruction</i>

FIGURE 5.32 Version of Figure 5.23 (page 230) using the dual-format microinstruction in Figure 5.31. Note that ALU/Jump microinstructions check the S1 and S2 fields for a constant specifier to see if the next address is sequential (as in microinstruction 2); otherwise they go to the Jump address (as in microinstructions 4 and 5). The microprogrammer changed the last microinstruction to generate a zero by subtracting a register from itself rather than through straightforward use of constant 0. Using the constant would have required an additional microinstruction since this format goes to the next sequential instruction if a constant is used. (See Figure 5.31.)

Sometimes performance can be improved by finding faster sequences of microcode, but normally it requires changes to the hardware. The branch equal zero instruction takes one extra clock cycle when the branch is not taken with hardwired control, but two with microcoded control; while branch not equal zero has the same performance for hardwired and microcoded control. Why would the former differ in performance? Figure 5.29 shows that microinstruction 52 branches on zero to fetch the next microinstruction, which is correct for the branch on not equal zero macroinstruction. Microinstruction 50 also tests for zero for the branch on zero macroinstruction and branches to the microinstruction that loads the new PC. The not zero case is handled by the following microinstruction (51), which jumps to fetch the next instruction—hence, one clock cycle for untaken branch on not equal zero and two for untaken branch on equal zero. One solution is simply to add “not zero” to the microcode branch conditions in Figure 5.22 (page 229) and change the branch on equal microcode to the version in Figure 5.33. Since there are only ten branch conditions, adding the eleventh would not require more than the four bits needed for an encoded version of that field.

Loc	Label	Dest	ALU	S1	S2	C	Misc	Cond	Jump label	Comment
50	Beq:		SUB	A	Constant	0		not 0?	Ifetch	<i>Branch =0</i>
51		PC	ADD	PC	imm16			Uncond	Ifetch	<i>=0: taken</i>

FIGURE 5.33 Branch not equal microcode from Figure 5.29 (page 234) rewritten by using a not zero condition in microinstruction 44.

This change drops the CPI from 7.68 to 7.63 for microcoded control, yet this is still higher than the CPI for hardwired control.

Example

Let's improve microcoded control so that the CPI for GCC is closer to the original CPI under hardwired control.

Answer

The main performance culprit is the separate test for interrupts in Figure 5.23. By modifying the hardware, `decode1` can kill two birds with one stone: In addition to jumping to the appropriate microinstructions corresponding to the opcode, it also jumps to the interrupt microcode if an interrupt is pending. Figure 5.34 shows the revised microcode. This modification saves one clock cycle from each instruction, reducing the CPI to 6.63.

Loc	Label	Dest	ALU	S1	S2	C	Misc	Cond	Jump label	Comment
0	Ifetch:						Instr Read	Mem?	Ifetch	$IR \leftarrow M[PC]$; wait for memory
1		PC	ADD	PC	Constant	4	$AB \leftarrow RF$	Decode1		Also go to interrupt if pending interrupt
2	Intrpt:	IAR	SUB	PC	Constant	4				Interrupt: undo PC increment
3		PC	Pass S2		Constant	0		Uncond	Ifetch	$PC \leftarrow 0$ & go fetch next instruction

FIGURE 5.34 Revised microcode that takes advantage of a change of the hardware to have `decode1` go to microinstruction 2 if there is a pending interrupt. This microinstruction must reverse the increment of PC in the prior microinstruction so that the correct value is saved.

5.8 Fallacies and Pitfalls

Pitfall: Microcode implementing a complex instruction may not be faster than macrocode.

At one time, microcode had the advantage of being fetched from a much faster memory than macrocode. Since caches came into use in 1968, microcode no longer has such a consistent edge in fetch time. Microcode does, however, still have the advantage of using internal temporary registers in the computation, which can be helpful on machines with few general-purpose registers. The disadvantage of microcode is that the algorithms must be selected before the machine is announced and can't be changed until the next model of the archi-

texture; macrocode, on the other hand, can utilize improvements in its algorithms at any time during the life of the machine.

The VAX Index instruction provides an example: The instruction checks to see if the index is between two bounds, one of which is usually zero. The VAX-11/780 microcode uses two compares and two branches to do this, while macrocode can perform the same check in one compare and one branch. The macrocode checks the index against the upper limit using **unsigned** comparisons, rather than two's complement comparisons. This treats a negative index (less than zero and so failing the comparison) as if it were a very large number, thus exceeding the upper limit. (The algorithm can be used with nonzero lower bounds by first subtracting the lower bound from the index.) Replacing the index instruction by this VAX macrocode always improves performance on the VAX-11/780.

Fallacy: If there is space in control store, new instructions are free of cost.

Since the length of control store is usually a power of two, at times there may be unused control store available to expand the instruction set. The analogy here is that of building a house and discovering, near completion, that you have enough land and materials left to add a room. This room wouldn't be free, however, since there would be the costs of labor and maintenance for the life of the home. The temptation to add "free" instructions can only occur when the instruction set is not fixed, as is likely to be the case in the first model of a computer. Because instruction set compatibility is a long-term requirement, all future models of this machine will be forced to include these "free" instructions, even if space is later at a premium. This expansion also ignores the cost of a longer development time to test the added instructions, as well as the possibility of costs of repairing bugs in them after the hardware is shipped.

Fallacy: Users find writable control store helpful.

Bugs in microcode persuaded designers of minicomputers and mainframes that it would be wiser to use RAM than ROM for control store. Doing so would enable microcode bugs to be repaired by shipping customers floppy disks rather than by having the field engineer pull boards and replace chips. Some customers and some manufacturers also decided that users should be allowed to write microcode; this opportunity became known as *writable control store* (WCS). Yet by the time WCS was offered, the world had changed to make WCS less attractive than originally envisioned:

- The tools for writing microcode were much poorer than those for writing macrocode. (The authors and many others stepped into that breach to provide better microprogramming tools.)
- At a time when main memory was expanding, WCS was limited to 1–4KB microinstructions. (Few programming tasks are harder than forcing code into too small a memory.)

- Microcoded control became increasingly tailored to the native macroinstruction set, making microprogramming less useful for tasks other than that for which it was intended.
- With the advent of timesharing, programs might run for only milliseconds before switching to other tasks. This meant that WCS would have to be swapped if more than one program needed it, and reloading WCS could easily take longer than a few milliseconds.
- Timesharing also meant that programs had to be protected from each other. Because, at such a low level, microprograms can circumvent all protection barriers, microprograms written by users were notoriously untrustworthy.
- The increasing demand for virtual memory meant that microprograms had to be restartable—any memory access could force the computation to be shelved.
- Finally, companies like DEC that offered WCS provided no customer support for those who wanted to write microcode.

Many customers ordered WCS, but few benefited from it. The death of WCS has been by a thousand small cuts, and WCS is not an option on current computers.

5.9

Concluding Remarks

In his first paper [1953] Wilkes identified advantages of microprogramming that still hold true today. One of these advantages is that microprogramming helps accommodate change. This can happen late in the development cycle, where simply changing some 0s to 1s in the control store can sometimes save redesigning hardware. A related advantage is that by emulating other instruction sets in microcode, software compatibility is simplified. Microprogramming also reduces the cost of adding more complex instructions to a standard micro-architecture to just the cost of a few more words of control store (although there is the pitfall that once an instruction set is created assuming microprogrammed control, it is difficult to ever build a machine without using it). This flexibility allows hardware construction to begin before the instruction set and microcode have been completely written, because specifying control is just a matter of programming. Finally, microprogramming now has the further advantage of having a large set of tools that have been developed to help write, edit, assemble, and debug microcode.

The drawback of microcode has always been performance. This is because microprogramming is a slave to memory technology: The clock cycle time is limited by the time to read microinstructions from control store. In the 1950s, microprogramming was impractical since virtually the only technology available for control store was the same one used for main memory. In the late 1960s and

early 1970s, semiconductor memory was available for control store, while main memory was constructed from core. The factor of ten in cycle time that differentiated the two technologies opened the door for microcode. The popularity of cache memory in the 1970s once again closed this gap, and machines were again built with the same technology for control store and memory.

For these reasons instruction sets invented since 1985 have not relied on microcode. Though no one likes to predict the future—least of all in writing—it is the authors' opinion that microprogramming is bound to memory technology. If in some future technology ROM becomes much faster than RAM, or if caches are no longer effective, microcode may regain its popularity.

5.10

Historical Perspective and References

Interrupts go back to computer industry pioneers Eckert and Mauchly. Interrupts were first used to signal arithmetic overflow on the UNIVAC I and later to alert a UNIVAC 1103 to start online data collection for a wind tunnel (see Codd [1962]). After the success of the first commercial computer, the UNIVAC 1101 in 1953, the first commercial computer to have interrupts, the 1103, was brought out. Interrupts were first used for I/O by A.L. Leiner in the National Bureau of Standards DYSEAC [Smotherman 1989].

Maurice Wilkes learned computer design in a summer workshop from Eckert and Mauchly and then went on to build the first full-scale, operational, stored-program computer—the EDSAC. From that experience he realized the difficulty of control. He thought of a more centralized control using a diode matrix and, after visiting the Whirlwind computer in the U.S., wrote:

I found that it did indeed have a centralized control based on the use of a matrix of diodes. It was, however, only capable of producing a fixed sequence of 8 pulses—a different sequence for each instruction, but nevertheless fixed as far as a particular instruction was concerned. It was not, I think, until I got back to Cambridge that I realized that the solution was to turn the control unit into a computer in miniature by adding a second matrix to determine the flow of control at the microlevel and by providing for conditional micro-instructions. [Wilkes 1985, 178]

Wilkes [1953] was ahead of his time in recognizing that problem. Unfortunately, the solution was also ahead of its time: To provide control, microprogramming relies on fast memory that was not available in the 1950s. Thus, Wilkes's ideas remained primarily academic conjecture for a decade, although he did construct the EDSAC 2 using microprogrammed control in 1958 with ROM made from magnetic cores.

IBM brought microprogramming into the spotlight in 1964 with the IBM 360 family. Before this event, IBM saw itself as many small businesses selling different machines with their own price and performance levels, but also with their own instruction sets. (Recall that little programming was done in high-level languages, so that programs written for one IBM machine would not run on another.) Gene Amdahl, one of the chief architects of the IBM 360, said that managers of each subsidiary agreed to the 360 family of computers only because they were convinced that microprogramming made it feasible—if you could take the same hardware and microprogram it with several different instruction sets, they reasoned, then you must also be able to take different hardware and microprogram them to run the same instruction set. To be sure of the viability of microprogramming, the IBM vice president of engineering even visited Wilkes surreptitiously and had a “theoretical” discussion of the pros and cons of microcode. IBM believed the idea was so important to their plans that they pushed the memory technology inside the company to make microprogramming feasible.

Stewart Tucker of IBM was saddled with the responsibility of porting software from the IBM 7090 to the new IBM 360. Thinking about the possibilities of microcode, he suggested expanding the control store to include simulators, or interpreters, for older machines. Tucker [1967] coined the term *emulation* for this, meaning full simulation at the microprogrammed level. Occasionally, emulation on the 360 was actually faster than the original hardware. Emulation became so popular with customers in the early years of the 360 that it was sometimes hard to tell which instruction set ran more programs.

Once the giant of the industry began using microcode, the rest soon followed. A difficulty in adopting microcode was that the necessary memory technology was not widely available, but that was soon solved by semiconductor ROM and later RAM. The microprocessor industry followed the same history, with limited resources of the earliest chips forcing hardwired control. But as the resources increased, the advantages of simpler design and ease of change persuaded many to use microprogramming.

With the increasing popularity of microprogramming came more sophisticated instruction sets, including virtual memory. Microprogramming may well have aided the spread of virtual memory, since microcode made it easier to cope with the difficulties that arose from mapping addresses and restarting instructions. The IBM 370 model 138, for example, implemented virtual memory entirely in microcode without any hardware support.

Over the years, most microarchitectures became more and more dedicated to support the intended instruction set, so that reprogramming for a different instruction set failed to offer satisfactory performance. With the passage of time came much larger control stores, and it became possible to consider a machine as elaborate as the VAX. To offer a single chip VAX in 1984 DEC reduced the instructions interpreted by microcode by trapping some instructions and performing them in software: 20% of VAX instructions are responsible for 60% of the microcode, yet are only executed 0.2% of the time. Figure 5.35 shows the

reduction in control store by subsetting the instruction set. (The VAX is so tied to microcode that we venture to predict it will be impossible to build a full-instruction-set VAX without microcode.) The microarchitecture of one of the simpler subsetted VAXes, the MicroVAX-I, is described in Levy and Eckhouse [1989].

	Full instruction set (VLSI VAX)	Subset instruction set (MicroVAX 32)
% instructions implemented	100%	80%
Size of control store (bits)	480 K	64 K
Number of chips in processor	9	2
% performance of VAX-11/780	100%	90%

FIGURE 5.35 By trapping some VAX instructions and addressing modes, control store was reduced almost eight-fold. The second chip of the subset VAX is for floating point.

While this book was being written, a landmark legal precedent concerning microcode was set. The question under litigation in *NEC v. Intel* was whether microcode is like writing, and thereby deserves copyright protection (Intel), or whether it is like hardware, which can be patented but not copyrighted (NEC). The importance of this matter lies in the fact that while it is trivial to get a copyright, getting a patent can take as long as a college education. A program can be copyrighted, so the question then follows: What is and isn't a program? Here is the legislated definition:

A 'computer program' is a set of statements or instructions to be used directly or indirectly in a computer in order to bring about a certain result.

After years of preparation and trial, a judge did declare that a microprogram was a program. The lawyers for the losing side then asked him to rescind his decision on grounds of partiality. They had discovered that through an investment club, the judge owned \$80 of stock belonging to the client he ruled for. (The tempting sum really was only \$80, highly frustrating to one of the authors who acted as an expert witness on the case!) The case was retried, and the new judge ruled that "microcode ... comes squarely within the definition of a 'computer program'..." [Gray 1989, 4]. Of course, the fact that two judges in two different trials made the same decision doesn't mean that the matter is closed—there are still higher levels of appeal available.

References

- CLARK, D. W., P. J. BANNON, AND J. B. KELLER [1988]. "Measuring VAX 8800 performance with a histogram hardware monitor," *Proc. 15th Annual Symposium on Computer Architecture* (May–June), Honolulu, Hawaii, 176–185.
- CODD, E. F. [1962]. "Multiprogramming," in F.L. Alt and M. Rubinoff, *Advances in Computers*, vol. 3, Academic Press, New York, 82.
- EMER, J. S. AND D. W. CLARK [1984]. "A characterization of processor performance in the VAX-11/780," *Proc. 11th Symposium on Computer Architecture* (June), Ann Arbor, Mich., 301–310.
- GRAY, W. P. [1989]. Memorandum of Decision, No. C-84-20799-WPG, U.S. District Court for the Northern District of California (February 7, 1989).
- LEVY, H. M. AND R. H. ECKHOUSE, JR. [1989]. *Computer Programming and Architecture: The VAX*, 2nd ed., Digital Press, Bedford, Mass. 358–372
- MCKEVITT, J., ET AL. [1977]. *8086 Design Report*, internal memorandum.
- PATTERSON, D. A. [1983]. "Microprogramming," *Scientific American* 248:3 (March), 36–43.
- REIGEL, E. W., U. FABER, AND D. A. FISCHER, [1972]. "The Interpreter—a microprogrammable building block system," *Proc. AFIPS 1972 Spring Joint Computer Conf.* 40, 705–723.
- SMOTHERMAN, M. [1989]. "A sequencing-based taxonomy of I/O systems and review of historical machines," *Computer Architecture News* 17:5 (September), 5–15.
- TUCKER, S. G. [1967]. "Microprogram control for the System/360," *IBM Systems Journal* 6:4, 222–241.
- WILKES, M. V. [1953]. "The best way to design an automatic calculating machine," in *Manchester University Computer Inaugural Conf.*, 1951, Ferranti, Ltd., London. (Not published until 1953.) Reprinted in "The Genesis of Microprogramming" in *Annals of the History of Computing* 8:116.
- WILKES, M. V. [1985]. *Memoirs of a Computer Pioneer*, The MIT Press, Cambridge, Mass.
- WILKES, M. V. AND J. B. STRINGER [1953]. "Microprogramming and the design of the control circuits in an electronic digital computer," *Proc. Cambridge Philosophical Society* 49:230–238. Also reprinted in D. P. Siewiorek, C. G. Bell, and A. Newell, *Computer Structures: Principles and Examples* (1982), McGraw-Hill, New York, 158–163, and in "The Genesis of Microprogramming" in *Annals of the History of Computing* 8:116.

EXERCISES

If finite-state diagrams and microprogramming are review topics, you may want to skip over questions 5.5 through 5.14.

5.1 [15/10/15/15] <5.5> One technique that tries to get the best of both the worlds of vertical and horizontal microarchitectures is a *two-level* control store, as illustrated by Figure 5.36. It tries to combine small control-store size with wide instructions. To avoid confusion the bottom level uses the prefix *nano-*, yielding the terms "nanoinstruction," "nanocode," and so forth. This technique was used in the Motorola 68000, 68010, and 68020, but it was originated in the Burroughs D-machine [Reigel, Faber, and Fischer 1972]. The idea is that the first level has many vertical instructions that point to the few unique horizontal instructions in the second level. The Burroughs D-machine was a general-purpose computer offering writable control store. Its microinstructions were 16 bits wide, with 12 of those bits specifying a nanoaddress, and the nanoinstructions were 56 bits wide. One instruction set interpreter used 1124 microinstructions and 123 nanoinstructions.

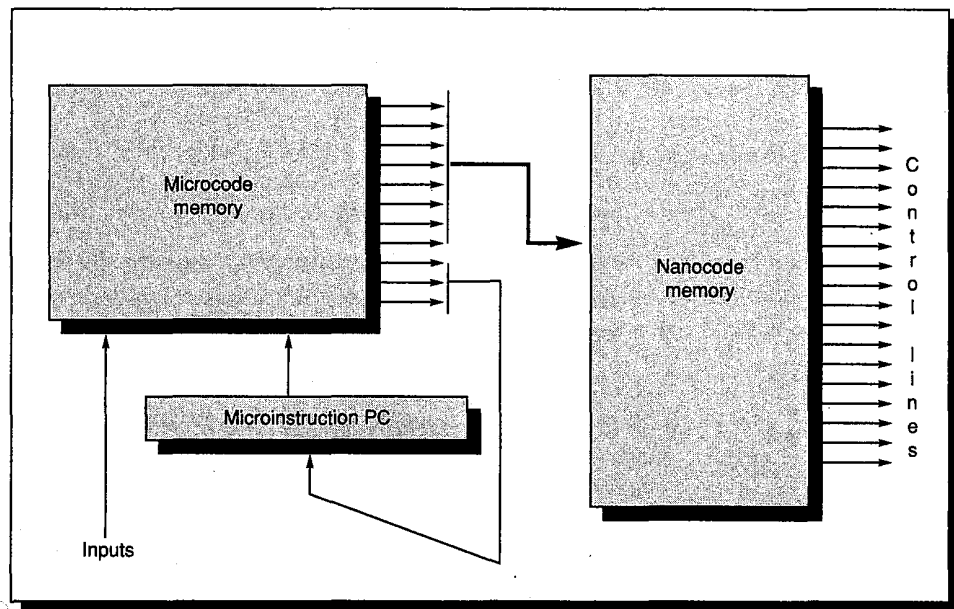


FIGURE 5.36 Two-level microprogrammed implementation showing relationship of microcode and nanocode.

- a. [15] <5.5> What is the general formula showing when a two-level control store scheme like Burroughs D-machine uses fewer bits than a single-level control store? Assume there are M microinstructions each a bits wide and N nanoinstructions each b bits wide.
- b. [10] Was the two-level control store of the D-machine successful in reducing control-store size versus a single-level control store for the interpreter?
- c. [15] After the code was optimized to improve CPI by 10%, the resulting code had 940 microinstructions and 161 nanoinstructions. Was the two-level control store of the D-machine successful in reducing control-store size versus a single-level control store for the **optimized** interpreter?
- d. [15] Did optimization increase or decrease the total number of bits needed to specify control? Why would the number of microinstructions decrease and the number of nanoinstructions increase?

5.2 [15] <5.5,5.6> One advantage of microcode is that it can handle rare cases without having the overhead of invoking the operating system before executing the trap routine. Suppose a machine with a CPI of 1.5 has an operating system that takes 100 clock cycles on a trap before it can execute the appropriate code. Suppose the trap code takes 10 clock cycles whether it is microcode or macrocode. For an instruction occurring 5% of the time, what percentage of the time must it trap before a microcode implementation is 1% faster overall than a macrocode implementation?

5.3 [20/20/30] <4.2,5.5,5.6> Let's explore the impact of subsetting an architecture as described in Figure 5.35. Suppose the MOV_{C3} instruction were left out of a VAX.

- a. [20] Write the VAX macrocode to replace MOV C3.
- b. [20] Assume the operands are placed in registers R0, R1, and R2 after a trap. Using the data for COBOLX in Figure C.1 in Appendix C on instruction usage (assuming all MOV C_ are MOV C3) and assuming the average MOV C3 moves 15 bytes, what would be the percentage change in instruction count if MOV C3 were not interpreted by microcode? (Ignore the cost of traps for this instruction.)
- c. [30] If you have access to a VAX, time the speed of MOV C3 versus a macrocode version of the routine from part a. Assuming that the trap overhead is 20 clock cycles, what is the impact on performance of trapping to software for MOV C3?

5.4 [15] <5.6> Assume we have a machine with a clock cycle time of 10 ns and a base CPI of 5. Because of the possibilities of interrupts we must have extra registers containing copies of the values of the registers at the beginning of the instruction. These registers are usually called *shadow registers*. Assume that the average instruction has two register operands that must be restored on an interrupt. The interrupt rate is 100 interrupts per second, and the interrupt cost is 30 cycles plus the time to restore the shadowed registers, each of which takes 10 cycles. What is the effective CPI after accounting for interrupts? What is the performance lost from interrupts?

5.5-5.7 Given the processor design and finite-state diagram for DLX as modified in the end of the hardwired-control portion of Section 5.7, explore the impact of performance of the following changes. In each case show the modified portion of the finite-state machine, describe the changes to the processor (if necessary), the change in the number of states, and calculate the change in CPI using the DLX instruction mix statistics in Figure C.4 for GCC. Show the reasons for the change.

5.5 [12] <5.7> Like the change to the ALU instructions in the second example in Section 5.7 and shown in Figures 5.20 and 5.21, remove the states that load Temp for the Set instructions in Figure 5.16 first by adding the "X" option and then by increasing the number of states.

5.6 [15] <5.7> Suppose that the memory interface was optimized so that it was not necessary to load MAR before a memory access, nor did the data have to be transferred in MDR for a read or write. Instead, any register on the S1 bus could specify the address, any register on the S2 bus could supply the data on a write, and any register on the Dest bus could receive data on a read.

5.7 [22] <5.7> Most computers overlap the fetching of the next instruction with the execution of the current instruction. Propose a scheme that overlaps all instruction fetches except jumps, branches, and stores. You must reorganize the finite-state machine so that the instruction is already fetched, possibly even partially decoded.

5.8 [15] <5.7> The example in Section 5.7 on page 228 assumes everything but the shifter can scale to 9 ns. Alas, the memory system can rarely scale as easily as the CPU. Reperform the analysis in this example, but this time assume that average number of memory wait states is 2 at the 9-ns clock cycle versus 1 at 10 ns in addition to the slowdown for shifts.

5.9-5.14 These questions address use of the microcoded control of DLX as shown in Figures 5.23, 5.25, and 5.27–5.29. In each case show the modified portion of the microcode; describe the changes to the processor (if necessary), the microinstruction fields (if necessary), and the change in the number of microinstructions; and calculate the change in CPI using the DLX instruction-mix statistics in Appendix C for GCC. Show the reasons for the change.

5.9 [15] <5.7> Like the change to the ALU instructions in the second example in Section 5.7, remove the microinstructions that load Temp for the Set instructions in Figure 5.28 (page 233) first by adding the “X” option and then by increasing the number of microinstructions.

5.10 [25] <5.7> Continuing the example in Figure 5.32 (page 237), rewrite the microcode found in Figure 5.29 (page 234) using the dual-format microinstructions of Figure 5.31 (page 236). What is the relative frequency of each type of microinstruction? What is the savings in control-store size versus the original DLX format? What is the change in CPI?

5.11 [20] <3.4, 5.7> Load byte and Load half take a clock cycle longer than Load word because of the alignment of data (see Figure 3.10 on page 97 and Figure 5.25 on page 231). Propose a change that eliminates the extra clock for these instructions. How does this change affect the CPI of GCC? How does it affect the CPI of TeX?

5.12 [20] <5.6, 5.7> Change the microcode to perform the following interrupt tests: page fault, arithmetic overflow or underflow, misaligned memory accesses, and using undefined instructions. Make whatever changes are needed to the microarchitecture and microinstruction format. What is the change in size and performance to perform these tests?

5.13 [20] <5.7> The computer designer must be careful not to tailor her design too closely to a particular, single program. Reevaluate the performance impact of all the example performance improvements in Exercises 5.9 to 5.12 this time using the average instruction mix data in Figure C.4. How do the programs affect the evaluations?

5.14 [20] <5.6, 5.7> Starting with the microcode in Figures 5.27 (page 233) and 5.34 (page 238), revise the microcode so that the next macroinstruction is fetched as early as possible during the ALU instructions. Assume a “perfect” memory system, taking one clock cycle per memory reference. Although technically this improvement speeds up instructions that **follow** ALU instructions, the easiest way to account for higher performance is as faster ALU instructions. How much faster are the ALU instructions? How does it affect overall performance according to GCC statistics?

5.15 [30] <4,5.6> If you have access to a machine that uses one of the instruction sets in Chapter 4, determine the worst-case interrupt latency for that implementation of the architecture. Be sure you are measuring the raw machine latency and **not** the operating system overhead.

5.16 [30] <5.6> Computer architects have sometimes been forced to support instructions that were never published in the original instruction set manual. This situation arises

because some programs are created that inadvertently set unused instruction fields to values other than the architect expected, which raises havoc when the architect tries to use those values to extend the instruction set. IBM solved that problem in the System 370 by trapping on every possible undefined field. Try executing instructions with undefined fields on a computer to see what happens. Do your new instructions compute anything useful? If so, would you use these new instructions in programs?

5.17 [35] <5.4, 5.5, 5.7> Take the datapath in Figure 5.1 and build a simulator that can perform any of the operations needed to implement the DLX instruction set. Now implement the DLX instruction set using:

Microprogrammed control, and

Hardwired control.

For hardwired control see if you can find PLA minimization and state-assignment programs to reduce the cost of control. From these two designs, determine the performance of each implementation and the cost in terms of gates or in terms of silicon area.

5.18 [35] <2.2, 5.5, 5.7> The similarities between the microinstructions and the macroinstructions of DLX suggest that performance can be gained by writing a program that translates from DLX macrocode to DLX microcode. (This is the insight that inspired WCS.) Write such a program and benchmark it. What is the resulting expansion of code size?

5.19 [50] <2.2, 4.4, 5.10> Recent attempts have been made to run existing software on hardwired control machines by building hand-tuned simulators for popular machines. Write such a simulator for the 8086 instruction set. Run some existing IBM PC programs, and see how fast your simulator is relative to an 8-MHz 8086.

5.20 [Discussion] <4.5,5.5,5.10> Hypothesis: If the first implementation of an architecture uses microprogramming, it affects the instruction set architecture. Why might this be true? Looking at examples in Chapter 4 or elsewhere, give supporting or contradicting evidence from real machines. Which machines will always use microcode? Why? Which machines will never use microcode? Why? What control implementation do you think the architect had in mind when designing the instruction set architecture?

5.21 [Discussion] <5.5,5.10> Wilkes invented microprogramming in large to simplify construction of control. Since 1980 there has been an explosion of computer-aided design software whose goal is also to simplify construction of control. Hypothesis: The advances in computer-aided design software have rendered microprogramming unnecessary. Find evidence to support and refute the hypothesis.

5.22 [Discussion] <5.10> The DLX instructions and the DLX microinstructions have many similarities. What would make it difficult for a compiler to produce DLX microcode rather than macrocode? What changes to the microarchitecture would make the DLX microcode more useful for this application?

It is quite a three-pipe problem.

Sir Arthur Conan Doyle, *The Adventures of Sherlock Holmes*

6.1	What Is Pipelining?	251
6.2	The Basic Pipeline for DLX	252
6.3	Making the Pipeline Work	255
6.4	The Major Hurdle of Pipelining—Pipeline Hazards	257
6.5	What Makes Pipelining Hard to Implement	278
6.6	Extending the DLX Pipeline to Handle Multicycle Operations	284
6.7	Advanced Pipelining—Dynamic Scheduling in Pipelines	290
6.8	Advanced Pipelining—Taking Advantage of More Instruction-Level Parallelism	314
6.9	Putting It All Together: A Pipelined VAX	328
6.10	Fallacies and Pitfalls	334
6.11	Concluding Remarks	337
6.12	Historical Perspective and References	338
	Exercises	343

6

Pipelining

6.1 What Is Pipelining?

Pipelining is an implementation technique whereby multiple instructions are overlapped in execution. Today, pipelining is the key implementation technique used to make fast CPUs.

A pipeline is like an assembly line: Each step in the pipeline completes a part of the instruction. As in a car assembly line, the work to be done in an instruction is broken into smaller pieces, each of which takes a fraction of the time needed to complete the entire instruction. Each of these steps is called a *pipe stage* or a *pipe segment*. The stages are connected one to the next to form a pipe—instructions enter at one end, are processed through the stages, and exit at the other end.

The throughput of the pipeline is determined by how often an instruction exits the pipeline. Because the pipe stages are hooked together, all the stages must be ready to proceed at the same time. The time required between moving an instruction one step down the pipeline is a machine cycle. The length of a machine cycle is determined by the time required for the slowest pipe stage (because all stages proceed at the same time). Often the machine cycle is one clock cycle (sometimes it is two, or rarely more), though the clock may have multiple phases.

The pipeline designer's goal is to balance the length of the pipeline stages. If the stages are perfectly balanced, then the time per instruction on the pipelined machine—assuming ideal conditions (i.e., no stalls)—is equal to

$$\frac{\text{Time per instruction on nonpipelined machine}}{\text{Number of pipe stages}}$$

Under these conditions, the speedup from pipelining equals the number of pipe stages. Usually, however, the stages will not be perfectly balanced; furthermore, pipelining does involve some overhead. Thus, the time per instruction on the pipelined machine will not have its minimum possible value, though it can be close (say within 10%).

Pipelining yields a reduction in the average execution time per instruction. This reduction can be obtained by decreasing the clock cycle time of the pipelined machine or by decreasing the number of clock cycles per instruction, or by both. Typically, the biggest impact is in the number of clock cycles per instruction, though the clock cycle is often shorter in a pipelined machine (especially in pipelined supercomputers). In the advanced pipelining sections of this chapter we will see how deep pipelines can be used to both decrease the clock cycle and maintain a low CPI.

Pipelining is an implementation technique that exploits parallelism among the instructions in a sequential instruction stream. It has the substantial advantage that, unlike some speedup techniques (see Chapters 7 and 10), it is not visible to the programmer. In this chapter we will first cover the concept of pipelining using DLX and a simplified version of its pipeline. We will then look at the problems pipelining introduces and the performance attainable under typical situations. Later in the chapter we will examine advanced techniques that can be used to overcome the difficulties that are encountered in pipelined machines and that may lower the performance attainable from pipelining.

We use DLX largely because its simplicity makes it easy to demonstrate the principles of pipelining. The same principles apply to more complex instruction sets, though the corresponding pipelines are more complex. We will see an example of such a pipeline in the Putting It All Together section.

6.2 The Basic Pipeline for DLX

Remember that in Chapter 5 (Section 5.3) we discussed how DLX could be implemented with five basic execution steps:

1. IF—instruction fetch
2. ID—instruction decode and register fetch
3. EX—execution and effective address calculation
4. MEM—memory access
5. WB—write back

Instruction number	Clock number								
	1	2	3	4	5	6	7	8	9
Instruction <i>i</i>	IF	ID	EX	MEM	WB				
Instruction <i>i</i> +1		IF	ID	EX	MEM	WB			
Instruction <i>i</i> +2			IF	ID	EX	MEM	WB		
Instruction <i>i</i> +3				IF	ID	EX	MEM	WB	
Instruction <i>i</i> +4					IF	ID	EX	MEM	WB

FIGURE 6.1 Simple DLX pipeline. On each clock cycle another instruction is fetched and begins its five-step execution. If an instruction is started every clock cycle, the performance will be five times that of a machine that is not pipelined.

We can pipeline DLX by simply fetching a new instruction **on each clock cycle**. Each of the steps above becomes a *pipe stage*—a step in the pipeline—resulting in the execution pattern shown in Figure 6.1. While each instruction still takes five clock cycles, during each clock cycle the hardware is executing some part of five different instructions.

Pipelining increases the CPU instruction throughput—the number of instructions completed per unit of time—but it does not reduce the execution time of an individual instruction. In fact, it usually slightly increases the execution time of each instruction due to overhead in the control of the pipeline. The increase in instruction throughput means that a program runs faster and has lower total execution time, even though no single instruction runs faster!

The fact that the execution time of each instruction remains unchanged puts limits on the practical depth of a pipeline, as we will see in the next section. Other design considerations limit the clock rate that can be attained by deeper pipelining. The most important consideration is the combined effect of latch delay and clock skew. Latches are required between pipe stages, adding setup time plus the delay through those latches to each clock period. Clock skew also contributes to the lower limit on the clock cycle. Once the clock cycle is as small as the sum of the clock skew and latch overhead, no further pipelining is useful.

Example

Consider a nonpipelined machine with five execution steps of lengths 50 ns, 50 ns, 60 ns, 50 ns, and 50 ns. Suppose that due to clock skew and setup, pipelining the machine adds 5 ns of overhead to each execution stage. Ignoring any latency impact, how much speedup in the instruction execution rate will we gain from a pipeline?

Answer

Figure 6.2 shows the execution pattern on the nonpipelined machine and on the pipelined machine.

The average instruction execution time on the nonpipelined machine is

$$\text{Average instruction execution time} = 50+50+60+50+50 \text{ ns} = 260 \text{ ns}$$

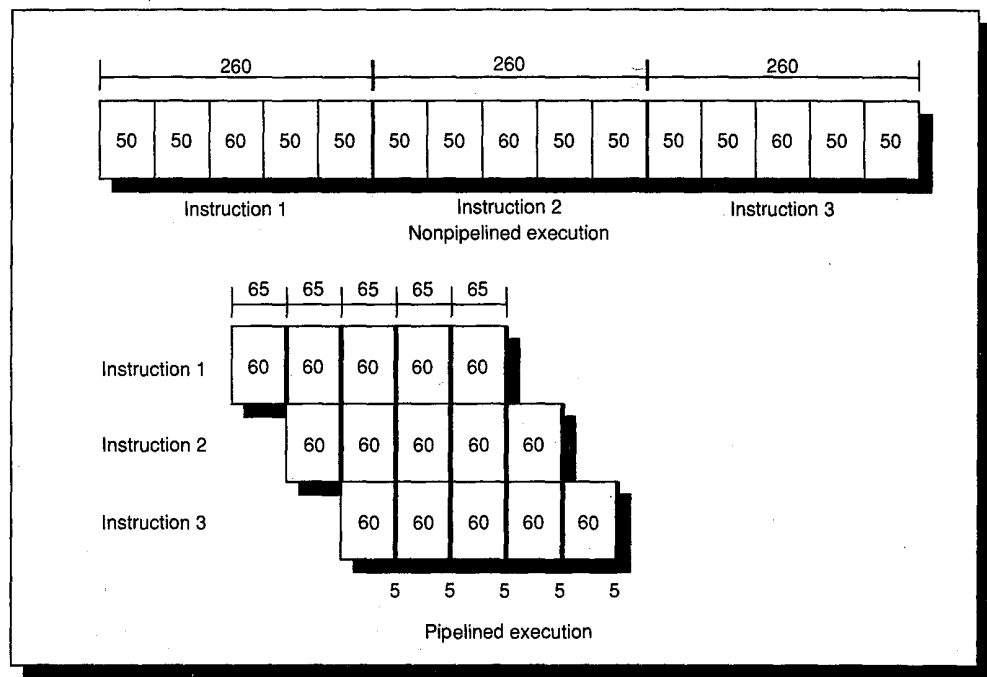


FIGURE 6.2 The execution pattern for three instructions shown for both the nonpipelined and pipelined versions. In the nonpipelined version, the three instructions are executed sequentially. In the pipelined version, the shaded areas represent the overhead of 5 ns per pipestage. The length of the pipestages must all be the same: 60 ns plus the 5-ns overhead. The latency of an instruction increases from 260 ns in the nonpipelined machine to 325 ns in the pipelined machine.

In the pipelined implementation, the clock must run at the speed of the slowest stage plus overhead, which will be 60 + 5 or 65 ns; this is the average instruction execution time. Thus, the speedup from pipelining is

$$\begin{aligned} \text{Speedup} &= \frac{\text{Average instruction time without pipeline}}{\text{Average instruction time with pipeline}} \\ &= \frac{260}{65} = 4 \text{ times} \end{aligned}$$

The 5-ns overhead essentially establishes a limit on the effectiveness of pipelining. If the overhead is not affected by changes in the clock cycle, Amdahl's Law tells us that the overhead limits the speedup.

Because the latches in a pipelined design can have a significant impact on the clock speed, designers have looked for latches that permit the highest possible clock rate. The Earle latch (invented by J. G. Earle [1965]) has three properties that make it especially useful in pipelined machines. First, it is relatively insensitive to clock skew. Second, the delay through the latch is always a constant two-gate delay, avoiding the introduction of skew in the data passing through the latch. Finally, two levels of logic can be done in the latch without increasing the latch delay time. This means that two levels of logic in the pipeline can be overlapped with the latch, so the majority of the overhead from the latch can be

hidden. We will not be analyzing the pipeline designs in this chapter at this level of detail. The interested reader should see Kunkel and Smith [1986].

The next two sections will add refinements and address some problems that can occur in this pipeline. In this discussion (up to the last segment of Section 6.5) we will focus on the pipeline for the integer portion of DLX. The complications that arise in the floating-point pipeline will be treated in Section 6.6.

6.3 Making the Pipeline Work

Your instinct is right if you find it hard to believe that pipelining is as simple as this, because it's not. In this and the following three sections, we will make our DLX pipeline "real" by dealing with problems that pipelining introduces.

To begin with, we have to determine what happens on every clock cycle of the machine and make sure that overlapping instructions doesn't overcommit resources. For example, a single ALU cannot be asked to compute an effective address and perform a subtract operation at the same time. As we will see, the simplicity of the DLX instruction set makes resource evaluation relatively easy.

The operations that occur during instruction execution, which were discussed in Section 5.3 of Chapter 5, are modified to execute in a pipeline as shown in Figure 6.3. The figure lists the major functional units in our DLX implementation, the pipe stages, and what has to happen in each stage of the pipeline. The vertical axis is labeled with the pipeline stages, while the horizontal axis shows major resources. Each intersection shows what happens for that resource in that stage. In Figure 6.4 we will show similar information using the instruction type as the horizontal axis. The combination of instructions that may be in the pipeline at any one time is arbitrary. Thus, the combined needs of all instruction types at any pipe stage determine what resources are needed at that stage.

Every pipe stage is active on every clock cycle. This requires all operations in a pipe stage to complete in one clock and any combination of operations to be able to occur at once. Here are the most important implications for the data path, as specified in Chapter 5:

1. The PC must be incremented on each clock. This must be done in IF rather than ID. This will require an additional incremter, since the ALU is already busy on every cycle and cannot be used to increment the PC.
2. A new instruction must be fetched on every clock—this is also done in IF.
3. A new data word is needed on every clock cycle—this is done in MEM.
4. There must be a separate MDR for loads (LMDR) and stores (SMDR), since when they are back-to-back, they overlap in time.
5. Three additional latches are needed to hold values that are needed later in the pipeline, but may be modified by a subsequent instruction. The values latched are the instruction, the ALU output, and the next PC.

Stage	PC unit	Memory	Data path
IF	$PC \leftarrow PC + 4;$	$IR \leftarrow Mem[PC];$	
ID	$PC1 \leftarrow PC$	$IR1 \leftarrow IR$	$A \leftarrow Rs1; B \leftarrow Rs2;$
EX			$DMAR \leftarrow A + (IR1_{16})^{16} \# \# IR1_{16..31}; SMDR \leftarrow B;$ or $ALUoutput \leftarrow A \text{ op } (B \text{ or } (IR1_{16})^{16} \# \# IR1_{16..31});$ or $ALUoutput \leftarrow PC1 + (IR1_{16})^{16} \# \# IR1_{16..31};$ $cond \leftarrow (A \text{ op } 0);$
MEM	if (cond) $PC \leftarrow ALUoutput$	$LMDR \leftarrow Mem[DMAR]$ or $Mem[DMAR] \leftarrow SMDR$	$ALUoutput1 \leftarrow ALUoutput$
WB			$Rd \leftarrow ALUoutput1 \text{ or } LMDR$

FIGURE 6.3 The table shows the major functional units and what may happen in every pipe stage in each unit. In several of the stages not all of the actions listed can occur, because they apply under different assumptions about the instruction. For example, there are three operations within the ALU during the EX stage. The first occurs only on a load or store; the second on ALU operations (with the input being B or the lower 16 bits of the IR, according to whether the instruction is register-register or register-immediate); the third operation occurs only on branches. For simplicity, we have shown the branch case only—jumps will add a 26-bit offset to the PC. The variables ALUoutput1, PC1, and IR1 save values for use in later stages of the pipeline. Designing the memory system to support a data load or store on every clock cycle is challenging; see Chapter 8 for an in-depth discussion. This type of table and that in Figure 6.4 are loosely based on Davidson's [1971] pipeline reservation tables.

Stage	ALU instruction	Load or store instruction	Branch instruction
IF	$IR \leftarrow Mem[PC];$ $PC \leftarrow PC + 4;$	$IR \leftarrow Mem[PC];$ $PC \leftarrow PC + 4;$	$IR \leftarrow Mem[PC];$ $PC \leftarrow PC + 4;$
ID	$A \leftarrow Rs1; B \leftarrow Rs2; PC1 \leftarrow PC$ $IR1 \leftarrow IR$	$A \leftarrow Rs1; B \leftarrow Rs2; PC1 \leftarrow PC$ $IR1 \leftarrow IR$	$A \leftarrow Rs1; B \leftarrow Rs2; PC1 \leftarrow PC$ $IR1 \leftarrow IR$
EX	$ALUoutput \leftarrow A \text{ op } B;$ or $ALUoutput \leftarrow A \text{ op } ((IR1_{16})^{16} \# \# IR1_{16..31});$	$DMAR \leftarrow A + ((IR1_{16})^{16} \# \# IR1_{16..31});$ $SMDR \leftarrow B;$	$ALUoutput \leftarrow PC1 + ((IR1_{16})^{16} \# \# IR1_{16..31});$ $cond \leftarrow (A \text{ op } 0);$
MEM	$ALUoutput1 \leftarrow ALUoutput$	$LMDR \leftarrow Mem[DMAR];$ or $Mem[DMAR] \leftarrow SMDR;$	if (cond) $PC \leftarrow ALUoutput;$
WB	$Rd \leftarrow ALUoutput1;$	$Rd \leftarrow LMDR;$	

FIGURE 6.4 Events on every pipe stage of the DLX pipeline. Because the instruction is not yet decoded, the first two pipe stages are always identical. Note that it was critical to be able to fetch the registers before decoding the instruction; otherwise another pipeline stage would be required. Due to the fixed instruction format, both register fields are always decoded and the registers accessed (though they are sometimes not needed); the PC and immediate fields can be sent to the ALU as well. At the beginning of the ALU operation the correct inputs are multiplexed in, based on the opcode. With this organization all instruction-dependent operations occur in the EX stage or later. As in Figure 6.3, we include the case for branches, but not jumps, which will have a 26-bit offset rather than a 16-bit offset. Jumps are essentially like branches.

Probably the biggest impact of pipelining on the machine resources is in the memory system. Although the memory-access time has not changed, the peak memory bandwidth must be increased by five times over the nonpipelined machine because two memory accesses are required on every clock in the pipelined machine versus two accesses every five clock cycles in a nonpipelined machine with the same number of steps per instruction. To provide two memory accesses every clock, most machines will use separate instruction and data caches (see Chapter 8, Section 8.3).

During the EX stage, the ALU can be used for three different functions: an effective data-address calculation, a branch-address calculation, or an ALU instruction. Fortunately, the DLX instructions are simple; an instruction in EX does at most one of these, so no conflict arises.

The pipeline we now have for DLX would function just fine if every instruction were independent of every other instruction in the pipeline. In reality, instructions in the pipeline can be dependent on one another; this is the topic of the next section.

6.4 The Major Hurdle of Pipelining— Pipeline Hazards

There are situations, called *hazards*, that prevent the next instruction in the instruction stream from executing during its designated clock cycle. Hazards reduce the performance from the ideal speedup gained by pipelining. There are three classes of hazards:

1. *Structural hazards* arise from resource conflicts when the hardware cannot support all possible combinations of instructions in simultaneous overlapped execution.
2. *Data hazards* arise when an instruction depends on the results of a previous instruction in a way that is exposed by the overlapping of instructions in the pipeline.
3. *Control hazards* arise from the pipelining of branches and other instructions that change the PC.

Hazards in pipelines can make it necessary to stall the pipeline. The major difference between stalls in a pipelined machine and stalls in a nonpipelined machine (such as those we saw in DLX in Chapter 5) occurs because there are multiple instructions under execution at once. A stall in a pipelined machine often requires that some instructions be allowed to proceed, while others are delayed. Typically, when an instruction is stalled, all instructions later in the pipeline than the stalled instruction are also stalled. Instructions earlier than the stalled instruction can continue, but no new instructions are fetched during the stall. We will see several examples of how stalls operate in this section—don't worry, they aren't as complex as they might sound!

A stall causes the pipeline performance to degrade from the ideal performance. Let's look at a simple equation for finding the actual speedup from pipelining, starting with the formula from the previous section.

$$\begin{aligned} \text{Pipeline speedup} &= \frac{\text{Average instruction time without pipeline}}{\text{Average instruction time with pipeline}} \\ &= \frac{\text{CPI without pipelining} * \text{Clock cycle without pipelining}}{\text{CPI with pipelining} * \text{Clock cycle with pipelining}} \\ &= \frac{\text{Clock cycle without pipelining}}{\text{Clock cycle with pipelining}} * \frac{\text{CPI without pipelining}}{\text{CPI with pipelining}} \end{aligned}$$

Remember that pipelining can be thought of as decreasing the CPI or the clock cycle time; let's treat it as decreasing the CPI. The ideal CPI on a pipelined machine is usually

$$\text{Ideal CPI} = \frac{\text{CPI without pipelining}}{\text{Pipeline depth}}$$

Rearranging this and substituting into the speedup equation yields:

$$\text{Speedup} = \frac{\text{Clock cycle without pipelining}}{\text{Clock cycle with pipelining}} * \frac{\text{Ideal CPI} * \text{Pipeline depth}}{\text{CPI with pipelining}}$$

If we confine ourselves to pipeline stalls,

$$\text{CPI with pipelining} = \text{Ideal CPI} + \text{Pipeline stall clock cycles per instruction}$$

We can substitute and obtain:

$$\text{Speedup} = \frac{\text{Clock cycle without pipelining}}{\text{Clock cycle with pipelining}} * \frac{\text{Ideal CPI} * \text{Pipeline depth}}{\text{Ideal CPI} + \text{Pipeline stall cycles}}$$

While this gives a general formula for pipeline speedup (ignoring stalls other than from the pipeline), in most instances a simpler equation can be used. Often, we choose to ignore the potential increase in the clock cycle due to pipelining overhead. This makes the clock rates equal and allows us to drop the first term. A simpler formula can now be used:

$$\text{Pipeline speedup} = \frac{\text{Ideal CPI} * \text{Pipeline depth}}{\text{Ideal CPI} + \text{Pipeline stall cycles}}$$

While we will use this simpler form for evaluating the DLX pipeline, a designer must be careful not to discount the potential impact on clock rate in evaluating pipelining strategies.

Structural Hazards

When a machine is pipelined, the overlapped execution of instructions requires pipelining of functional units and duplication of resources to allow all possible combinations of instructions in the pipeline. If some combination of instructions

cannot be accommodated due to resource conflicts, the machine is said to have a *structural hazard*. The most common instances of structural hazards arise when some functional unit is not fully pipelined. Then a sequence of instructions that all use that functional unit cannot be sequentially initiated in the pipeline. Another common way that structural hazards appear is when some resource has not been duplicated enough to allow all combinations of instructions in the pipeline to execute. For example, a machine may have only one register-file write port, but under certain circumstances, the pipeline might want to perform two writes in a clock cycle. This will generate a structural hazard. When a sequence of instructions encounters this hazard, the pipeline will stall one of the instructions until the required unit is available.

Many pipelined machines share a single memory pipeline for data and instructions. As a result, when an instruction contains a data-memory reference, the pipeline must stall for one clock cycle; the machine cannot fetch the next instruction because the data reference is using the memory port. Figure 6.5 shows what a one-memory-port pipeline looks like when it stalls during a load. We will see another type of stall when we talk about data hazards.

Instruction	Clock cycle number								
	1	2	3	4	5	6	7	8	9
Load instruction	IF	ID	EX	MEM	WB				
Instruction <i>i</i> +1		IF	ID	EX	MEM	WB			
Instruction <i>i</i> +2			IF	ID	EX	MEM	WB		
Instruction <i>i</i> +3				stall	IF	ID	EX	MEM	WB
Instruction <i>i</i> +4						IF	ID	EX	MEM

FIGURE 6.5 A pipeline stalled for a structural hazard—a load with one memory port. With only one memory port, the pipeline cannot initiate a data fetch and instruction fetch in the same cycle. A load instruction effectively steals an instruction-fetch cycle, causing the pipeline to stall—no instruction is initiated on clock cycle 4 (which normally would be instruction *i*+3). Because the instruction being fetched is stalled, all other instructions in the pipeline can proceed normally. The stall cycle will continue to pass through the pipeline.

Example

Suppose that data references constitute 30% of the mix and that the ideal CPI of the pipelined machine, ignoring the structural hazard, is 1.2. Disregarding any other performance losses, how much faster is the ideal machine without the memory structural hazard, versus the machine with the hazard?

Answer

The ideal machine will be faster by the ratio of the speedup of the ideal machine over the real machine. Since the clock rates are unaffected, we can use the following for speedup:

$$\text{Pipeline speedup} = \frac{\text{Ideal CPI} * \text{Pipeline depth}}{\text{Ideal CPI} + \text{Pipeline stall cycles}}$$

Since the ideal machine has no stalls, its speedup is simply $\frac{1.2 * \text{Pipeline depth}}{1.2}$.

The speedup of the real machine is $\frac{1.2 * \text{Pipeline depth}}{1.2 + 0.3 * 1} = \frac{1.2 * \text{Pipeline depth}}{1.5}$.

$$\frac{\text{Speedup}_{\text{ideal}}}{\text{Speedup}_{\text{real}}} = \frac{\left(\frac{1.2 * \text{Pipeline depth}}{1.2} \right)}{\left(\frac{1.2 * \text{Pipeline depth}}{1.5} \right)} = \frac{1.5}{1.2} = 1.25$$

Thus, the machine without the structural hazard is 25% faster.

If all other factors are equal, a machine without structural hazards will always have a lower CPI. Why, then, would a designer allow structural hazards? There are two reasons: to reduce cost and to reduce the latency of the unit. Pipelining all the functional units may be too costly. Machines that support one-clock-cycle memory references require twice as much total memory bandwidth and often have higher bandwidth at the pins. Likewise, fully pipelining a floating-point multiplier consumes lots of gates. If the structural hazard would not occur often, it may not be worth the cost to avoid it. It is also usually possible to design a nonpipelined unit, or one that isn't fully pipelined, with a shorter total delay than a fully pipelined unit. For example, both the CDC 7600 and the MIPS R2010 floating-point unit choose shorter latency (fewer clocks per operation) versus full pipelining. As we will see shortly, reducing latency has other performance benefits and can frequently overcome the disadvantage of the structural hazard.

Example

Many recent machines do not have fully pipelined floating-point units. For example, suppose we had an implementation of DLX with a 5-clock-cycle latency for floating-point multiply, but no pipelining. Will this structural hazard have a large or small performance impact on Spice running on DLX? For simplicity, assume that the floating-point multiplies are uniformly distributed.

Answer

The data in Figure C.4 show that floating-point multiply has a frequency of 6% in Spice. Our proposed pipeline can handle up to a 20% frequency of floating-point multiplies—one every five clock cycles. This means that the performance benefit of fully pipelining the floating-point multiply is likely to be low, as long as the floating-point multiplies are not clustered but are distributed uniformly. If they were clustered, the impact could be much larger.

Data Hazards

A major effect of pipelining is to change the relative timing of instructions by overlapping their execution. This introduces data and control hazards. Data

hazards occur when the order of access to operands is changed by the pipeline versus the normal order encountered by sequentially executing instructions. Consider the pipelined execution of these instructions:

```
ADD    R1, R2, R3
SUB    R4, R1, R5
```

The SUB instruction has a source, R1, that is the destination of the ADD instruction. As shown in Figure 6.6, the ADD instruction writes the value of R1 in the WB pipe stage, but the SUB instruction reads the value during its ID stage. This problem is called a *data hazard*. Unless precautions are taken to prevent it, the SUB instruction will read the wrong value and try to use it. In fact, the value used by the SUB instruction is not even deterministic: Though we might think it logical to assume that SUB would always use the value of R1 that was assigned by an instruction prior to ADD, this is not always the case. If an interrupt should occur between the ADD and SUB instructions, the WB stage of the ADD will complete, and the value of R1 at that point **will** be the result of the ADD. This unpredictable behavior is obviously unacceptable.

Instruction	Clock cycle					
	1	2	3	4	5	6
ADD instruction	IF	ID	EX	MEM	WB—data written here	
SUB instruction		IF	ID—data read here	EX	MEM	WB

FIGURE 6.6 The ADD instruction writes a register that is a source operand for the SUB instruction. But the ADD doesn't finish writing the data into the register file until three clock cycles after SUB begins reading it!

The problem posed in this example can be solved with a simple hardware technique called *forwarding* (also called *bypassing* and sometimes *short-circuiting*). This technique works as follows: The ALU result is always fed back to the ALU input latches. If the forwarding hardware detects that the previous ALU operation has written the register corresponding to a source for the current ALU operation, control logic selects the forwarded result as the ALU input rather than the value read from the register file. Notice that with forwarding, if the SUB is stalled, the ADD will be completed, and the bypass will not be activated, causing the value from the register to be used. This is also true for the case of an interrupt between the two instructions.

In our DLX pipeline, we must pass results to not only the instruction that immediately follows, but also to the instruction after that. By the third instruction down the line, the ID and WB stages overlap; however, as the write is not finished until the end of WB, we must continue to forward the result. Figure 6.7 shows a set of instructions in the pipeline and the forwarding operations that can occur.

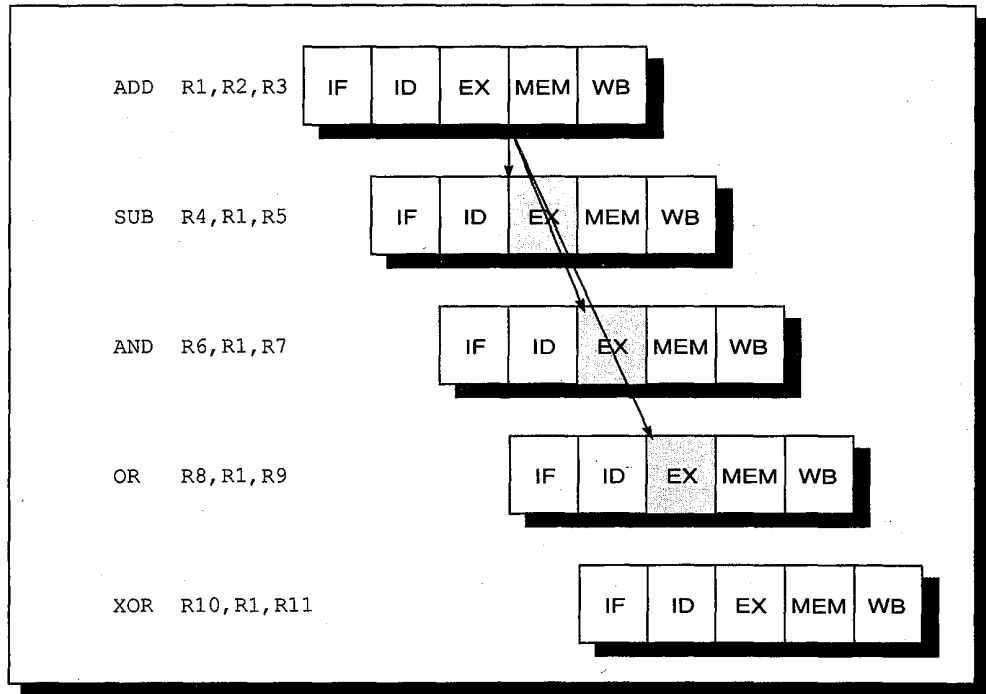


FIGURE 6.7 A set of instructions in the pipeline that need to forward results. The ADD instruction sets R1, and the next four instructions use it. The value of R1 must be bypassed to the SUB, AND, and OR instructions. By the time the XOR instruction goes to read R1 in the ID phase, the ADD instruction has completed WB, and the value is available.

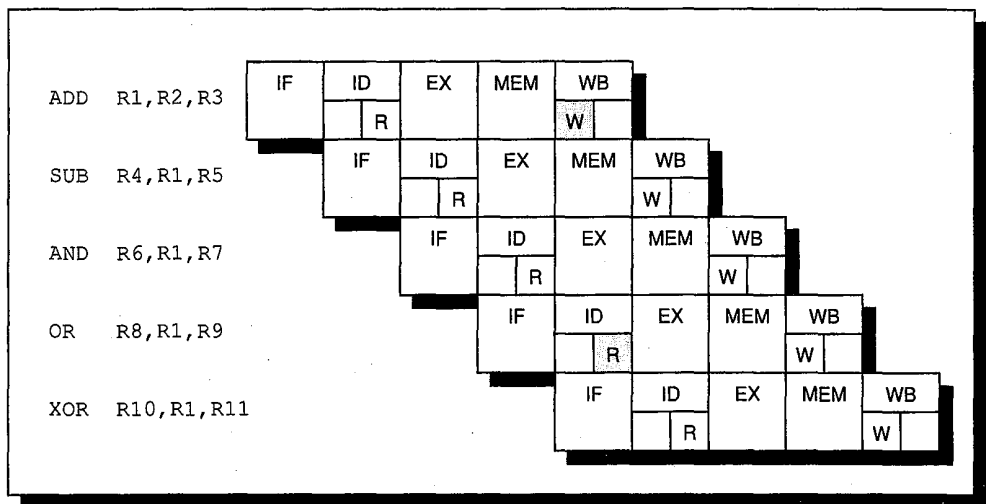


FIGURE 6.8 The same instruction sequence as shown in Figure 6.7, with register reads and writes occurring in opposite halves of the ID and WB stages. The SUB and AND instructions will still require the value of R1 to be bypassed to them, and this will happen as they enter their EX stage. However, by the time of the OR instruction, which also uses R1, the write of R1 has completed, and no forwarding is required. The XOR depends on the ADD, but the value of R1 from the ADD is always written back the cycle before XOR reaches its ID stage and reads it.

It is desirable to cut down the number of instructions that must be bypassed, since each level requires special hardware. Remembering that the register file is accessed twice in a clock cycle, it is possible to do the register writes in the first half of WB and the reads in the second half of ID. This eliminates the need to bypass to a third instruction, as shown in Figure 6.8.

Each level of bypass requires a latch and a pair of comparators to examine whether the adjacent instructions share a destination and a source. Figure 6.9 shows the structure of the ALU and its bypass unit as well as what values are in the bypass registers for the instruction sequence in Figure 6.7. Two ALU result buffers are needed to hold ALU results to be stored into the destination register in the next two WB stages. For ALU operations, the result is always forwarded when the instruction using the result as a source enters its EX stage. (The instruction that computed the value to be forwarded may be in its MEM or WB stages.) The results in the buffers can be inputs into either port on the ALU, via a pair of multiplexers. Multiplexer control can be done by either the control unit (which must then track the destinations and sources of all operations in the pipeline) or locally by logic associated with the bypass (in which case the bypass buffers will contain tags giving the register numbers the values are destined for). In either event, the logic must test if either of the two previous instructions wrote a register that is the input to the current instruction. If so, then the multiplexer select is set to choose from the appropriate result register rather than from the bus. Because the ALU operates in a single pipeline stage, there is no need for a pipeline stall with any combination of ALU instructions once the bypasses have been implemented.

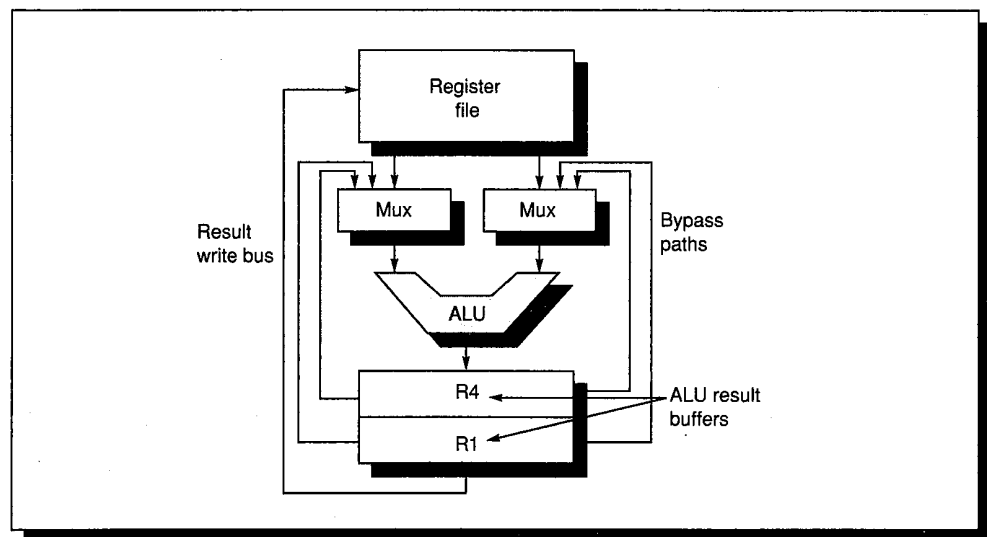


FIGURE 6.9 The ALU with its bypass unit. The contents of the buffer are shown at the point where the AND instruction of the code sequence in Figure 6.8 is about to begin the EX stage. The ADD instruction that computed R1 (in the second buffer) is in its WB stage, and the left input multiplexer is set to pass the just-computed value of R1 (not the value read from the register file) as the first operand to the AND instruction. The result of the subtract, R4, is in the first buffer. These buffers correspond to the variables ALUoutput and ALUoutput1 in Figures 6.3 and 6.4.

A hazard is created whenever there is a dependence between instructions, and they are close enough that the overlap caused by pipelining would change the order of access to an operand. Our example hazards have all been with register operands, but it is also possible for a pair of instructions to create a dependence by writing and reading the same memory location. In our DLX pipeline, however, memory references are always kept in order, preventing this type of hazard from arising. Cache misses could cause the memory references to get out of order if we allowed the processor to continue working on later instructions while an earlier instruction that missed the cache was accessing memory. For DLX's pipeline we just stall the entire pipeline, effectively making the instruction that contained the miss run for multiple clock cycles. In an advanced section of this chapter, Section 6.7, we will discuss machines that allow loads and stores to be executed in an order different from that in the program. All the data hazards discussed in this section, however, involve registers within the CPU.

Forwarding can be generalized to include passing a result directly to the functional unit that requires it: A result is forwarded from the output of one unit to the input of another, rather than just from the result of a unit to the input of the same unit. Take, for example, the following sequence:

```
ADD    R1, R2, R3
SW     25(R1), R1
```

To prevent a stall in this sequence, we would need to forward the value of R1 from the ALU both to the ALU, so that it can be used in the effective address calculation, and to the MDR (memory data register), so that it can be stored without any stall cycles.

Data hazards may be classified as one of three types, depending on the order of read and write accesses in the instructions. By convention, the hazards are named by the ordering in the program that must be preserved by the pipeline. Consider two instructions i and j , with i occurring before j . The possible data hazards are:

- **RAW** (*read after write*) — j tries to read a source before i writes it, so j incorrectly gets the old value. This is the most common type of hazard and the one that appears in Figures 6.6 and 6.7.
- **WAR** (*write after read*) — j tries to write a destination before it is read by i , so i incorrectly gets the new value. This cannot happen in our example pipeline because all reads are early (in ID) and all writes are late (in WB). This hazard occurs when there are some instructions that write results early in the instruction pipeline, and other instructions that read a source after a write of an instruction later in the pipeline. For example, autoincrement addressing can create a WAR hazard.
- **WAW** (*write after write*) — j tries to write an operand before it is written by i . The writes end up being performed in the wrong order, leaving the value written by i rather than the value written by j in the destination. This hazard is present only in pipelines that write in more than one pipe stage (or allow an

instruction to proceed even when a previous instruction is stalled). The DLX pipeline writes a register only in WB and avoids this class of hazards.

Note that the RAR (*read after read*) case is not a hazard.

Not all data hazards can be handled without a performance effect. Consider the following sequence of instructions:

```
LW  R1, 32(R6)
ADD  R4, R1, R7
SUB  R5, R1, R8
AND  R6, R1, R7
```

This case is different from the situation with back-to-back ALU operations. The LW instruction does not have the data until the end of the MEM cycle, while the ADD instruction needs to have the data by the beginning of that clock cycle. Thus, the data hazard from using the result of a load instruction cannot be completely eliminated with simple hardware. We can forward the result immediately to the ALU from the MDR, and for the SUB instruction—which begins two clock cycles after the load—the result arrives in time, as shown in Figure 6.10. However, for the ADD instruction, the forwarded result arrives too late—at the end of a clock cycle, though it is needed at the beginning.

LW R1, 32(R6)	IF	ID	EX	MEM	WB
ADD R4, R1, R7		IF	ID	EX	MEM
SUB R5, R1, R8			IF	ID	EX
AND R6, R1, R7				IF	ID

FIGURE 6.10 Pipeline hazard occurring when the result of a load instruction is used by the next instruction as a source operand and is forwarded. The value is available when it returns from memory at the end of the load instruction’s MEM cycle. However, it is needed at the beginning of that clock cycle for the ADD (the EX stage of the add). The load value can be forwarded to the SUB instruction and will arrive in time for that instruction (EX). The AND can simply read the value during ID since it reads the registers in the second half of the cycle and the value is written in the first half.

The load instruction has a delay or latency that cannot be eliminated by forwarding alone—to do so would require the data-access time to be zero. The most common solution to this problem is a hardware addition called a pipeline interlock. In general, a *pipeline interlock* detects a hazard and stalls the pipeline until the hazard is cleared. In this case, the interlock stalls the pipeline beginning with the instruction that wants to use the data until the sourcing instruction produces it. This delay cycle, called a *pipeline stall* or *bubble*, allows the load data to arrive from memory; it can now be forwarded by the hardware. The CPI for the stalled instruction increases by the length of the stall (one clock cycle in this case). The stalled pipeline is shown in Figure 6.11.

Any instruction	IF	ID	EX	MEM	WB					
LW R1, 32(R6)		IF	ID	EX	MEM	WB				
ADD R4, R1, R7			IF	ID	stall	EX	MEM	WB		
SUB R5, R1, R8				IF	stall	ID	EX	MEM	WB	
AND R6, R1, R7					stall	IF	ID	EX	MEM	WB

FIGURE 6.11 The effect of the stall on the pipeline. All instructions starting with the instruction that has the dependence are delayed. With the delay, the value of the load that returns in MEM can now be forwarded to the EX cycle of the ADD instruction. Because of the stall, the SUB instruction will now read the value from the registers during its ID cycle rather than having it forwarded from the MDR.

The process of letting an instruction move from the instruction decode stage (ID) into the execution stage (EX) of this pipeline is usually called *instruction issue*; and an instruction that has made this step is said to have *issued*. For the DLX integer pipeline, all the data hazards can be checked during the ID phase of the pipeline. If a data hazard exists, the instruction is stalled before it is issued. Later in this chapter, we will look at situations where instruction issue is much more complex. Detecting interlocks early in the pipeline reduces the hardware complexity because the hardware never has to suspend an instruction that has updated the state of the machine, unless the entire machine is stalled.

Example

Suppose that 20% of the instructions are loads, and half the time the instruction following a load instruction depends on the result of the load. If this hazard creates a single-cycle delay, how much faster is the ideal pipelined machine (with a CPI of 1) that does not delay the pipeline, compared to a more realistic pipeline? Ignore any stalls other than pipeline stalls.

Answer

The ideal machine will be faster by the ratio of the CPIs. The CPI for an instruction following a load is 1.5, since they stall half the time. Since loads are 20% of the mix, the effective CPI is $(0.8 \cdot 1 + 0.2 \cdot 1.5) = 1.1$. This yields a performance ratio of $\frac{1.1}{1}$. Hence, the ideal machine is 10% faster.

Many types of stalls are quite frequent. The typical code-generation pattern for a statement such as $A=B+C$ produces a stall for a load of the second data value. Figure 6.12 shows that the store need not result in another stall, since the result of the addition can be forwarded to the MDR. Machines where the operands may come from memory for arithmetic operations will need to stall the pipeline in the middle of the instruction to wait for memory to complete its access.

LW	R1, B	IF	ID	EX	MEM	WB				
LW	R2, C		IF	ID	EX	MEM	WB			
ADD	R3, R1, R2			IF	ID	stall	EX	MEM	WB	
SW	A, R3				IF	stall	ID	EX	MEM	WB

FIGURE 6.12 The DLX code sequence for $A=B+C$. The ADD instruction must be stalled to allow the load of C to complete. The SW need not be delayed further because the forwarding hardware passes the result from the ALU directly to the MDR for storing.

Rather than just allow the pipeline to stall, the compiler could try to schedule the pipeline to avoid these stalls, by rearranging the code sequence to eliminate the hazard. For example, the compiler would try to avoid generating code with a load followed by an immediate use of the load destination register. This technique, called *pipeline scheduling* or *instruction scheduling*, was first used in the 1960s, and became an area of major interest in the 1980s as pipelined machines became more widespread.

Example

Generate DLX code that avoids pipeline stalls for the following sequence:

```
a = b + c;
d = e - f;
```

Assume loads have a latency of one clock cycle.

Answer

Here is the scheduled code:

```
LW Rb, b
LW Rc, c
LW Re, e           ; swapped with next instruction to avoid stall
ADD Ra, Rb, Rc
LW Rf, f
SW a, Ra           ; store/load interchanged to avoid stall in SUB
SUB Rd, Re, Rf
SW d, Rd
```

Both load interlocks (LW Rc, c/ADD Ra, Rb, Rc and LW Rf, f/SUB Rd, Re, Rf) have been eliminated. There is a dependence between the ALU instruction and the store, but the pipeline structure allows the result to be forwarded. Notice that the use of different registers for the first and second statements was critical for this schedule to be legal. In particular, if the variable e were loaded into the same register as b or c, this schedule would not be legal. In

general, pipeline scheduling can increase the register count required. In Section 6.8, we will see that this increase can be substantial for machines that can issue multiple instructions in one clock.

This technique works sufficiently well that some machines rely on software to avoid this type of hazard. A load requiring that the following instruction not use its result is called a *delayed load*. The pipeline slot after a load is often called the *load delay* or *delay slot*. When the compiler cannot schedule the interlock, a no-op instruction may be inserted. This does not affect running time, but only increases the code space versus a machine with the interlock. Whether or not the hardware detects this interlock and stalls the pipeline, performance will be enhanced if the compiler schedules instructions. If the stall occurs, the performance impact will be the same, whether the machine executes an idle cycle or executes a no-op. Figure 6.13 shows that scheduling can eliminate the majority of these delays. It is clear from this figure that load delays in GCC are significantly harder to schedule than in Spice or TeX.

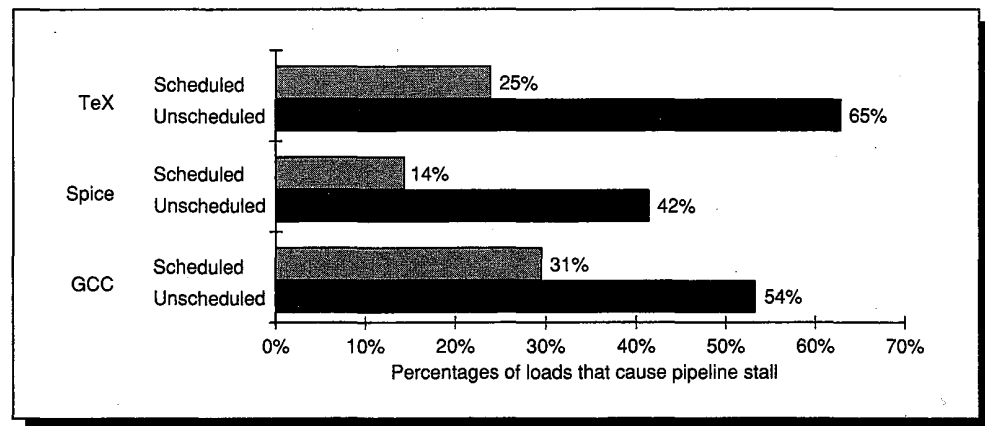


FIGURE 6.13 Percentage of the loads that result in a stall with the DLX pipeline. The black bars show the amount without compiler scheduling; the gray bars show the effect of a good, but simple, scheduling algorithm. These data show scheduling effectiveness after global optimization (see Chapter 3, Section 3.7). Global optimization actually makes scheduling relatively harder because there are fewer candidates available for scheduling into delay slots. For example, on GCC and TeX, when the programs are scheduled but not globally optimized, the percentage of load delays that result in a stall drops to 22% and 19%, respectively.

Implementing Data Hazard Detection in Simple Pipelines

How pipeline interlocks are implemented depends quite heavily on the length and complexity of the pipeline. For a complex machine with long-running instructions and multicycle interdependences, a central table that keeps track of the availability of operands and the outstanding writes may be needed (see Sec-

tion 6.7). For the DLX integer pipeline, the only interlock we need to enforce is load followed by immediate use. This can be done with a simple comparator that looks for this pattern of load destination and source. The hardware required to detect and control the load data hazard and to forward the load result is as follows:

- Additional multiplexers on the inputs to the ALU (just as was required for the bypass hardware for register–register instructions)
- Extra paths from the MDR to both multiplexer inputs to the ALU
- A buffer to save the destination-register numbers from the prior two instructions (the same as for register–register forwarding)
- Four comparators to compare the two possible source register fields with the destination fields of the prior instructions and look for a match

The comparators check for a load interlock at the beginning of the EX cycle. The four possibilities and the required actions are shown in Figure 6.14.

For DLX, the hazard detection and forwarding hardware is reasonably simple; we will see that things become much more complicated when the pipelines are very deep (Section 6.6). But before we do that, let’s see what happens with branches in our DLX pipeline.

Situation	Example code sequence	Action
No dependence	LW R1 , 45 (R2) ADD R5, R6, R7 SUB R8, R6, R7 OR R9, R6, R7	No hazard possible because no dependence exists on R1 in the immediately following three instructions.
Dependence requiring stall	LW R1 , 45 (R2) ADD R5, R1 , R7 SUB R8, R6, R7 OR R9, R6, R7	Comparators detect the use of R1 in the ADD and stall the ADD (and SUB and OR) before the ADD begins EX.
Dependence overcome by forwarding	LW R1 , 45 (R2) ADD R5, R6, R7 SUB R8, R1 , R7 OR R9, R6, R7	Comparators detect use of R1 in SUB and forward result of load to ALU in time for SUB to begin EX.
Dependence with accesses in order	LW R1 , 45 (R2) ADD R5, R6, R7 SUB R8, R6, R7 OR R9, R1 , R7	No action required because the read of R1 by OR occurs in the second half of the ID phase, while the write of the loaded data occurred in the first half. See Figure 6.8 (page 262).

FIGURE 6.14 Situations that the pipeline hazard detection hardware can see by comparing the destination and sources of adjacent instructions. This table indicates that the only compare needed is between the destination and the sources on the two instructions following the instruction that wrote the destination. In the case of a stall, the pipeline dependences will look like the third case, once execution continues.

Control Hazards

Control hazards can cause a greater performance loss for our DLX pipeline than do data hazards. When a branch is executed, it may or may not change the PC to something other than its current value plus 4. (Recall that if a branch changes the PC to its target address, it is a *taken* branch; if it falls through, it is *not taken*, or *untaken*.) If instruction i is a taken branch, then the PC is normally not changed until the end of MEM, after the completion of the address calculation and comparison, as shown in Figure 6.4 (page 256). This means stalling for three clock cycles, at the end of which the new PC is known and the proper instruction can be fetched. This effect is called a *control* or *branch hazard*. Figure 6.15 shows a three-cycle stall for a control hazard.

Branch instruction	IF	ID	EX	MEM	WB					
Instruction $i+1$		<i>stall</i>	<i>stall</i>	<i>stall</i>		IF	ID	EX	MEM	WB
Instruction $i+2$			<i>stall</i>	<i>stall</i>	<i>stall</i>	IF	ID	EX	MEM	WB
Instruction $i+3$				<i>stall</i>	<i>stall</i>	<i>stall</i>	IF	ID	EX	MEM
Instruction $i+4$					<i>stall</i>	<i>stall</i>	<i>stall</i>	IF	ID	EX
Instruction $i+5$						<i>stall</i>	<i>stall</i>	<i>stall</i>	IF	ID
Instruction $i+6$							<i>stall</i>	<i>stall</i>	<i>stall</i>	IF

FIGURE 6.15 Ideal DLX pipeline stalling after a control hazard. The instruction labeled instruction $i+k$ represents the k th instruction executed after the branch. There is a difficulty in that the branch instruction is not decoded until after instruction $i+1$ has been fetched. This figure shows the conceptual difficulty, while Figure 6.16 shows what really happens.

Branch instruction	IF	ID	EX	MEM	WB					
Instruction $i+1$		IF	<i>stall</i>	<i>stall</i>		IF	ID	EX	MEM	WB
Instruction $i+2$			<i>stall</i>	<i>stall</i>	<i>stall</i>	IF	ID	EX	MEM	WB
Instruction $i+3$				<i>stall</i>	<i>stall</i>	<i>stall</i>	IF	ID	EX	MEM
Instruction $i+4$					<i>stall</i>	<i>stall</i>	<i>stall</i>	IF	ID	EX
Instruction $i+5$						<i>stall</i>	<i>stall</i>	<i>stall</i>	IF	ID
Instruction $i+6$							<i>stall</i>	<i>stall</i>	<i>stall</i>	IF

FIGURE 6.16 What might really happen in the DLX pipeline. Instruction $i+1$ is fetched, but the instruction is ignored and the fetch is restarted once the branch target is known. It is probably obvious that if the branch is not taken, the second IF for instruction $i+1$ is redundant. This will be addressed shortly.

The pipeline in Figure 6.15 is not possible because we don't know that the instruction is a branch until after the fetch of the next instruction. Figure 6.16 fixes this by simply redoing the fetch once the target is known.

Three clock cycles wasted for every branch is a significant loss. With a 30% branch frequency and an ideal CPI of 1, the machine with branch stalls achieves

only about half the ideal speedup from pipelining. Thus, reducing the branch penalty becomes critical. The number of clock cycles in a branch stall can be reduced in two steps:

1. Find out whether the branch is taken or not earlier in the pipeline.
2. Compute the taken PC (address of the branch target) earlier.

To optimize the branch behavior, **both** of these must be done—it doesn't help to know the target of the branch without knowing whether the next instruction to execute is the target or the instruction at PC+4. Both steps should be taken as early in the pipeline as possible.

In DLX, the branches (BEQZ and BNEZ) require testing only equality to zero. Thus, it is possible to complete this decision by the end of the ID cycle using special logic devoted to this test. To take advantage of an early decision on whether the branch is taken, both PCs (taken and not taken) must be computed early. Computing the branch target address requires a separate adder, which can add during ID. With the separate adder and a branch decision made during ID, there is only a one-clock-cycle stall on branches. Figure 6.17 shows the branch portion of the revised resource allocation table from Figure 6.4 (page 256).

In some machines, branch hazards are even more expensive in clock cycles than in our example, since the time to evaluate the branch condition and compute the destination can be even longer. For example, a machine with separate

Pipe stage	Branch instruction
IF	IR←Mem[PC]; PC←PC+4;
ID	A←Rs1; B←Rs2; PC1←PC; IR1←IR; BTA←PC+(IR₁₆)¹⁶##IR_{16..31}) if (Rs1 op 0) PC←BTA
EX	
MEM	
WB	

FIGURE 6.17 Revised pipeline structure (see Figure 6.4, page 256) showing the use of a separate adder to compute the branch target address. The operations that are new or have changed are in bold. Because the branch target address (BTA) addition happens during ID, it will happen for all instructions; the branch condition (Rs1 op 0) will also be done for all instructions. The last operation in ID is to replace the PC. We must know that the instruction is a branch before we perform this step. This requires decoding the instruction before the end of ID, or doing this operation at the very beginning of EX when the PC is sent out. Because the branch is done by the end of ID, the EX, MEM, and WB stages are unused for branches. An additional complication arises for jumps that have a longer offset than branches. We can resolve this by using an additional adder that sums the PC and lower 26 bits of the IR. Alternatively, we could attempt a clever scheme that does a 16-bit add in the first half of the cycle and determines whether to add in 10 bits from IR in the second half of the cycle, by decoding the jump opcodes early.

decode and register fetch stages will probably have a *branch delay*—the length of the control hazard—that is at least one clock cycle longer. The branch delay, unless it is dealt with, turns into a branch penalty. Many VAXes have branch delays of four clock cycles or more, and large, deeply pipelined machines often have branch penalties of six or seven. In general, the deeper the pipeline, the worse the branch penalty in clock cycles. Of course, the relative performance effect of a longer branch penalty depends on the overall CPI of the machine. A high CPI machine can afford to have more expensive branches because the percentage of the machine's performance that will be lost from branches is less.

Before talking about methods for reducing the pipeline penalties that can arise from branches, let's take a brief look at the dynamic behavior of branches.

Branch Behavior in Programs

Since branches can dramatically affect pipeline performance, we should look at their behavior so as to get some ideas about how the penalties of branches and jumps might be reduced. We already know the branch frequencies for our programs from Chapter 4. Figure 6.18 reviews the overall frequency of control-flow operations for three of the machines and gives the breakdown between branches and jumps.

All of the machines show a conditional branch frequency of 11%–17%, while the frequency of unconditional branches varies between 2% and 8%. An obvious

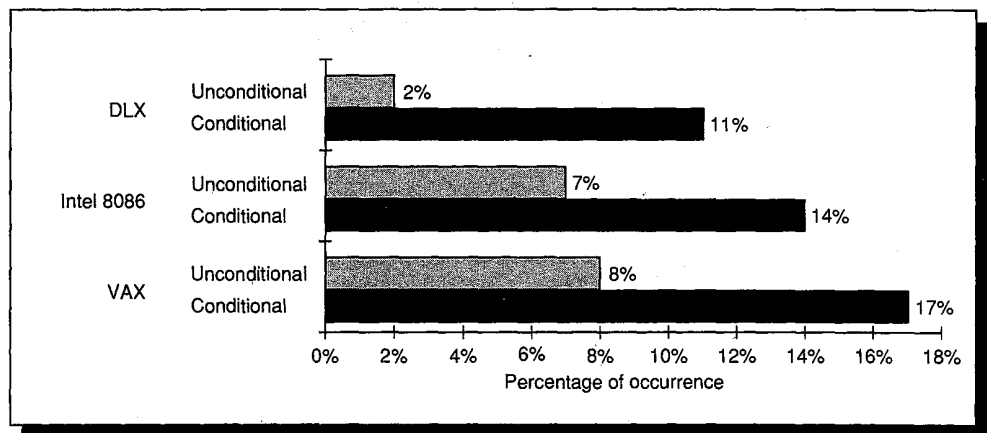


FIGURE 6.18 The frequency of instructions (branches, jumps, calls, and returns) that may change the PC. These data represent the average over the programs measured in Chapter 4. Instructions are divided into two classes: branches, which are conditional (including loop branches), and those that are unconditional (jumps, calls, and returns). The 360 is omitted because the ordinary unconditional branches are not separated from the conditional branches. Emer and Clark [1984] reported that 38% of the instructions executed in their measurements of the VAX were instructions that could change the PC. They measured that 67% of these instructions actually cause a branch in control flow. Their data were taken on a timesharing workload and reflect many uses; their measurement of branch frequency is much higher than the one in this chart.

question is, how many of the branches are taken? Knowing the breakdown between taken and untaken branches is important because this will affect strategies for reducing the branch penalties. For the VAX, Clark and Levy [1984] measured simple conditional branches to be taken with a frequency of just about 50%. Other branches, which occur much less often, have different ratios. Most bit-testing branches are not taken, and loop branches are taken with about 90% probability.

For DLX, we measured the branch behavior in Chapter 3 and summarized it in Figure 3.22 (page 107). That data showed 53% of the conditional branches are taken. Finally, 75% of the branches executed are forward-going branches. With this data in mind, let's look at ways to reduce branch penalties.

Reducing Pipeline Branch Penalties

There are several methods for dealing with the pipeline stalls due to branch delay, and four simple compile-time schemes are discussed in this section. In these schemes the predictions are static—they are fixed for each branch during the entire execution, and the predictions are compile-time guesses. More ambitious schemes using hardware to predict branches dynamically are discussed in Section 6.7.

The easiest scheme is to freeze the pipeline, holding any instructions after the branch until the branch destination is known. The attractiveness of this solution lies primarily in its simplicity. It is the solution used earlier in the pipeline shown in Figures 6.15 and 6.16.

A better and only slightly more complex scheme is to predict the branch as not taken, simply allowing the hardware to continue as if the branch were not

Untaken branch instruction	IF	ID	EX	MEM	WB				
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				
Instruction $i+1$		IF	IF	ID	EX	MEM	WB		
Instruction $i+2$			<i>stall</i>	IF	ID	EX	MEM	WB	
Instruction $i+3$				<i>stall</i>	IF	ID	EX	MEM	WB
Instruction $i+4$					<i>stall</i>	IF	ID	EX	MEM

FIGURE 6.19 The predict-not-taken scheme and the pipeline sequence when the branch is untaken (on the top) and taken (on the bottom). When the branch is untaken, determined during ID, we have fetched the fall through and just continue. If the branch is taken during ID, we restart the fetch at the branch target. This causes all instructions following the branch to stall one clock cycle.

executed. Here, care must be taken not to change the machine state until the branch outcome is definitely known. The complexity that arises from this—that is, knowing when the state might be changed by an instruction and how to “back out” a change—might cause us to reconsider the simpler solution of flushing the pipeline. In the DLX pipeline, this *predict-not-taken* scheme is implemented by continuing to fetch instructions as if the branch were a normal instruction. The pipeline looks as if nothing out of the ordinary is happening. If the branch is taken, however, we need to stop the pipeline and restart the fetch. Figure 6.19 shows both situations.

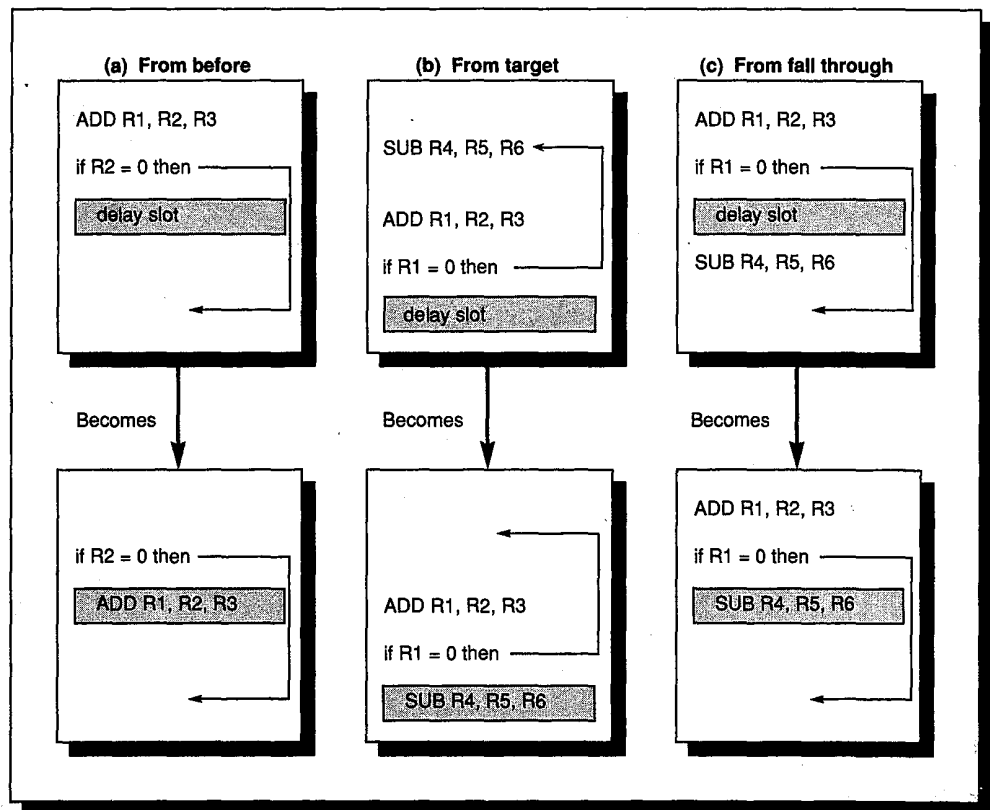


FIGURE 6.20 Scheduling the branch-delay slot. The top picture in each pair shows the code before scheduling, and the bottom picture shows the scheduled code. In (a) the delay slot is scheduled with an independent instruction from before the branch. This is the best choice. Strategies (b) and (c) are used when (a) is not possible. In the code sequences for (b) and (c), the use of R1 in the branch condition prevents the ADD instruction (whose destination is R1) from being moved after the branch. In (b) the branch-delay slot is scheduled from the target of the branch; usually the target instruction will need to be copied because it can be reached by another path. Strategy (b) is preferred when the branch is taken with high probability, such as a loop branch. Finally, the branch may be scheduled from the not-taken fall through, as in (c). To make this optimization legal for (b) or (c), it must be “OK” to execute the SUB instruction when the branch goes in the unexpected direction. By “OK” we mean that the work is wasted, but the program will still execute correctly. This is the case, for example, if R4 were a temporary register unused when the branch goes in the unexpected direction.

An alternative scheme is to predict the branch as taken. As soon as the branch is decoded and the target address is computed, we assume the branch to be taken and begin fetching and executing at the target. Since in our DLX pipeline we don't know the target address any earlier than we know the branch outcome, there is no advantage in this approach. However, in some machines—especially those with condition codes or more powerful (and hence slower) branch conditions—the branch target is known before the branch outcome, and this scheme makes sense.

Some machines have used another technique called delayed branch, which has been used in many microprogrammed control units. In a *delayed branch*, the execution cycle with a branch delay of length n is:

```

branch instruction
sequential successor1
sequential successor2
.....
sequential successorn
branch target if taken
    
```

The sequential successors are in the *branch-delay slots*. As with load-delay slots, the job of the software is to make the successor instructions valid and useful. A number of optimizations are used. Figure 6.20 shows the three ways in which the branch delay can be scheduled. Figure 6.21 shows the different constraints for each of these branch-scheduling schemes, as well as situations in which they win.

The primary limitations on delayed-branch scheduling arise from the restrictions on the instructions that are scheduled into the delay slots and from our ability to predict at compile time whether a branch is likely to be taken or not. Figure 6.22 shows the effectiveness of the branch scheduling in DLX with a single branch-delay slot using a simple branch-scheduling algorithm. It shows that

Scheduling strategy	Requirements	Improves performance when?
(a) From before branch	Branch must not depend on the rescheduled instructions.	Always.
(b) From target	Must be OK to execute rescheduled instructions if branch is not taken. May need to duplicate instructions.	When branch is taken. May enlarge program if instructions are duplicated.
(c) From fall through	Must be OK to execute instructions if branch is taken.	When branch is not taken.

FIGURE 6.21 Delayed-branch-scheduling schemes and their requirements. The origin of the instruction being scheduled into the delay slot determines the scheduling strategy. The compiler must enforce the requirements when looking for instructions to schedule the delay slot. When the slots cannot be scheduled, they are filled with no-op instructions. In strategy (b), if the branch target is also accessible from another point in the program—as it would be if it were the head of a loop—the target instructions must be copied and not just moved.

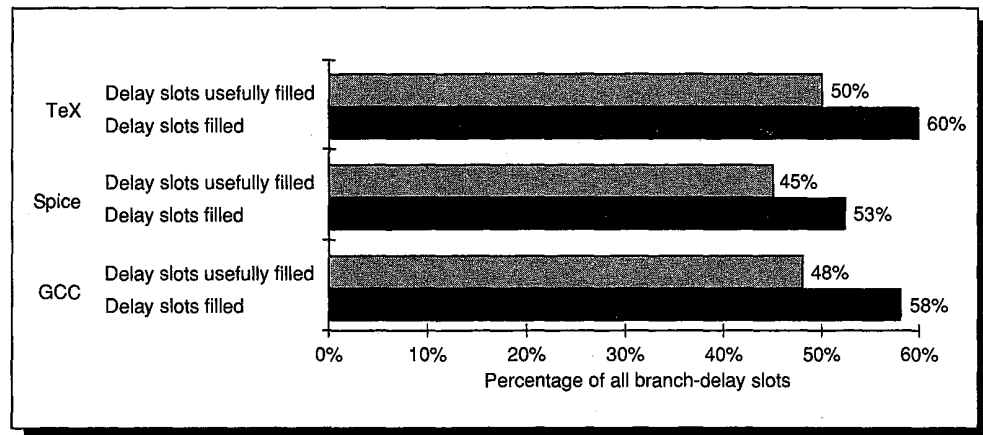


FIGURE 6.22 Frequency with which a single branch-delay slot is filled and how often the instruction is useful to the computation. The solid bar shows the percentage of the branch-delay slots occupied by some instruction other than a no-op. The difference between 100% and the dark column represents those branches that are followed by a no-op. The shaded bar shows how often those instructions do useful work. The difference between the shaded and solid bars is the percentage of instructions executed in a branch delay but not contributing to the computation. These instructions occur because optimization (b) is only useful when the branch is taken. If optimization (c) were used it would also contribute to this difference, since it is only useful when the branch is not taken.

slightly more than half the branch-delay slots are filled, and most of the filled slots do useful work. On average about 80% of the filled delay slots contribute to the computation. This number seems surprising, since branches are only taken about 53% of the time. The success rate is high because about one-half of the branch delays are being filled with an instruction from before the branch (strategy (a)), which is useful independent of whether the branch is taken.

When the scheduler in Figure 6.22 cannot use strategy (a)—moving an instruction from before the branch to fill the branch-delay slot—it uses only strategy (b)—moving it from the target. (For simplicity reasons, the schedule does not use strategy (c).) In total, nearly half the branch-delay slots are dynamically useful, eliminating one-half the branch stalls. Looking at Figure 6.22 we see that the primary limitation is the number of empty slots—those filled with no-ops. It is unlikely that the ratio of useful slots to filled slots, about 80%, can be improved, since this would require much better accuracy in predicting branches. In the Exercises we consider an extension of the delayed-branch idea that tries to fill more slots.

There is a small additional hardware cost for delayed branches. Because of the delayed effect of branches, multiple PCs (one plus the length of the delay) are needed to correctly restore the state when an interrupt occurs. Consider when the interrupt occurs after a taken-branch instruction is completed, but before all the instructions in the delay slots and the branch target are completed. In this case, the PC's of the delay slots and the PC of the branch target must be saved, since they are not sequential.

What is the effective performance of each of these schemes? The effective pipeline speedup with branch penalties is

$$\text{Pipeline speedup} = \frac{\text{Ideal CPI} * \text{Pipeline depth}}{\text{Ideal CPI} + \text{Pipeline stall cycle}}$$

If we assume that the ideal CPI is 1, then we can simplify this:

$$\text{Pipeline speedup} = \frac{\text{Pipeline depth}}{1 + \text{Pipeline stall cycles from branches}}$$

Since

$$\text{Pipeline stall cycles from branches} = \text{Branch frequency} * \text{Branch penalty}$$

we obtain:

$$\text{Pipeline speedup} = \frac{\text{Pipeline depth}}{(1 + \text{Branch frequency} * \text{Branch penalty})}$$

Using the DLX measurements in this section, Figure 6.23 shows several hardware options for dealing with branches, along with their performances (assuming a base CPI of 1).

Scheduling scheme	Branch penalty	Effective CPI	Pipeline speedup over nonpipelined machine	Pipeline speedup over stall pipeline on branch
Stall pipeline	3	1.42	3.52	1.00
Predict taken	1	1.14	4.39	1.25
Predict not taken	1	1.09	4.59	1.30
Delayed branch	0.5	1.07	4.67	1.33

FIGURE 6.23 Overall costs of a variety of branch schemes with the DLX pipeline. These data are for our DLX pipeline using the measured control-instruction frequency of 14% and the measurements of delay-slot filling from Figure 6.22. In addition, we know that 65% of the control instructions actually change the PC (taken branches plus unconditional changes). Shown are both the resultant CPI and the speedup over a nonpipelined machine, which we assume would have a CPI of 5 without any branch penalties. The last column of the table gives the speedup over a scheme that always stalls on branches.

Remember that the numbers in this section are **dramatically** affected by the length of the pipeline delay and the base CPI. A longer pipeline delay will cause an increase in the penalty and a larger percentage of wasted time. A delay of only one clock cycle is small—many machines have minimum delays of five or more. With a low CPI, the delay must be kept small, while a higher base CPI would reduce the relative penalty from branches.

Summary: Performance of the DLX Integer Pipeline

We close this section on hazard detection and elimination by showing the total distribution of idle clock cycles for our benchmarks when run on the DLX integer pipeline with software for pipeline scheduling. Figure 6.24 shows the distribution of clock cycles lost to load delays and branch delays in our three programs, by combining the separate measurements shown in Figures 6.13 (page 268) and 6.22.

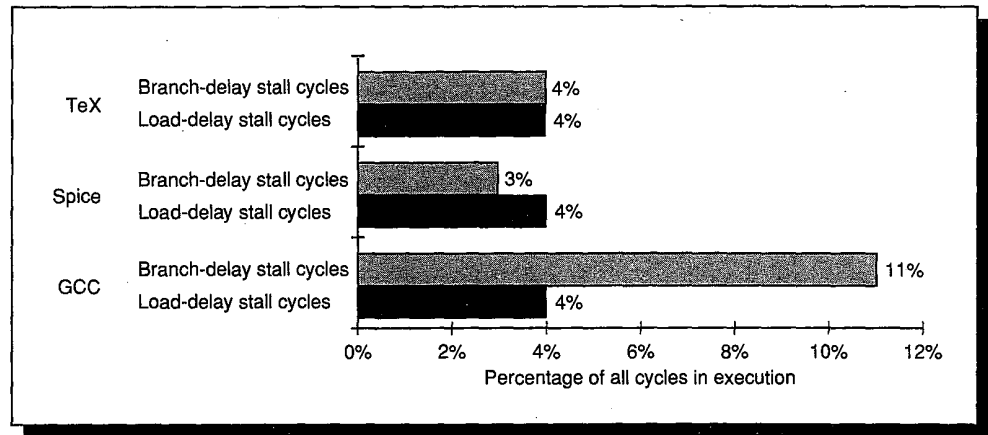


FIGURE 6.24 Percentage of the clock cycles spent on delays versus executing instructions. This assumes a perfect memory system; the clock-cycle count and instruction count would be identical if there were no integer pipeline stalls. This graph says that from 7% to 15% of the clock cycles are stalls; the remaining 85% to 93% are clock cycles that issue instructions. The Spice clock cycles do not include stalls in the FP pipeline, which will be shown at the end of Section 6.6. The pipeline scheduler fills load delays before branch delays and this affects the distribution of delay cycles.

For the GCC and TeX programs, the effective CPI (ignoring any stalls except those from pipeline hazards) on this pipelined version of DLX is 1.1. Compare this to the CPI for the complete nonpipelined, hardwired version of DLX described in Chapter 5 (Section 5.7), which is 5.8. Ignoring all other sources of stalls and assuming that the clock rates will be the same, the performance improvement from pipelining is 5.3 times.

6.5 What Makes Pipelining Hard to Implement

Now that we understand how to detect and resolve hazards, we can deal with some complications that we have avoided so far. In Chapter 5 we saw that interrupts are among the most difficult aspects of implementing a machine; pipelining increases that difficulty. In the second part of this section, we discuss some of the challenges raised by different instruction sets.

Dealing with Interrupts

Interrupts are harder to handle in a pipelined machine because the overlapping of instructions makes it more difficult to know whether an instruction can safely change the state of the machine. In a pipelined machine, an instruction is executed piece by piece and is not completed for several clock cycles. Yet in the process of executing it may need to update the machine state. Meanwhile, an interrupt can force the machine to abort the instruction's execution before it is completed.

As in nonpipelined implementations, the most difficult interrupts have two properties: (1) they occur within instructions, and (2) they must be restartable. In our DLX pipeline, for example, a virtual memory page fault resulting from a data fetch cannot occur until sometime in the MEM cycle of the instruction. By the time that fault is seen, several other instructions will be in execution. Since a page fault must be restartable and requires the intervention of another process, such as the operating system, the pipeline must be safely shut down and the state saved so that the instruction can be restarted in the correct state. This is usually implemented by saving the PC of the instruction (during IF) to restart it. If the restarted instruction is not a branch then we will continue to fetch the sequential successors and begin their execution in the normal fashion. If the restarted instruction is a branch, then we will evaluate the branch condition and begin fetching from either the target or the fall through. When an interrupt occurs, we can take the following steps to save the pipeline state safely:

1. Force a trap instruction into the pipeline on the next IF.
2. Until the trap is taken, turn off all writes for the faulting instruction and for all instructions that follow in the pipeline. This prevents any state changes for instructions that will not be completed before the interrupt is handled.
3. After the interrupt-handling routine in the operating system receives control, it immediately saves the PC of the faulting instruction. This value will be used to return from the interrupt later.

When we use delayed branches it is no longer possible to re-create the state of the machine with the single PC of the interrupted instruction, because the instructions in the pipeline may not be sequentially related. In particular, when the instruction that causes the interrupt is a branch-delay slot, and the branch was taken, then the instructions to restart are those in the slot plus the instruction at the branch target. The branch itself has completed execution and is not restarted. The addresses of the instructions in the branch-delay slot and the target are not sequential. So we need to save and restore a number of PCs that is one more than the length of the branch delay. This is done in the third step above.

After the interrupt has been handled, special instructions return the machine from the interrupt by reloading the PCs and restarting the instruction stream (using RFE in DLX). If the pipeline can be stopped so that the instructions just before the faulting instruction are completed and those after it can be restarted

from scratch, the pipeline is said to have *precise interrupts*. Ideally, the faulting instruction would not have changed the state, and correctly handling some interrupts requires that the faulting instruction have no effects. For other interrupts, such as floating-point exceptions, the faulting instruction on some machines writes its result before the interrupt can be handled. In such cases, the hardware must be prepared to retrieve the source operands, even if the destination is identical to one of the source operands.

Supporting precise interrupts is a requirement in many systems, while in others it is valuable because it simplifies the operating system interface. At a minimum, any machine with demand paging or IEEE arithmetic trap handlers must make its interrupts precise, either in the hardware or with some software support.

Precise interrupts are challenging because of the same problems that make instructions difficult to restart. As we saw in the last chapter, restarting is complicated by the fact that instructions can change the state of the machine before they are **guaranteed** to complete (sometimes called *committed* instructions). Because instructions in the pipeline may have dependences, not updating the machine state is impractical if the pipeline is to keep going. Thus, as a machine is more heavily pipelined, it becomes necessary to be able to back out of any state changes made before the instruction is committed (as discussed in Chapter 5). Fortunately, DLX has no such instructions, given the pipeline we have used.

Figure 6.25 (page 281) shows the DLX pipeline stages and which “problem” interrupts might occur in each stage. Because in pipelining there are multiple instructions in execution, multiple interrupts may occur on the same clock cycle. For example, consider this instruction sequence:

LW	IF	ID	EX	MEM	WB	
ADD		IF	ID	EX	MEM	WB

This pair of instructions can cause a data page fault and an arithmetic interrupt at the same time, since the LW is in MEM while the ADD is in EX. This case can be handled by dealing with only the data page fault and then restarting the execution. The second interrupt will reoccur (but not the first, if the software is correct), and when it does it can be handled independently.

In reality, the situation is not all this straightforward. Interrupts may occur out of order; that is, an instruction may cause an interrupt before an earlier instruction causes one. Consider again the above sequence of instructions LW; ADD. The LW can get a data page fault, seen when the instruction is in MEM, and the ADD can get an instruction page fault, seen when the ADD instruction is in IF. The instruction page fault will actually occur first, even though it is caused by a later instruction! This situation can be resolved in two ways. To explain them, let’s call the instruction in the position of the LW “instruction i ” and the instruction in the position of the ADD “instruction $i+1$.”

Pipeline stage	Problem interrupts occurring
IF	Page fault on instruction fetch; misaligned memory access; memory-protection violation
ID	Undefined or illegal opcode
EX	Arithmetic interrupt
MEM	Page fault on data fetch; misaligned memory access; memory-protection violation
WB	None

FIGURE 6.25 Interrupts from Chapter 5 that cause stop and restart of the DLX pipeline in a transparent fashion. The pipeline stage where these interrupts occur is also shown. Interrupts raised from instruction or data-memory access account for six out of seven cases. These interrupts and their corresponding names in other processors are in Figures 5.9 and 5.11.

The first approach is completely precise and is the simplest to understand for the user of the architecture. The hardware posts each interrupt in a status vector carried along with each instruction as it goes down the pipeline. When an instruction enters WB (or is about to leave MEM), the interrupt status vector is checked. If any interrupts are posted, they are handled in the order in which they would occur in time—the interrupt corresponding to the earliest instruction is handled first. This guarantees that all interrupts will be seen on instruction i before any are seen on $i+1$. Of course, any action taken on behalf of instruction i may be invalid, but because no state is changed until WB, this is not a problem in the DLX pipeline. Nevertheless, pipeline control may want to disable any actions on behalf of an instruction i (and its successors) as soon as the interrupt is recognized. For pipelines that could update state earlier than WB, this disabling is required.

The second approach is to handle an interrupt as soon as it appears. This could be regarded as slightly less precise because interrupts occur in an order different from the order they would occur in if there were no pipelining. Figure 6.26 shows two interrupts occurring in the DLX pipeline. Because the interrupt at instruction $i+1$ is handled when it appears, the pipeline must be stopped immediately without completing any instructions that have yet to change state. For the DLX pipeline, this will be $i-2$, $i-1$, i , and $i+1$, assuming the interrupt is recognized at the end of the IF stage of the ADD instruction. The pipeline is then restarted with instruction $i-2$. Since the instruction causing the interrupt can be any of $i-2$, ..., $i+1$, the operating system must determine which instruction faulted. This is easy to figure out if the type of interrupt and its corresponding pipe stage are known. For example, only $i+1$ (the ADD instruction) could get an instruction page fault at this point, and only $i-2$ could get a data page fault. After handling the fault for $i+1$ and restarting at $i-2$, the data page fault will be encountered on instruction i , which will cause i , ..., $i+3$ to be interrupted. The data page fault can then be handled.

Instruction $i-3$	IF	ID	EX	MEM	WB				
Instruction $i-2$		IF	ID	EX	MEM	WB			
Instruction $i-1$			IF	ID	EX	<i>MEM</i>	<i>WB</i>		
Instruction i (LW)				IF	ID	<i>EX</i>	<i>MEM</i>	<i>WB</i>	
Instruction $i+1$ (ADD)					IF	<i>ID</i>	<i>EX</i>	<i>MEM</i>	<i>WB</i>
Instruction $i+2$						<i>IF</i>	<i>ID</i>	<i>EX</i>	<i>MEM</i> <i>WB</i>

Instruction $i-3$	IF	ID	EX	MEM	WB					
Instruction $i-2$		IF	ID	EX	MEM	WB				
Instruction $i-1$			IF	ID	EX	MEM	WB			
Instruction i (LW)				IF	ID	EX	MEM	<i>WB</i>		
Instruction $i+1$ (ADD)					IF	ID	<i>EX</i>	<i>MEM</i>	<i>WB</i>	
Instruction $i+2$						IF	<i>ID</i>	<i>EX</i>	<i>MEM</i> <i>WB</i>	
Instruction $i+3$							<i>IF</i>	<i>ID</i>	<i>EX</i> <i>MEM</i>	
Instruction $i+4$								IF	ID	EX

FIGURE 6.26 The actions taken for interrupts occurring at different points in the pipeline and handled **immediately**. This shows the instructions interrupted when an instruction page fault occurs in instruction $i+1$ (in the top diagram), and a data page fault in instruction i in the bottom diagram. The pipe stages in bold are the cycles during which the interrupt is recognized. The pipe stages in italics are the instructions that will not be completed due to the interrupt, and will need to be restarted. Because the earliest effect of the interrupt is on the pipe stage after it occurs, instructions that are in the WB stage when the interrupt occurs will complete, while those that have not yet reached WB will be stopped and restarted.

Instruction Set Complications

Another set of difficulties arises from odd bits of state that may create additional pipeline hazards or may require extra hardware to save and restore. Condition codes are a good example of this. Many machines set the condition codes implicitly as part of the instruction. At first glance, this looks like a good idea, since condition codes decouple the evaluation of the condition from the actual branch. However, implicitly set condition codes can cause difficulties in making branches fast. They limit the effectiveness of branch scheduling because most operations will modify the condition code, making it hard to schedule instructions between the setting of the condition code and the branch. Furthermore, in machines with condition codes, the processor must decide when the branch condition is fixed. This involves finding out when the condition code has been set for the last time prior to the branch. On the VAX, most instructions set the condition code, so that an implementation will have to stall if it tries to determine the branch condition early. Alternatively, the branch condition can be evaluated by the branch late in the pipeline, but this still leads to a long branch delay. On the 360/370 many, but not all, instructions set the condition codes. Figure 6.27 shows how the situation differs on the DLX, the VAX, and the 360 for the fol-

lowing C code sequence, assuming that b and d are initially in registers R2 and R3 (and should not be destroyed):

```
a = b + d;
if (b==0) ...
```

DLX	VAX	IBM 360
ADD R1,R2,R3	ADDL3 a,R2,R3	LR R1,R2
...	...	AR R1,R3
SW a,R1	CL R2,0	ST a,R1
...	BEQL label	...
BEQZ R2,label		LTR R2,R2
		BZ label

FIGURE 6.27 Code sequence for the above two statements. Because the ADD computes the sum of b and d, and the branch condition depends only on b, an explicit compare (on R2) is needed on the VAX and 360. On DLX, the branch depends only on R2 and can be arbitrarily far away from it. (In addition the sw could be moved into the branch-delay slot.) On the VAX all ALU operations and moves set the condition codes, so that a compare must be right before the branch. On the 360, for this example the instruction load and test register (LTR) is used to set the condition code. However, most loads on the 360 do not set the condition codes; thus, a load (or a store) could be moved between the LTR and the branch.

Provided there is lots of hardware to spare, **all** instructions before the branch in the pipeline can be examined to decide when the branch is determined. Of course, architectures with explicitly set condition codes avoid this difficulty. However, pipeline control must still track the last instruction that sets the condition code to know when the branch condition is decided. In effect, the condition code must be treated as an operand requiring hazard detection for RAW hazards on branches, just as DLX must do on the registers.

A final thorny area in pipelining is multicycle operations. Imagine trying to pipeline a sequence of VAX instructions such as this:

```
MOVL R1,R2
ADDL3 42(R1),56(R1)+,@(R1)
SUBL2 R2,R3
MOVC3 @(R1)[R2],74(R2),R3
```

These instructions differ radically in the number of clock cycles they will require, from as low as one up to hundreds of clock cycles. They also require different numbers of data memory accesses, from zero to possibly hundreds. Data hazards are very complex and occur both between and within instructions.

The simple solution of making all instructions execute for the same number of clock cycles is unacceptable because it introduces an enormous number of hazards and bypass conditions, and makes an immensely long pipeline. Pipelining the VAX at the instruction level is difficult (as we will see in Section 6.9), but a clever solution was found by the VAX 8800 designers. They pipeline the microinstruction execution; because the microinstructions are simple (they look a lot like DLX), the pipeline control is much easier. While it is not clear that this approach can achieve quite as low a CPI as an instruction-level pipeline for the VAX, it is much simpler, possibly leading to a shorter clock cycle time.

Load/store machines that have simple operations with similar amounts of work pipeline more easily. If architects realize the relationship between instruction set design and pipelining, they can design architectures for more efficient pipelining. In the next section we will see how the DLX pipeline deals with long-running instructions.

6.6

Extending the DLX Pipeline to Handle Multicycle Operations

We now want to explore how our DLX pipeline can be extended to handle floating-point operations. This section concentrates on the basic approach and the design alternatives, and closes with some performance measurements of a DLX floating-point pipeline.

It is impractical to require that all DLX floating-point operations complete in one clock cycle, or even in two. Doing so would mean either accepting a slow clock or using enormous amounts of logic in the floating-point units, or both. Instead, the floating-point pipeline will allow for a longer latency for operations. This is easier to grasp if we imagine the floating-point instructions as having the same pipeline as the integer instructions, with two important changes. First, the EX cycle may be repeated as many times as needed to complete the operation; the number of repetitions can vary for different operations. Second, there may be multiple floating-point functional units. A stall will occur if the instruction to be issued will either cause a structural hazard for the functional unit it uses or cause a data hazard.

For this section let's assume that there are four separate functional units in our DLX implementation:

1. The main integer unit
2. FP and integer multiplier
3. FP adder
4. FP and integer divider

The integer unit handles all loads and stores to either register set, all the integer operations (except multiply and divide), and branches. For now we will also

assume that the execution stages of the other functional units are not pipelined, so that no other instruction using the functional unit may issue until the previous instruction leaves EX. Moreover, if an instruction cannot proceed to the EX stage, the entire pipeline behind that instruction will be stalled. Figure 6.28 shows the resulting pipeline structure. In the next section we will deal with schemes that allow the pipeline to progress when there are more functional units or when the functional units are pipelined.

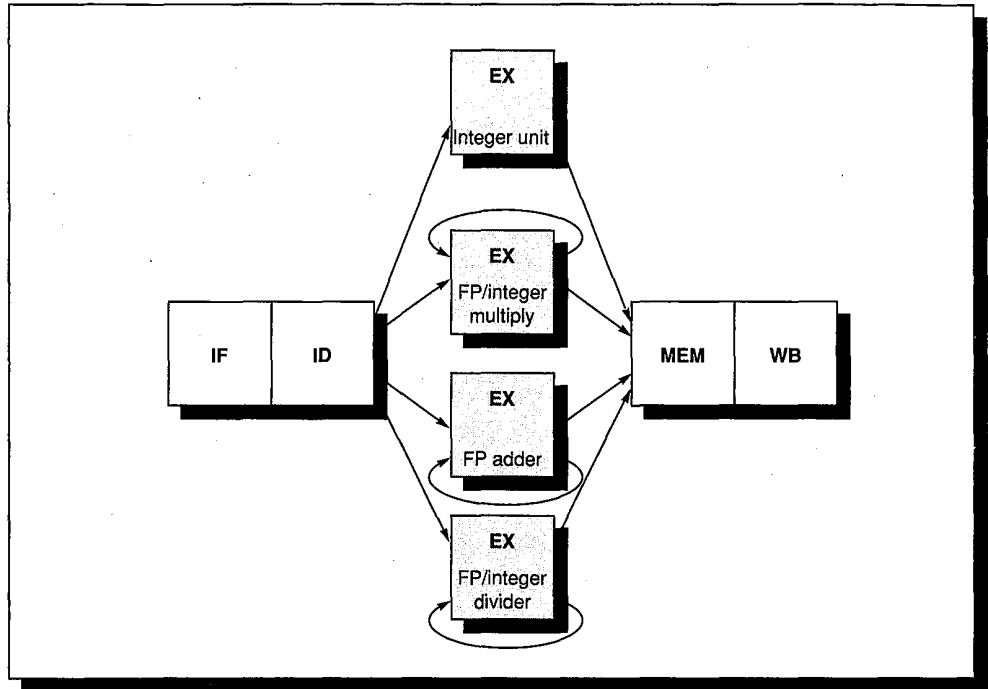


FIGURE 6.28 The DLX pipeline with three additional nonpipelined, floating-point, functional units. Because only one instruction issues on every clock cycle, all instructions go through the standard pipeline for integer operations. The floating-point operations simply loop when they reach the EX stage. After they have finished the EX stage, they proceed to MEM and WB to complete execution.

Since the EX stage may be repeated many times—30 to 50 repetitions for a floating-point divide would not be unreasonable—we must find a way to track long potential dependences and resolve hazards that last over tens of clock cycles, rather than just one or two. There is also the overlap between integer and floating-point instructions to deal with. However, overlapped integer and FP instructions do not complicate hazard detection, except on floating-point memory references and moves between the register sets. This is because, except for these memory references and moves, the FP and integer registers are distinct, and all integer instructions operate on the integer registers while the floating-point operations operate only on their own registers. This simplification of pipeline control is a major advantage of having separate register files for integer and floating-point data.

For now, let's assume that all floating-point operations take the same number of clock cycles—say 20 in the EX stage. What kind of hazard-detection circuitry will we need? Because all operations take the same amount of time, and register reads and writes always occur in the same stage, only RAW hazards are possible; no WAR or WAW hazards can occur. Thus, all we need to track is the destination register of each active functional unit. When we want to issue a new floating-point instruction, we take the following steps:

1. *Check for structural hazard*—Wait until the required functional unit is not busy.
2. *Check for a RAW data hazard*—Wait until the source registers are not listed as destinations by any of the EX stages in the functional units.
3. *Check for forwarding*—Test if the destination register of an instruction in MEM or WB is one of the source registers of the floating-point instruction; if so, enable the input multiplexer to use that result, rather than the register contents.

There is a small complication arising from conflicts between floating-point loads and floating-point operations when they both reach the WB stage simultaneously. We will deal presently with this situation in a more general fashion.

The above discussion assumes that the FP-functional-unit execution times were all the same. However, this does not hold up under practical scrutiny: Floating-point adds can typically be done in less than 5 clock cycles, multiplies in less than 10, and divides in about 20 or more. What we want is to allow the execution times of the functional units to differ, while still allowing the functional units to overlap execution. This would not change the basic structure of the pipeline in Figure 6.28, though it may cause the number of iterations around the loops to vary. Overlapping the execution of instructions whose running times differ, however, creates three complications: contention for register access at the end of the pipeline, the possibility of WAR and WAW hazards, and greater difficulty in providing precise interrupts.

We have already seen that FP loads and FP operations can contend for the floating-point register file on writes. When floating-point operations vary in execution time, they can also collide when trying to write results. This problem can be resolved by establishing a static priority for use of the WB stage. If multiple instructions wish to enter the MEM stage simultaneously, all instructions except the one with the highest priority are stalled in their EX stage. A simple, though sometimes suboptimal, heuristic is to give priority to the unit with the longest latency, since that is the one most likely to be the cause of the bottleneck. Although this scheme is reasonably simple to implement, this change to the DLX pipeline is quite significant. In the integer pipeline, all hazards were checked before the instruction issued to the EX stage. With this scheme for determining access to the result write port, instructions can stall after they issue.

Overlapping instructions with different execution times could introduce WAR and WAW hazards into our DLX pipeline, because the time at which

instructions write is no longer fixed. If all instructions still read their registers at the same time, no WAR hazards will be introduced.

WAW hazards are introduced because instructions can write their results in a different order than they appear. For example, consider the following code sequence:

```
DIVF    F0, F2, F4
SUBF    F0, F8, F10
```

A WAW hazard occurs between the divide and the subtract operations: The subtract will complete first, writing its result before the divide writes its result. Note that this hazard only occurs when the result of the divide will be overwritten **without** any instruction ever using it! If there were a use of F0 between the DIVF and the SUBF, the pipeline would stall because of a data dependence, and the SUBF would not issue until the DIVF was completed. We could argue that, for our pipeline, WAW hazards only occur when a useless instruction is executed, but we must still detect them and make sure that the result of the SUBF appears in F0 when we are done. (As we will see in Section 6.10, such sequences sometimes do occur in reasonable code.)

There are two possible ways to handle this WAW hazard. The first approach is to delay the issue of the subtract instruction until the DIVF enters MEM. The second approach is to stamp out the result of the divide by detecting the hazard and telling the divide unit not to write its result. Then, the SUBF can issue right away. Because this hazard is rare, either scheme will work fine—you can pick whatever is simpler to implement. As a pipeline gets more complex, however, we will need to devote increasing resources to determining when an instruction can issue.

Another problem caused by these long-running instructions can be illustrated with a very similar sequence of code:

```
DIVF    F0, F2, F4
ADDF    F10, F10, F8
SUBF    F12, F12, F14
```

This code sequence looks straightforward; there are no dependences. The problem with which we are concerned arises because an instruction issued early may complete after an instruction issued later. In this example, we can expect ADDF and SUBF to complete **before** the DIVF completes. This is called *out-of-order completion* and is common in pipelines with long-running operations. Since hazard detection will prevent any dependence among instructions from being violated, why is out-of-order completion a problem? Suppose that the SUBF causes a floating-point arithmetic interrupt at a point where the ADDF has completed but the DIVF has not. The result will be an imprecise interrupt, something we are trying to avoid. It may appear that this could be handled by letting the floating-point pipeline drain, as we do for the integer pipeline. But the interrupt may be in a position where this is not possible. For example, if the

DIVF decided to take a floating-point–arithmetic interrupt after the add completed, we could not have a precise interrupt at the hardware level. In fact, since the ADDF destroys one of its operands, we could not restore the state to what it was before the DIVF, even with software help.

This problem is being created because instructions are completing in a different order from the order in which they were issued. There are four possible approaches to dealing with out-of-order completion. The first is to ignore the problem and settle for imprecise interrupts. This approach was used in the 1960s and early 1970s. It is still used in some supercomputers, where certain classes of interrupts are not allowed or are handled by the hardware without stopping the pipeline. But it is difficult to use this approach in most machines built today, due to features such as virtual memory and the IEEE floating-point standard, which essentially require precise interrupts, through a combination of hardware and software.

A second approach is to queue the results of an operation until all the operations that were issued earlier are complete. Some machines actually use this solution, but it becomes expensive when the difference in running times among operations is long, since the number of results to queue can become large. Furthermore, results from the queue must be bypassed so as to continue issuing instructions while waiting for the longer instruction. This requires a large number of comparators and a very large multiplexer. There are two viable variations on this basic approach. The first is a *history file*, used in the CYBER 180/990. The history file keeps track of the original values of registers. When an interrupt occurs and the state must be rolled back earlier than some instruction that completed out of order, the original value of the register can be restored from the history file. A similar technique is used for autoincrement and autodecrement addressing on machines like VAXes. Another approach, the *future file*, proposed by J. Smith and Plezkun [1988], keeps the newer value of a register; when all earlier instructions have completed, the main register file is updated from the future file. On an interrupt, the main register file has the precise values for the interrupted state.

A third technique in use is to allow the interrupts to become somewhat imprecise, but keep enough information so that the trap-handling routines can create a precise sequence for the interrupt. This means knowing what operations were in the pipeline and their PCs. Then, after handling a trap, the software finishes any instructions that precede the latest instruction completed, and the sequence can restart. Consider the following worst-case code sequence:

Instruction₁—a long-running instruction that eventually interrupts execution

Instruction₂, ..., instruction_{*n*-1}—a series of instructions that are not completed

Instruction_{*n*}—an instruction that is finished

Given the PCs of all the instructions in the pipeline and the interrupt return PC, the software can find the state of instruction₁ and instruction_{*n*}. Since instruction_{*n*} has completed, we will want to restart execution at instruction_{*n*+1}. After

handling the interrupt, the software must simulate the execution of instruction₁, ... , instruction_{n-1}. Then we can return from the interrupt and restart at instruction_{n+1}. The complexity of executing these instructions properly by the handler is the major difficulty of this scheme. There is an important simplification: If instruction₂, ... , instruction_n are all integer instructions, then we know that if instruction_n has completed, all of instruction₂, ... , instruction_{n-1} have also completed. Thus, only floating-point operations need to be handled. To make this scheme tractable the number of floating-point instructions that can be overlapped in execution can be limited. For example, if we only overlap two instructions, then only the interrupting instruction need be completed by software. This restriction may reduce the potential throughput if the FP pipelines are deep or if there is a significant number of FP functional units. This approach is used in the SPARC architecture to allow overlap of floating-point and integer operations.

The final technique is a hybrid scheme that allows the instruction issue to continue only if it is certain that all the instructions before the issuing instruction will complete without causing an interrupt. This guarantees that when an interrupt occurs, no instructions after the interrupting one will be completed, and all of the instructions before the interrupting one can be completed. This sometimes means stalling the machine to maintain precise interrupts. To make this scheme work, the floating-point functional units must determine if an interrupt is possible early in the EX stage (in the first three clock cycles in the DLX pipeline), so as to prevent further instructions from completing. This scheme is used in the MIPS R2000/3000 architecture and is discussed further in Appendix A, Section A.7.

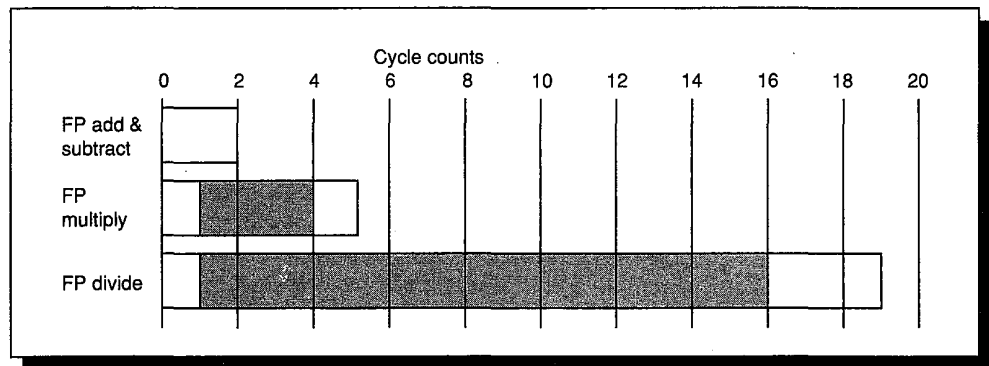


FIGURE 6.29 Total clock cycle count and permissible overlap among double-precision, floating-point operations on the MIPS R2010/3010 FP unit. The overall length of the bar shows the total number of EX cycles required to complete the operation. For example, after five clock cycles a multiply result is available. The shaded regions are times during which FP operations can be overlapped. As is common in most FP units, some of the FP logic is shared—the rounding logic, for example, is often shared. This means that FP operations with different running times cannot overlap arbitrarily. Also note that multiply and divide are not pipelined in this FP unit, so only one multiply or divide can be outstanding. The motivation for this pipeline design is discussed further in Appendix A (page A-31).

Performance of a DLX FP Pipeline

To look at the FP pipeline performance of DLX, we need to specify the latency and issue restrictions for the FP operations. We have chosen to use the pipeline structure of the MIPS R2010/3010 FP unit. While this unit has some structural hazards, it tends to have low-latency FP operations compared to most other FP units. The latencies and issue restrictions for DP floating-point operations are depicted in Figure 6.29 (page 289).

Figure 6.30 gives the breakdown of integer and floating-point stalls for Spice. There are four classes of stalls: load delays, branch delays, floating-point structural delays, and floating-point data hazards. The compiler tries to schedule both load and FP delays before it schedules branch delays. Interestingly, about 27% of the time in Spice is spent waiting for a floating-point result. Since the structural hazards are small, further pipelining of the floating-point unit would not gain much. In fact, the impact might easily be negative if the floating-point pipeline latency became longer.

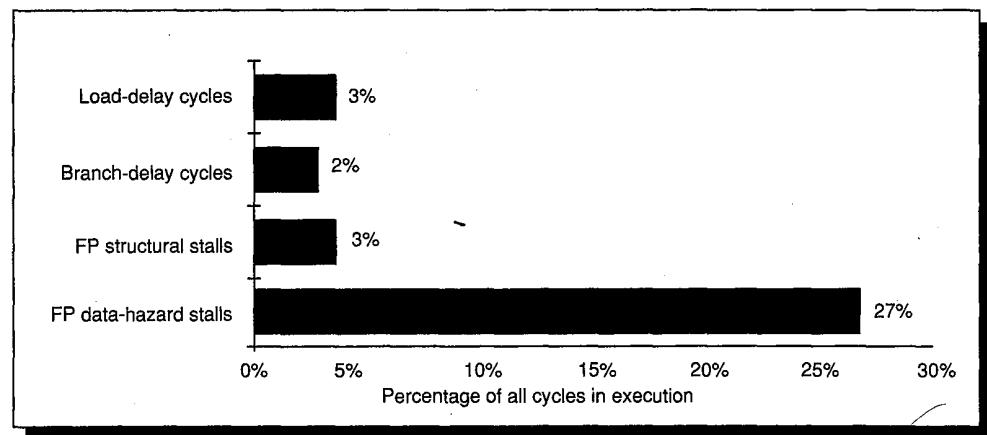


FIGURE 6.30 Percentage of clock cycles in Spice that are pipeline stalls. This again assumes a perfect memory system with no memory-system stalls. In total, 35% of the clock cycles in Spice are stalls, and without any stalls Spice would run about 50% faster. The percentage of stalls differs from Figure 6.24 (page 278) because this cycle count includes all the FP stalls, while the previous graph includes only the integer stalls.

6.7

Advanced Pipelining— Dynamic Scheduling in Pipelines

So far we have assumed that our pipeline fetches an instruction and issues it, unless there is a data dependence between an instruction already in the pipeline and the fetched instruction. If there is a data dependence, then we stall the instruction and cease fetching and issuing until the dependence is cleared. Software is responsible for scheduling the instructions to minimize these stalls. This

approach, which is called *static scheduling*, while first used in the 1960s, has become popular more recently. Many of the earlier, heavily pipelined machines used *dynamic scheduling*, whereby the hardware rearranges the instruction execution to reduce the stalls.

Dynamic scheduling offers a couple of advantages: It enables handling some cases when dependences are unknown at compile time, and it simplifies the compiler. It also allows code that was compiled with one pipeline in mind to run efficiently on a different pipeline. As we will see, these advantages are gained at a significant increase in hardware complexity. The first two parts of this section deal with reducing the cost of data dependences, especially in deeply pipelined machines. Corresponding to the dynamic hardware techniques for scheduling around data dependences are dynamic techniques for handling branches. These techniques are used for two purposes: to predict whether a branch will be taken, and to find the target more quickly. *Hardware branch prediction*, the name for these techniques, is the topic of the third part of this advanced section.

Dynamic Scheduling Around Hazards with a Scoreboard

The major limitation of the pipelining techniques we have used so far is that they all use in-order instruction issue. If an instruction is stalled in the pipeline, no later instructions can proceed. If there are multiple functional units, these units could lie idle. So, if instruction j depends on a long-running instruction i , currently in execution in the pipeline, then all instructions after j must be stalled until i is finished and j can execute. For example, consider this code:

```
DIVF   F0, F2, F4
ADDF   F10, F0, F8
SUBF   F6, F6, F14
```

The SUBF instruction cannot execute because the dependence of ADDF on DIVF causes the pipeline to stall; yet SUBF does not depend on anything in the pipeline. This is a performance limitation that can be eliminated by not requiring instructions to execute in order.

In the DLX pipeline, both structural and data hazards were checked at ID: When an instruction could execute properly, it was issued from ID. To allow us to begin executing the SUBF in the above example, we must separate the issue process into two parts: checking the structural hazards, and waiting for the absence of a data hazard. We can still check for structural hazards when we issue the instruction; thus, we still use in-order instruction issue. However, we want the instructions to begin execution as soon as their data operands are available. Thus, the pipeline will do *out-of-order execution*, which obviously implies *out-of-order completion*.

In introducing out-of-order execution, we have essentially split two pipe stages of DLX into three pipe stages. The two stages in DLX were:

1. ID—decode instruction, check for all hazards, and fetch operands
2. EX—execute instruction

In the DLX pipeline all instructions passed through issue stage in order, and a stalled instruction in ID caused a stall for all instructions behind it. The three stages we will need to allow out-of-order execution are:

1. Issue—decode instructions, check for structural hazards
2. Read operands—wait until no data hazards, then read operands
3. Execute

These three stages replace the ID and EX stages in the simple DLX pipeline.

While all instructions pass through the issue stage in order (in-order issue), they can be stalled or bypass each other in the second stage (read operands), and thus enter execution out of order. *Scoreboarding* is a technique for allowing instructions to execute out of order when there are sufficient resources and no data

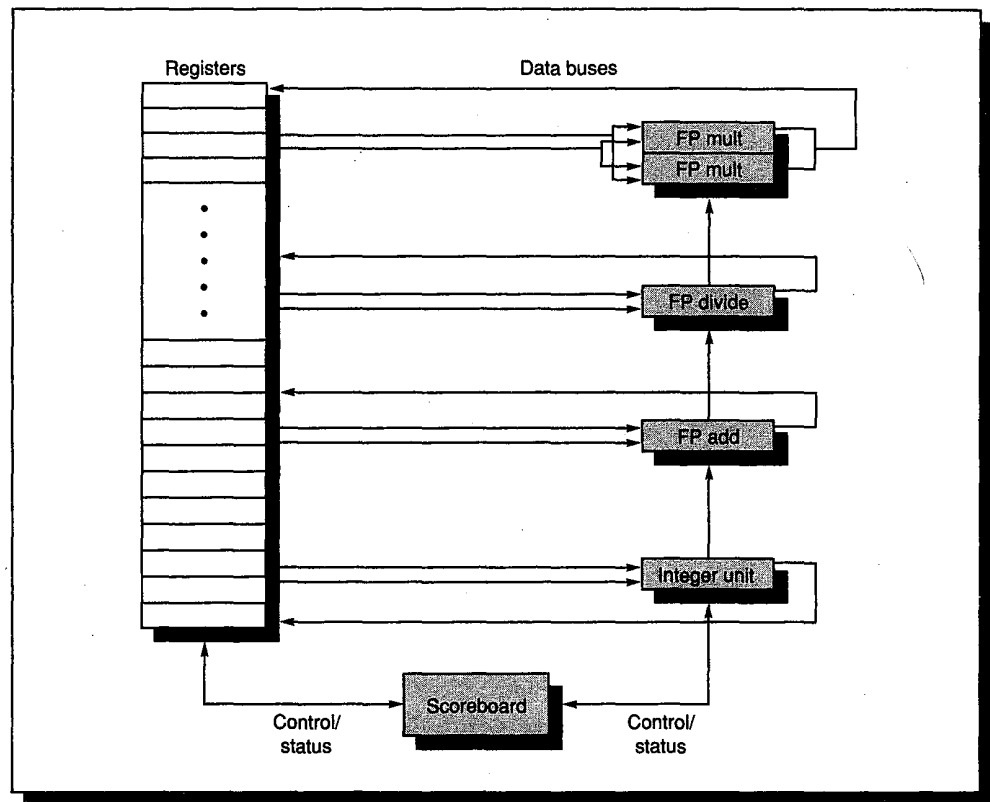


FIGURE 6.31 This shows the basic structure of a DLX machine with a scoreboard. The scoreboard's function is to control instruction execution (vertical control lines). All data flows between the register file and the functional units over the buses (the horizontal lines, called trunks in the CDC 6600). There are two FP multipliers, an FP divider, an FP adder, and an integer unit. One set of buses (two inputs and one output) serves a group of functional units. The details of the scoreboard are shown in Figures 6.32–6.35.

dependences; it is named after the CDC 6600 scoreboard, which developed this capability.

Before we see how scoreboarding could be used in the DLX pipeline, it is important to observe that WAR hazards, which did not exist in the DLX floating-point or integer pipelines, may exist when instructions are executed out of order. Assume our earlier example has changed so that the SUBF destination is F8. If ADDF and SUBF use two different functional units, then it is possible to execute the SUBF before the ADDF, but it will yield an incorrect result if ADDF has not read F8 before SUBF writes its result. The hazard for this case can be avoided by two rules: (1) read registers only during Read Operands, and (2) queue both the ADDF operation **and** copies of its operands. Of course, WAW hazards must still be detected, such as would occur if the destination of the SUBF were F10. This WAW hazard can be eliminated by stalling the issue of the SUBF instruction.

The goal of a scoreboard is to maintain an execution rate of one instruction per clock cycle (when there are no structural hazards) by executing an instruction as early as possible. Thus, when the instruction at the front of the queue is stalled, other instructions can be issued and executed if they do not depend on any active or stalled instruction. The scoreboard takes full responsibility for instruction issue and execution, including all hazard detection. Taking advantage of out-of-order execution requires multiple instructions to be in their EX stage simultaneously. This can be achieved with either multiple functional units or with pipelined functional units. Since these two capabilities—pipelined functional units and multiple functional units—are essentially equivalent for the purposes of pipeline control, we will assume the machine has multiple functional units.

The CDC 6600 had 16 separate functional units, including 4 floating-point units, 5 units for memory references, and 7 units for integer operations. On DLX, scoreboards make sense only on the floating-point unit. Let's assume that there are two multipliers, one adder, one divide unit, and a single integer unit for all memory references, branches, and integer operations. Although this example is much smaller than the CDC 6600, it is sufficiently powerful to demonstrate the principles. Because both DLX and the CDC 6600 are load/store, the techniques are nearly identical for the two machines. Figure 6.31 shows what the machine looks like.

Every instruction goes through the scoreboard, where a picture of the data dependences is constructed; this step corresponds to instruction issue and replaces part of the ID step in the DLX pipeline. This picture then determines when the instruction can read its operands and begin execution. If the scoreboard decides the instruction cannot execute immediately, it monitors every change in the hardware and decides when the instruction can execute. The scoreboard also controls when an instruction can write its result into the destination register. Thus, all hazard detection and resolution is centralized in the scoreboard. We will see a picture of the scoreboard later (Figure 6.32 on page 296), but first we need to understand the steps in the issue and execution segment of the pipeline.

Each instruction undergoes four steps in executing. (Since we are concentrating on the FP operations, we will not consider a step for memory access.) Let's first examine the steps informally and then look in detail at how the scoreboard keeps the necessary information that determines when to progress from one step to the next. The four steps, which replace the ID, EX, and WB steps in the standard DLX pipeline, are as follows:

1. **Issue**—If a functional unit for the instruction is free and no other active instruction has the same destination register, the scoreboard issues the instruction to the functional unit and updates its internal data structure. By ensuring that no other active functional unit wants to write its result into the destination register, we guarantee that WAW hazards cannot be present. If a structural or WAW hazard exists, then the instruction issue stalls, and no further instructions will issue until these hazards are cleared. This step replaces a portion of the ID step in the DLX pipeline.
2. **Read operands**—The scoreboard monitors the availability of the source operands. A source operand is available if no active instruction is going to write it, or if the register containing the operand is being written by a currently active functional unit. When the source operands are available, the scoreboard tells the functional unit to proceed to read the operands from the registers and begin execution. The scoreboard resolves RAW hazards dynamically in this step, and instructions may be sent into execution out of order. This step, together with Issue, completes the function of the ID step in the simple DLX pipeline.
3. **Execution**—The functional unit begins execution upon receiving operands. When the result is ready, it notifies the scoreboard that it has completed execution. This step replaces the EX step in the DLX pipeline and takes multiple cycles in the DLX FP pipeline.
4. **Write result**—Once the scoreboard is aware that the functional unit has completed execution, the scoreboard checks for WAR hazards. A WAR hazard exists if there is a code sequence like the following:

```

DIVF  F0, F2, F4
ADDF  F10, F0, F8
SUBF  F8, F8, F14

```

ADDF has a source operand F8, which is the same register as the destination of SUBF. But ADDF actually depends on an earlier instruction. The scoreboard will still stall the SUBF until ADDF reads its operands. In general, then, a completing instruction cannot be allowed to write its results when

- there is an instruction that has not read its operands,
- one of the operands is the same register as the result of the completing instruction, and
- the other operand was the result of an earlier instruction.

If this WAR hazard does not exist, or when it clears, the scoreboard tells the functional unit to store its result to the destination register. This step replaces the WB step in the simple DLX pipeline.

Based on its own data structure, the scoreboard controls the instruction progression from one step to the next by communicating with the functional units. But there is a small complication: There is only a limited number of source operands and result buses to the register file. The scoreboard must guarantee that the number of functional units allowed to proceed into steps 2 and 4 do not exceed the number of buses available. We will not go into further detail on this, other than to mention that the CDC 6600 solved this problem by grouping the 16 functional units together into four groups and supplying a set of buses, called *data trunks*, for each group. Only one unit in a group could read its operands or write its result during a clock.

Now let's look at the detailed data structure maintained by a DLX scoreboard with five functional units. Figure 6.32 (page 296) shows what the scoreboard's information looks like for a simple sequence of instructions:

LF	F6, 34 (R2)
LF	F2, 45 (R3)
MULTF	F0, F2, F4
SUBF	F8, F6, F2
DIVF	F10, F0, F6
ADDF	F6, F8, F2

There are three parts to the scoreboard:

1. Instruction status—Indicates which of the four steps the instruction is in.
2. Functional unit status—Indicates the state of the functional unit (FU). There are nine fields for each functional unit:
 - Busy—Indicates whether the unit is busy or not
 - Op—Operation to perform in the unit (e.g., add or subtract)
 - Fi—Destination register
 - Fj,Fk—Source-register numbers
 - Qj,Qk—Number of the units producing source registers Fj, Fk
 - Rj,Rk—Flags indicating when Fj, Fk are ready; fields are reset when new values are read so that the scoreboard knows that the source operand has been read (this is required to handle WAR hazards)
3. Register result status—Indicates which functional unit will write a register, if an active instruction has the register as its destination.

Instruction status										
Instruction		Issue	Read operands	Execution complete	Write result					
LF	F6, 34 (R2)	√	√	√	√					
LF	F2, 45 (R3)	√	√	√						
MULTF	F0, F2, F4	√								
SUBF	F8, F6, F2	√								
DIVF	F10, F0, F6	√								
ADDF	F6, F8, F2									

Functional unit status										
FU no.	Name	Busy	Op	Fi	Fj	Fk	Qj	Qk	Rj	Rk
1	Integer	Yes	Load	F2	R3				No	No
2	Mult1	Yes	Mult	F0	F2	F4	1		No	Yes
3	Mult2	No								
4	Add	Yes	Sub	F8	F6	F2		1	Yes	No
5	Divide	Yes	Div	F10	F0	F6	2		No	Yes

Register result status										
	F0	F2	F4	F6	F8	F10	F12	...	F30	
FU no.	2	1			4	5				

FIGURE 6.32 Components of the scoreboard. Each instruction that has issued or is pending issue has an entry in the instruction-status table. There is one entry in the functional-unit-status table for each functional unit. Once an instruction issues, the record of its operands is kept in the functional-unit-status table. Finally, the register-result table indicates which unit will produce each pending result; the number of entries is equal to the number of registers. The instruction-status register says that (1) the first LF has completed and written its result, and (2) the second LF has completed execution but has not yet written its result. The MULTF, SUBF, and DIVF have all issued but are stalled, waiting for their operands. The functional-unit status says that the first multiply unit is waiting for the integer unit, the add unit is waiting for the integer unit, and the divide unit is waiting for the first multiply unit. The ADDF instruction is stalled due to a structural hazard; it will clear when the SUBF completes. If an entry in one of these scoreboard tables is not being used, it is left blank. For example, the Rk field is not used on a load, and the Mult2 unit is unused, hence its fields have no meaning. Also, once an operand has been read, the Rj and Rk fields are set to No. These are left blank to minimize the complexity of the tables.

Now let's look at how the code sequence begun in Figure 6.32 continues execution. After that, we will be able to examine in detail the conditions that the scoreboard uses to control execution.

Example

Assume the following EX cycle latencies for the floating-point functional units: Add is 2 clock cycles, multiply is 10 clock cycles, and divide is 40 clock cycles. Using the code segment in Figure 6.32, and beginning with the point indicated by the instruction status in Figure 6.32, show what the status tables look like when MULTF and DIVF are each ready to go to the write-result state.

Answer

There are RAW data hazards from the second LF to MULTF and SUBF, from MULTF to DIVF, and from SUBF to ADDF. There is a WAR data hazard between DIVF and ADDF. Finally, there is a structural hazard on the add functional unit for ADDF. What the tables look like when MULTF and DIVF are ready to go to write result are shown in Figures 6.33 and 6.34, respectively.

Instruction status				
Instruction	Issue	Read operands	Execution complete	Write result
LF F6, 34 (R2)	√	√	√	√
LF F2, 45 (R3)	√	√	√	√
MULTF F0, F2, F4	√	√	√	
SUBF F8, F6, F2	√	√	√	√
DIVF F10, F0, F6	√			
ADDF F6, F8, F2	√	√	√	

Functional unit status										
FU no.	Name	Busy	Op	Fi	Fj	Fk	Qj	Qk	Rj	Rk
1	Integer	No								
2	Mult1	Yes	Mult	F0	F2	F4			No	No
3	Mult2	No								
4	Add	Yes	Add	F6	F8	F2			No	No
5	Divide	Yes	Div	F10	F0	F6	2		No	Yes

Register result status									
	F0	F2	F4	F6	F8	F10	F12	...	F30
FU no.	2			4		5			

FIGURE 6.33 Scoreboard tables just before the MULTF goes to write result. The DIVF has not yet read its operands, since it has a dependence on the result of the multiply. The ADDF has read its operands and is in execution, although it was forced to wait until the SUBF finished to get the functional unit. ADDF cannot proceed to write result because of the WAR hazard on F6, which is used by the DIVF.

Instruction status										
Instruction	Issue	Read operands	Execution complete	Write result						
LF F6, 34 (R2)	√	√	√	√						
LF F2, 45 (R3)	√	√	√	√						
MULTF F0, F2, F4	√	√	√	√						
SUBF F8, F6, F2	√	√	√	√						
DIVF F10, F0, F6	√	√	√							
ADDF F6, F8, F2	√	√	√	√						

Functional unit status										
FU no.	Name	Busy	Op	Fi	Fj	Fk	Qj	Qk	Rj	Rk
1	Integer	No								
2	Mult1	No								
3	Mult2	No								
4	Add	No								
5	Divide	Yes	Div	F10	F0	F6			No	No

Register Result status										
	F0	F2	F4	F6	F8	F10	F12	...	F30	
FU no.										5

FIGURE 6.34 Scoreboard tables just before the **DIVF** goes to write result. **ADDF** was able to complete as soon as **DIVF** passed through read operands and got a copy of **F6**. Only the **DIVF** remains to finish.

Instruction status	Wait until	Bookkeeping
Issue	Not busy (FU) and not result(D)	Busy(FU) ← yes; Result(D) ← FU; Op(FU) ← op; Fi(FU) ← D; Fj(FU) ← S1; Fk(FU) ← S2; Qj ← Result(S1); Qk ← Result(S2); Rj ← not Qj; Rk ← not Qk
Read operands	Rj and Rk	Rj ← No; Rk ← No
Execution complete	Functional unit done	
Write result	$\forall f((Fj(f) \neq Fi(FU) \text{ or } Rj(f) = \text{No}) \& (Fk(f) \neq Fi(FU) \text{ or } Rk(f) = \text{No}))$	$\forall f(\text{if } Qj(f) = \text{FU} \text{ then } Rj(f) \leftarrow \text{Yes});$ $\forall f(\text{if } Qk(f) = \text{FU} \text{ then } Rk(f) \leftarrow \text{Yes});$ Result(Fi(FU)) ← Clear; Busy(FU) ← No

FIGURE 6.35 Required checks and bookkeeping actions for each step in instruction execution. FU stands for the functional unit used by the instruction, D is the destination register, S1 and S2 are the source registers, and op is the operation to be done. To access the scoreboard entry named F_j for functional unit FU we use the notation F_j(FU). Result(D) is the value of the result register field for register D. The test on the write-result case prevents the write when there is a WAR hazard. For simplicity we assume that all of the bookkeeping operations are done in one clock cycle.

Now we can see how the scoreboard works in detail by looking at what has to happen for the scoreboard to allow each instruction to proceed. Figure 6.35 shows what the scoreboard requires for each instruction to advance and the bookkeeping action necessary when the instruction does advance.

The costs and benefits of scoreboarding are an interesting question. The CDC 6600 designers measured a performance improvement of 1.7 for FORTRAN programs and 2.5 for hand-coded assembly language. However, this was measured in the days before software pipeline scheduling, semiconductor main memory, and caches (which lower memory-access time). The scoreboard on the CDC 6600 had about as much logic as one of the functional units, which is surprisingly low. The main cost was in the large number of buses—about four times as many as would be required if the machine only executed instructions in order (or if it only initiated one instruction per Execute cycle).

The scoreboard does not handle a few situations as well as it might. For example, when an instruction writes its result, a dependent instruction in the pipeline must wait for access to the register file because all results are written through the register file and never forwarded. This increases the latency and limits the ability of multiple instructions waiting for a result to initiate. WAW hazards would be very infrequent, so the stalls they cause are probably not a significant concern in the CDC 6600. However, in the next section we will see that dynamic scheduling offers the possibility of overlapping the execution of multiple iterations of a loop. To do this effectively requires a scheme for handling WAW hazards, which are likely to increase in frequency when multiple iterations are overlapped.

Another Dynamic Scheduling Approach— The Tomasulo Algorithm

Another approach to parallel execution around hazards was used by the IBM 360/91 floating-point unit. This scheme was credited to R. Tomasulo and is named after him. The IBM 360/91 was completed about three years after the CDC 6600, before caches appeared in commercial machines. IBM's goal was to achieve high floating-point performance from an instruction set and from compilers designed for the entire 360 computer family, rather than for only floating-point-intensive applications. Remember that the 360 architecture has only four double-precision floating-point registers, which limits the effectiveness of compiler scheduling; this fact was another motivation for the Tomasulo approach. Lastly, the IBM 360/91 had long memory accesses and long floating-point delays, which the Tomasulo algorithm was designed to overcome. At the end of the section, we will see that Tomasulo's algorithm can also support the overlapped execution of multiple iterations of a loop.

We will explain the algorithm, which focuses on the floating-point unit, in the context of a pipelined, floating-point unit for DLX. The primary difference between DLX and the 360 is the presence of register-memory instructions in the latter machine. Because Tomasulo's algorithm uses a load functional unit, no

significant changes are needed to add register–memory addressing modes; the primary addition is another bus. The IBM 360/91 also had pipelined functional units, rather than multiple functional units. The only difference between these is that a pipelined unit can start at most one operation per clock cycle. Since there are really no fundamental differences, we describe the algorithm as if there were multiple functional units. The IBM 360/91 could accommodate three operations for the floating-point adder and two for the floating-point multiplier. In addition, up to six floating-point loads, or memory references, and up to three floating-point stores could be outstanding. Load data buffers and store data buffers are used for this function. Although we will not discuss the load and store units, we do need to include the buffers for operands.

Tomasulo's scheme shares many ideas with the CDC 6600 scoreboard, so we assume the reader has understood the scoreboard thoroughly. There are, however, two significant differences. First, hazard detection and execution control are distributed—*reservation stations* at each functional unit control when an instruction can begin execution at that unit. This function is centralized in the scoreboard on the CDC 6600. Second, results are passed directly to functional units rather than going through the registers. The IBM 360/91 has a common result bus (called the *common data bus*, or CDB) that allows all units waiting for an operand to be loaded simultaneously. The CDC 6600 writes results into registers, where waiting functional units may have to contend for them. Also, the CDC 6600 has multiple completion buses (two in the floating-point unit), while the IBM 360/91 has only one.

Figure 6.36 shows the basic structure of a Tomasulo-based floating-point unit for DLX; none of the execution control tables are shown. The reservation stations hold instructions that have been issued and are awaiting execution at a functional unit, as well as the information needed to control the instruction once it has begun execution to the unit. The load buffers and store buffers hold data coming from and going to memory. The floating-point registers are connected by a pair of buses to the functional units and by a single bus to the store buffers. All results from the functional units and from memory are sent on the common data bus, which goes everywhere except to the load buffer. All the buffers and reservation stations have tag fields, employed by hazard control.

Before we describe the details of the reservation stations and the algorithm, let's look at the steps an instruction goes through—just as we did for the scoreboard. Since operands are transmitted differently than in a scoreboard, there are only three steps:

1. Issue—Get an instruction from the floating-point operation queue. If the operation is a floating-point operation, issue it if there is an empty reservation station, and send the operands to the reservation station if they are in the registers. If the operation is a load or store, it can issue if there is an available buffer. If there is not an empty reservation station or an empty buffer, then there is a structural hazard and the instruction stalls until a station or buffer is freed.

2. Execute—If one or more of the operands is not yet available, monitor the CDB while waiting for the register to be computed. This step checks for RAW hazards. When both operands are available, execute the operation.
3. Write result—When the result is available, write it on the CDB and from there into the registers and any functional units waiting for this result.

Although these steps are fundamentally similar to those in the scoreboard, there are three important differences. First, there is no checking for WAW and WAR hazards—these are eliminated as a byproduct of the algorithm, as we will see shortly. Second, the CDB is used to broadcast results rather than waiting on the registers. Third, the loads and stores are treated as basic functional units.

The data structures used to detect and eliminate hazards are attached to the reservation stations, the register file, and the load and store buffers. Although different information is attached to different objects, everything except the load

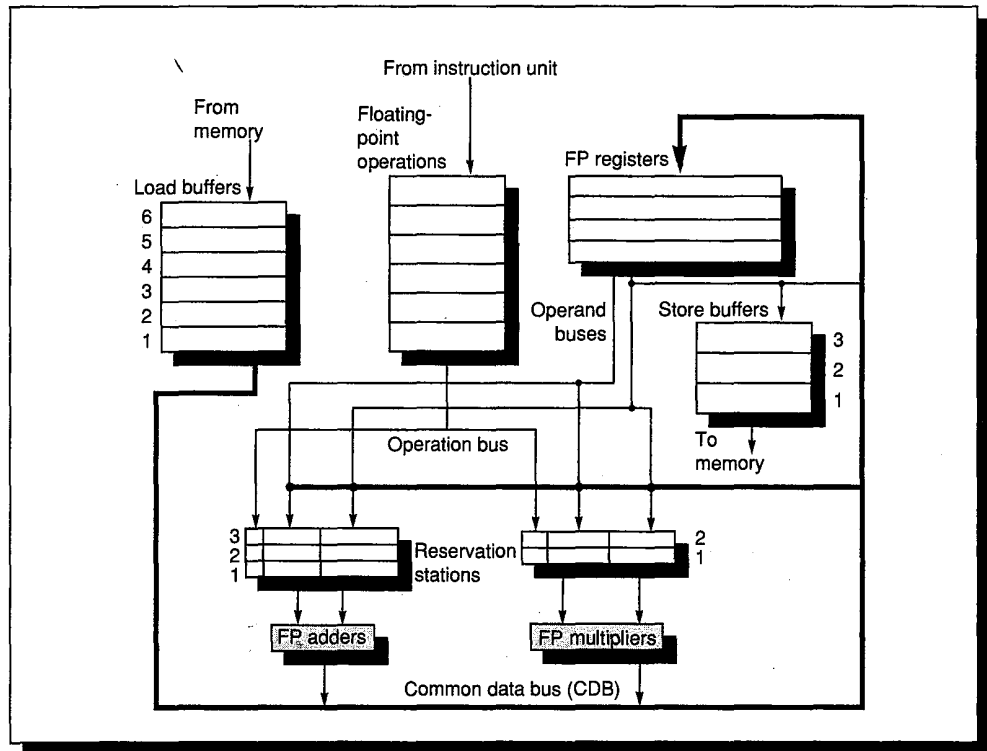


FIGURE 6.36 The basic structure of a DLX FP unit using Tomasulo's algorithm. Floating-point operations are sent from the instruction unit into a queue (called the FLOS, or floating-point operation stack, in the IBM 360/91) when they are issued. The reservation stations include the operation and the actual operands, as well as information used for detecting and resolving hazards. There are load buffers to hold the results of outstanding loads and store buffers to hold the addresses of outstanding stores waiting for their operands. All results from either the FP units or the load unit are put on the common data bus (CDB), which goes to the FP register file as well as the reservation stations and store buffers. The FP adders implement addition and subtraction, while the FP multipliers do multiplication and division.

buffers contains a tag field per entry. The tag field is a four-bit quantity that denotes one of the five reservation stations or one of the six load buffers. The tag field is used to describe which functional unit will produce a result needed as a source operand. Unused values, such as zero, indicate that the operand is already available. In describing the information, the scoreboard names are used wherever this will not lead to confusion. The names used by the IBM 360/91 are also shown. It is important to remember that the tags in the Tomasulo scheme refer to the buffer or unit that will produce a result; the register number is discarded when an instruction issues to a reservation station.

Each reservation station has six fields:

Op—The operation to perform on source operands S1 and S2.

Qj,Qk—The reservation stations that will produce the corresponding source operand; a value of zero indicates that the source operand is already available in Vi or Vj, or is unnecessary. The IBM 360/91 calls these SINKunit and SOURCEunit.

Vj,Vk—The value of the source operands. These are called SINK and SOURCE on the IBM 360/91. Note that only one of the V field or the Q field is valid for each operand.

Busy—Indicates that this reservation station and its accompanying functional unit are occupied.

The register file and store buffer each have a field, Qi:

Qi—The number of the functional unit that will produce a value to be stored into this register or into memory. If the value of Qi is zero, no currently active instruction is computing a result destined for this register or buffer. For a register, this means the value is given by the register contents.

The load and store buffers each require a busy field, indicating when a buffer is available due to completion of a load or store assigned there. The store buffer also has a field V, the value to be stored.

Before we examine the algorithm in detail, let's see what the system of tables looks like for the following code sequence:

1. LF F6, 34 (R2)
2. LF F2, 45 (R3)
3. MULTF F0, F2, F4
4. SUBF F8, F6, F2
5. DIVF F10, F0, F6
6. ADDF F6, F8, F2

We saw what the scoreboard looked like for this program when only the first load had written its result. Figure 6.37 depicts the reservation stations, load and

store buffers, and the register tags. The numbers appended to the names add, mult, and load stand for the tag for that reservation station—Add1 is the tag for the result from the first add unit. In addition we have included a central table called “Instruction status.” This table is included only to help the reader understand the algorithm; it is **not** actually a part of the hardware. Instead, the state of each operation that has issued is kept in a reservation station.

There are two important differences from scoreboards that are observable in these tables. First, the value of an operand is stored in the reservation station in one of the V fields as soon as it is available; it is not read from the register file once the instruction has issued. Second, the ADDF instruction has issued. This was blocked in the scoreboard by a structural hazard.

Instruction status				
Instruction		Issue	Execute	Write result
LF	F6, 34 (R2)	√	√	√
LF	F2, 45 (R3)	√	√	
MULTF	F0, F2, F4	√		
SUBF	F8, F6, F2	√		
DIVF	F10, F0, F6	√		
ADDF	F6, F8, F2	√		

Reservation stations						
Name	Busy	Op	Vj	Vk	Qj	Qk
Add1	Yes	SUB	(Load1)			Load2
Add2	Yes	ADD			Add1	Load2
Add3	No					
Mult1	Yes	MULT		(F4)	Load2	
Mult2	Yes	DIV		(Load1)	Mult1	

Register status									
Field	F0	F2	F4	F6	F8	F10	F12	...	F30
Qi	Mult1	Load2		Add2	Add1	Mult2			
Busy	Yes	Yes	No	Yes	Yes	Yes	No	...	No

FIGURE 6.37 Reservation stations and register tags. All of the instructions have issued, but only the first load instruction has completed and written its result to the CDB. The instruction-status table is not actually present, but the equivalent information is distributed throughout the hardware. The notation (X), where X is either a register number or a functional unit, indicates that this field contains the result of the functional unit X or the contents of register X at the time of issue. The other instructions are all at reservation stations or, as in the case of instruction 2, completing a memory reference. The load and store buffers are not shown. Load buffer 2 is the only busy load buffer and it is performing on behalf of instruction 2 in the sequence—loading from memory address R3 + 45. There are no stores, so the store buffer is not shown. Remember that an operand is specified by either the Q field or the V field at any time.

The big advantages of the Tomasulo scheme are (1) the distribution of the hazard detection logic, and (2) the elimination of stalls for WAW and WAR hazards. The first advantage arises from the distributed reservation stations and the use of the CDB. If multiple instructions are waiting on a single result, and each instruction already has its other operand, then the instructions can be released simultaneously by the broadcast on the CDB. In the scoreboard the waiting instructions must all read their results from the registers when register buses are available.

WAW and WAR hazards are eliminated by renaming registers using the reservation stations. For example, in our code sequence in Figure 6.37 we have issued both the `DIVF` and the `ADDF`, even though there is a WAR hazard involving `F6`. The hazard is eliminated in one of two ways. If the instruction providing the value for the `DIVF` has completed, then `Vk` will store the result, allowing `DIVF` to execute independent of the `ADDF` (this is the case shown). On the other hand, if the `LF` had not completed, then `Qk` would point to the `Load1` and the `DIVF` instruction would be independent of the `ADDF`. Thus, in either case, the `ADDF` can issue and begin executing. Any uses of the result of the `MULTF` would point to the reservation station, allowing the `ADDF` to complete and store its value into the registers without affecting the `DIVF`. We'll see an example of the elimination of a WAW hazard shortly. But let's first look at how our earlier example continues execution.

Example

Assume the same latencies for the floating-point functional units as we did for Figure 6.34: Add is 2 clock cycles, multiply is 10 clock cycles, and divide is 40 clock cycles. With the same code segment, show what the status tables look like when the `MULTF` is ready to go to write result.

Answer

The result is shown in the three tables in Figure 6.38. Unlike the example with the scoreboard, `ADDF` has completed since the operands of `DIVF` are copied, thereby overcoming the WAR hazard.

Instruction status				
Instruction		Issue	Execute	Write result
LF	F6, 34 (R2)	√	√	√
LF	F2, 45 (R3)	√	√	√
MULTF	F0, F2, F4	√	√	
SUBF	F8, F6, F2	√	√	√
DIVF	F10, F0, F6	√		
ADDF	F6, F8, F2	√	√	√

Reservation stations						
Name	Busy	Op	Vj	Vk	Qj	Qk
Add1	No					
Add2	No					
Add3	No					
Mult1	Yes	MULT	(Load2)	(F4)		
Mult2	Yes	DIV		(Load1)	Mult1	

Register status									
Field	F0	F2	F4	F6	F8	F10	F12	...	F30
Qi	Mult1					Mult2			
Busy	Yes	No	No	No	No	Yes	No	...	No

FIGURE 6.38 Multiply and divide are the only instructions not finished. This is different from the scoreboard case, because the elimination of WAR hazards allowed the ADDF to finish right after the SUBF on which it depended.

Figure 6.39 gives the steps for each instruction to go through. Load and stores are only slightly special. A load can be executed as soon as it is available. When execution is completed and the CDB is available, a load puts its result on the CDB like any functional unit. Stores receive their values from the CDB or from the register file and execute autonomously; when they are done they turn the busy field off to indicate availability, just like a load buffer or reservation station.

Instruction status	Wait until	Action or bookkeeping
Issue	Station or buffer empty	<pre> if (Register[S1].Qi ≠ 0) {RS[r].Qj ← Register[S1].Qi} else {RS[r].Vj ← S1; RS[r].Qj ← 0}; if (Register[S2].Qi ≠ 0) {RS[r].Qk ← Register[S2].Qi}; else {RS[r].Vk ← S2; RS[r].Qk ← 0} RS[r].Busy ← yes; Register[D].Qi = r; </pre>
Execute	(RS[r].Qj=0) and (RS[r].Qk=0)	None—operands are in Vj and Vk
Write result	Execution completed at <i>r</i> and CDB available	<pre> ∀x(if (Register[x].Qi=r) {Fx ← result; Register[x].Qi ← 0}); ∀x(if (RS[x].Qj=r) {RS[x].Vj ← result; RS[x].Qj ← 0}); ∀x(if (RS[x].Qk=r) {RS[x].Vk ← result; RS[x].Qk ← 0}); ∀x(if (Store[x].Qi=r) {Store[x].V ← result; Store[x].Qi ← 0}); RS[r].Busy ← No </pre>

FIGURE 6.39 Steps in the algorithm and what is required for each step. For the issuing instruction, D is the destination, S1 and S2 are the sources, and *r* is the reservation station or buffer that D is assigned to. RS is the reservation-station data structure. The value returned by a reservation station or by the load unit is called the “result.” Register is the register data structure, while Store is the store-buffer data structure. When an instruction is issued, the destination register has its Qi field set to the number of the buffer or reservation station to which the instruction is issued. If the operands are available in the registers, they are stored in the V fields. Otherwise, the Q fields are set to indicate the reservation station that will produce the values needed as source operands. The instruction waits at the reservation station until both its operands are available, indicated by zero in the Q fields. The Q fields are set to zero either when this instruction is issued, or when an instruction on which this instruction depends completes and does its write back. When an instruction has finished execution and the CDB is available, it can do its write back. All the buffers, registers, and reservation stations whose value of Qj or Qk is the same as the completing reservation station update their values from the CDB and mark the Q fields to indicate that values have been received. Thus, the CDB can broadcast its result to many destinations in a single clock cycle, and if the waiting instructions have their operands, they can all begin execution on the next clock cycle. For simplicity we assume that all bookkeeping actions are done in a single cycle.

To understand the full power of eliminating WAW and WAR hazards through dynamic renaming of registers, we must look at a loop. Consider the following simple sequence for multiplying the elements of a vector by a scalar in F2:

```

Loop: LD    F0, 0(R1)
      MULTD F4, F0, F2
      SD    0(R1), F4
      SUB   R1, R1, #8
      BNEZ  R1, Loop ; branches if R1≠0

```

With a branch-taken strategy, using reservation stations will allow multiple executions of this loop to proceed at once. This advantage is gained without unrolling the loop—in effect, the loop is unrolled dynamically by the hardware. In

the 360 architecture, the presence of only 4 FP registers would severely limit the use of unrolling. (We will see shortly, when we unroll a loop and schedule it to avoid interlocks, many more registers are required.) Tomasulo's algorithm supports the overlapped execution of multiple copies of the same loop with only a small number of registers used by the program.

Let's assume we have issued all the instructions in two successive iterations of the loop, but none of the floating-point loads/stores or operations has completed. The reservation stations, register-status tables, and load and store buffers at this point are shown in Figure 6.40. (The integer ALU operation is ignored, and it is assumed the branch was predicted as taken.) Once the system reaches this state, two copies of the loop could be sustained with a CPI close to one provided the multiplies could complete in four clock cycles. We will see how compiler techniques can achieve a similar result in Section 6.8.

An additional element that is critical to making Tomasulo's algorithm work is shown in this example. The load instruction from the second loop iteration could easily complete before the store from the first iteration, although the normal sequential order is different. The load and store can safely be done in a different order, provided the load and store access different addresses. This is checked by examining the addresses in the store buffer whenever a load is issued. If the load address matches the store-buffer address, we must stop and wait until the store buffer gets a value; we can then access it or get the value from memory.

This scheme can yield very high performance, provided the cost of branches can be kept small—this is a problem we will look at later in this section. There are also limitations imposed by the complexity of the Tomasulo scheme, which requires a large amount of hardware. In particular, there are many associative stores that must run at high speed, as well as complex control logic. Lastly, the performance gain is limited by the single completion bus (CDB). While additional CDBs can be added, each CDB must interact with all the pipeline hardware, including the reservation stations. In particular, the associative tag-matching hardware would need to be duplicated at all stations for each CDB.

While Tomasulo's scheme may be appealing if the designer is forced to pipeline an architecture that is difficult to schedule code for or has a shortage of registers, the authors believe that the advantages of the Tomasulo approach are limited for architectures that can be efficiently pipelined and statically scheduled with software. However, as available gate counts grow and the limits of software scheduling are reached, we may see dynamic scheduling employed. One possible direction is a hybrid organization that uses dynamic scheduling for loads and stores, while statically scheduling register-register operations.

Reducing Branch Penalties with Dynamic Hardware Prediction

The previous section describes techniques for overcoming data hazards. If control hazards are not addressed, Amdahl's Law predicts, they will limit pipelined-execution performance. Earlier, we looked at simple hardware schemes for

Instruction status					
Instruction		From iteration	Issue	Execute	Write result
LD	F0, 0 (R1)	1	√	√	
MULTD	F4, F0, F2	1	√		
SD	0 (R1), F4	1	√		
LD	F0, 0 (R1)	2	√	√	
MULTD	F4, F0, F2	2	√		
SD	0 (R1), F4	2	√		

Reservation stations						
Name	Busy	Fm	Vj	Vk	Qj	Qk
Add1	No					
Add2	No					
Add3	No					
Mult1	Yes	MULT		(F2)	Load1	
Mult2	Yes	MULT		(F2)	Load2	

Register status									
Field	F0	F2	F4	F6	F8	F10	F12	...	F30
Qi	Load2		Mult2						
Busy	yes	no	yes	no	no	no			

Store buffers			
Field	Store 1	Store 2	Store 3
Qi	Mult1	Mult2	
Busy	Yes	Yes	No
Address	(R1)	(R1)–8	

Load buffers			
Field	Load 1	Load 2	Load 3
Address	(R1)	(R1)–8	
Busy	Yes	Yes	No

FIGURE 6.40 Two active iterations of the loop with no instruction having yet completed. Load and store buffers are included, with addresses to be loaded from and stored to. The loads are in the load buffer; entries in the multiplier reservation stations indicate that the outstanding loads are the sources. The store buffers indicate that the multiply destination is their value to store.

dealing with branches (assume taken or not taken) and software-oriented approaches (delayed branches). This section focuses on using hardware to dynamically predict the outcome of a branch—the prediction will change if the branch changes its behavior while the program is running.

The simplest dynamic branch-prediction scheme is a *branch-prediction buffer*. A branch-prediction buffer is a small memory indexed by the lower por-

tion of the branch instruction address. The memory contains a bit that says whether the branch was recently taken or not. This is the simplest sort of buffer; it has no tags and is useful only to reduce the branch delay when it is longer than the time to compute the possible target PCs. We don't know, in fact, if the prediction is correct—it may have been put there by another branch that has the same low-order address bits. But this doesn't matter. It is assumed to be correct, and fetching begins in the predicted direction. If the branch prediction turns out to be wrong, the prediction bit is inverted.

This simple one-bit prediction scheme has a performance shortcoming: If a branch is almost always taken, then when it is not taken, we will predict incorrectly twice, rather than once. Consider a loop branch whose behavior is taken nine times sequentially, then not taken once. If the next time around it is predicted not taken, the prediction will be wrong. Thus, the prediction accuracy will only be 80%, even on branches that are 90% taken. To remedy this, two-bit prediction schemes are often used. In a two-bit scheme, a prediction must miss twice in a row before it is changed. Figure 6.41 shows the finite-state machine for the two-bit prediction scheme.

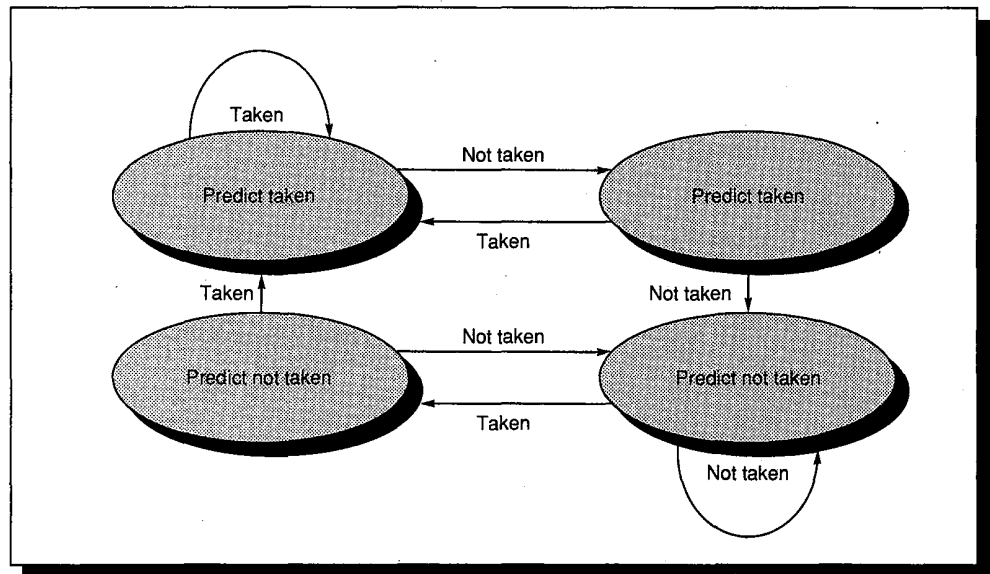


FIGURE 6.41 This shows the states in a two-bit prediction scheme. By using two bits rather than one, a branch that strongly favors taken or not taken—as many branches do—will be mispredicted only once. The two bits are used to encode the four states in the system.

The branch-prediction buffer can be implemented as a small, special cache accessed with the instruction address during the IF pipe stage, or as a pair of bits attached to each block in the instruction cache and fetched with the instruction (see Section 8.3 in Chapter 8). If the instruction is predicted as a branch and if the branch is predicted as taken, fetching begins from the target as soon as the

PC is known. Otherwise, fetching and sequential executing continue. If the prediction turns out to be wrong, the prediction bits are changed as shown in Figure 6.41. While this scheme is useful for most pipelines, the DLX pipeline finds out both whether the branch is taken and what the target of the branch is at the same time. Thus, this scheme does not help for the simple DLX pipeline; we will explore a scheme that can work for DLX a little later. First, let's see how well a prediction buffer works with a longer pipeline.

The accuracy of a two-bit prediction scheme is affected by how often the prediction for each branch is correct and by how often the entry in the prediction buffer matches the branch being executed. When the entry does not match, the prediction bit is used anyway because no better information is available. Even if the entry was for another branch, the guess could be a lucky one. In fact, there is about a 50% probability of being correct, even if the prediction is for some other branch. Studies of branch-prediction schemes have found that two-bit prediction has an accuracy of about 90% when the entry in the buffer is the branch entry. A buffer of between 500 and 1000 entries has a hit rate of 90%. The overall prediction accuracy is given by

$$\text{Accuracy} = (\% \text{ predicted correctly} * \% \text{ that prediction is for this instruction}) +$$

$$(\% \text{ lucky guess}) * (1 - \% \text{ that prediction is for this instruction})$$

$$\text{Accuracy} = (90\% * 90\%) + (50\% * 10\%) = 86\%$$

This number is higher than our success rate for filling delayed branches and would be useful in a pipeline with a longer branch delay. Now let's look at a dynamic prediction scheme that is useable for DLX and see how it compares to our branch-delay scheme.

To reduce the branch penalty on DLX, we need to know from what address to fetch by the end of IF. This means we must know whether the as yet undecoded instruction is a branch and, if it is a branch, what the next PC should be. If the instruction is a branch and we know what the next PC should be, we can have a branch penalty of zero. A branch-prediction cache that stores the predicted address for the next instruction after a branch is called a *branch-target buffer*. Because we are predicting the next instruction address and will send it out **before** decoding the instruction, we **must** know whether the fetched instruction is predicted as a taken branch. We also want to know whether the address in the target buffer is for a taken or not-taken prediction, so that we can reduce the time to find a mispredicted branch. Figure 6.42 shows what the branch-target buffer looks like. If the PC of the fetched instruction matches a PC in the buffer, then the corresponding predicted PC is used as the next PC. In Chapter 8 we will discuss caches in much more detail; we will see that the hardware for this branch-target buffer is similar to the hardware for a cache.

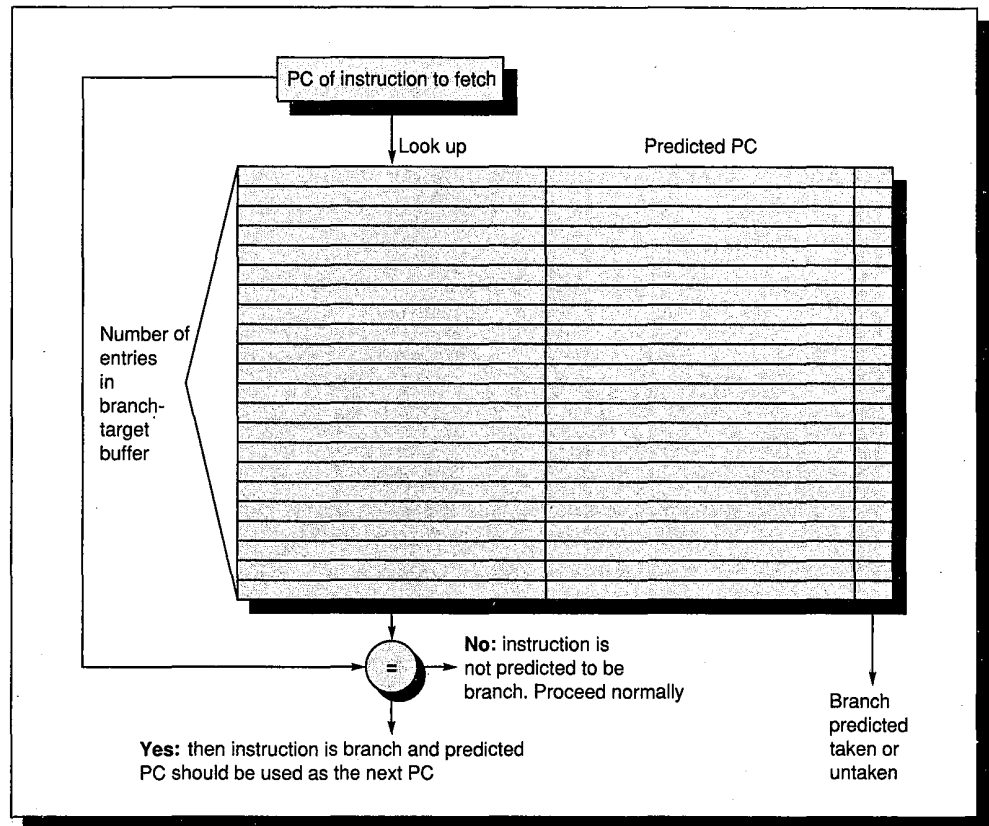


FIGURE 6.42 A branch-target buffer. The PC of the instruction being fetched is matched against a set of instruction addresses stored in the first column; these represent the addresses of known branches. If the PC matches one of these entries, then the instruction being fetched is a branch. If it is a branch, then the second field, predicted PC, contains the prediction for the next PC after the branch. Fetching begins immediately at that address. The third field just tracks whether the branch was predicted taken or untaken and helps keep the misprediction penalty small.

If a matching entry is found in the branch-target buffer, fetching begins immediately at the predicted PC. Note that (unlike a branch-prediction buffer) the entry must be for this instruction, because the predicted PC will be sent out before it is known whether this instruction is even a branch. If we did not check whether the entry matched this PC, then the wrong PC would be sent out for instructions that were not branches, resulting in a slower machine. Figure 6.43 shows the steps followed when using a branch-target buffer and when these steps occur in the pipeline. From this we can see that there will be no branch delay if a branch-prediction entry is found in the buffer and is correct. Otherwise, there will be a penalty of at least one clock cycle. In practice, there could be a penalty of two clock cycles because the branch-target buffer must be updated. We could assume that the instruction following a branch or at the branch target is not a branch, and do the update during that instruction time. However, this does complicate the control. Instead, we will take a two-clock-cycle penalty when the branch is not correctly predicted.

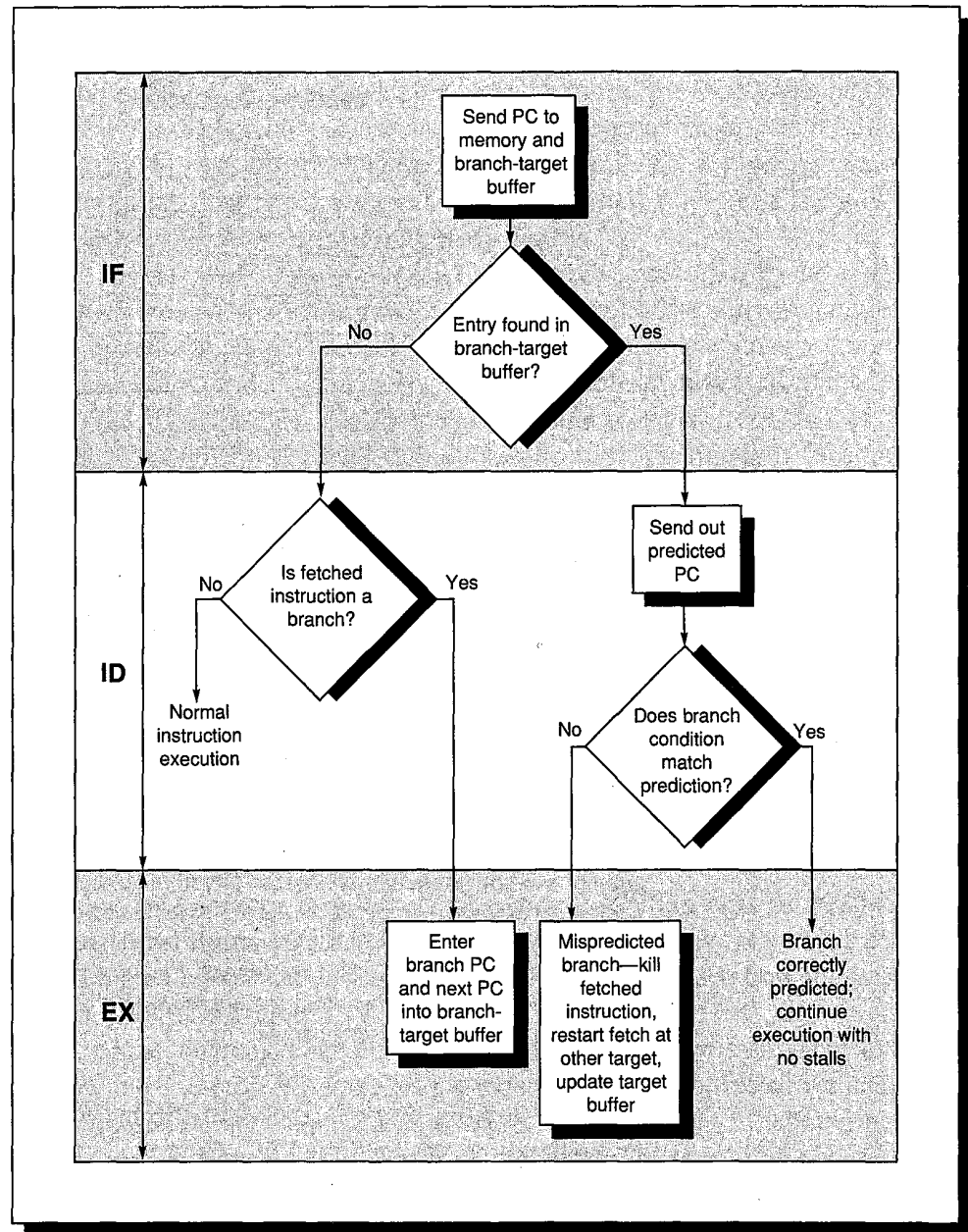


FIGURE 6.43 The steps involved in handling an instruction with a branch-target buffer. If the PC of an instruction is found in the buffer, then the instruction must be a branch, and fetching immediately begins from the predicted PC in ID. If the entry is not found and it subsequently turns out to be a branch, it is entered in the buffer along with the target, which is known at the end of ID. If the instruction is a branch, is found, and is correctly predicted, then execution proceeds with no delays. If the prediction is incorrect, we suffer a one-clock-cycle delay fetching the wrong instruction and restart the fetch one clock cycle later. If the branch is not found in the buffer and the instruction turns out to be a branch, we will have proceeded as if the instruction were a branch and can turn this into an assume-not-taken strategy; the penalty will differ depending on whether the branch is actually taken or not.

To evaluate how well a branch-target buffer works, we first must determine what the penalties are in all possible cases. Figure 6.44 contains this information.

Instruction in buffer	Prediction	Actual branch	Penalty cycles
Yes	Taken	Taken	0
Yes	Taken	Not taken	2
Yes	Not taken	Not taken	0
Yes	Not taken	Taken	2
No		Taken	2
No		Not taken	1

FIGURE 6.44 Penalties for all possible combinations of whether the branch is in the buffer, how it is predicted, and what it actually does. There is no branch penalty if everything is correctly predicted and the branch is found in the target buffer. If the branch is not correctly predicted, the penalty is equal to one clock cycle to update the buffer with the correct information (during which an instruction cannot be fetched) and one clock cycle, if needed, to restart fetching the next correct instruction for the branch. If the branch is not found and not taken, the penalty is only one clock cycle because the pipeline assumes not taken when it is not aware that the instruction is a branch. Other mismatches cost two clock cycles, since we must restart the fetch and update the buffer.

Using the same probabilities as for a branch-prediction buffer—90% probability of finding the entry and 90% probability of correct prediction—and the taken/not taken percentage taken from earlier in this chapter, we can find the total branch penalty:

$$\begin{aligned} \text{Branch penalty} &= \% \text{ branches found in buffer} * \% \text{ incorrect predictions} * 2 + \\ &\quad (1-\% \text{ branches found in buffer}) * \% \text{ taken branches} * 2 + \\ &\quad (1-\% \text{ branches found in buffer}) * \% \text{ untaken branches} * 1 \end{aligned}$$

$$\text{Branch penalty} = 90\% * 10\% * 2 + 10\% * 60\% * 2 + 10\% * 40\% * 1$$

$$\text{Branch penalty} = 0.34 \text{ clock cycles}$$

This compares with a branch penalty for delayed branches of about 0.5 clock cycles per branch. Remember, though, that the improvement from dynamic branch prediction will grow as the branch delay grows.

Branch-prediction schemes are limited both by prediction accuracy and by the penalty for misprediction. It is unlikely that we can improve the effective branch-prediction success much above 80% to 90%. Instead, we can try to reduce the penalty for misprediction. This is done by fetching from both the predicted and unpredicted direction. This requires that the memory system be dual ported or have an interleaved cache. While this adds cost to the system, it may be the only way to reduce branch penalties below a certain point.

We have seen a variety of software-based static schemes and hardware-based dynamic schemes for trying to boost the performance of our pipelined machine. Pipelining tries to exploit the potential for parallelism among sequential instructions. In the ideal case all the instructions would be independent, and our DLX pipeline would exploit parallelism among the five instructions simultaneously in the pipeline. Both the static scheduling techniques of the last section and the dynamic techniques of this section focus on maintaining the throughput of the pipeline at one instruction per clock. In the next section we will look at techniques that attempt to exploit overlap more than by the factor of 5, to which we are restricted with the simple DLX pipeline.

6.8

Advanced Pipelining—Taking Advantage of More Instruction-Level Parallelism

To improve performance further we would like to decrease the CPI to less than one. But the CPI cannot be reduced below one if we issue only one instruction every clock cycle. The goal of the techniques discussed in this section is to allow multiple instructions to issue in a clock cycle.

As we know from earlier sections, to keep a pipeline full, parallelism among instructions must be exploited by finding sequences of unrelated instructions that can be overlapped in the pipeline. Two related instructions must be separated by a distance equal to the pipeline latency of the first of the instructions. Throughout this section we will assume the latencies shown in Figure 6.45. Branches still have a one-clock-cycle delay. We assume that the functional units are fully pipelined or replicated, and that an operation can be issued on every clock cycle.

As we try to execute more instructions on every clock cycle and try to overlap more instructions, we will need to find and exploit more instruction-level parallelism. Thus, before looking at pipeline organizations that require more parallelism among instructions, let's look at a simple compiler technique that will help create additional parallelism.

Instruction producing result	Destination instruction	Latency in clocks
FP ALU op	Another FP ALU op	3
FP ALU op	Store double	2
Load double	FP ALU op	1
Load double	Store double	0

FIGURE 6.45 Latencies of operations used in this section. The first column shows the originating instruction type. The second column is the type of the consuming instruction. The last column is the separation in clock cycles to avoid a stall. These numbers are similar to the average latencies we would see on an FP unit, like the one we described for DLX in Figure 6.29 (page 289).

Increasing Instruction-Level Parallelism with Loop Unrolling

To compare the approaches discussed in this section, we will use a simple loop that adds a scalar value to a vector in memory. The DLX code, not accounting for the pipeline, looks like this:

```

Loop: LD    F0,0(R1)    ; load the vector element
      ADDD  F4,F0,F2    ; add the scalar in F2
      SD    0(R1),F4    ; store the vector element
      SUB   R1,R1,#8    ; decrement the pointer by
                          ; 8 bytes (per DW)
      BNEZ  R1,LOOP    ; branch when it's zero
    
```

For simplicity, we assume the array starts at location 0. If it were located elsewhere, the loop would require one additional integer instruction.

Let's start by seeing how well this loop will run when it is scheduled on a simple pipeline for DLX with the latencies discussed above.

Example

Show how the vector add loop would look on DLX, both scheduled and unscheduled, including any stalls or idle clock cycles.

Answer

Without any scheduling the loop will execute as follows:

	Clock cycle issued
Loop: LD F0,0(R1)	1
stall	2
ADDD F4,F0,F2	3
stall	4
stall	5
SD 0(R1),F4	6
SUB R1,R1,#8	7
BNEZ R1,LOOP	8
stall	9

This requires 9 clock cycles per iteration. We can schedule the loop to obtain

```

Loop: LD    F0,0(R1)
      stall
      ADDD  F4,F0,F2
      SUB   R1,R1,#8
      BNEZ  R1,LOOP    ; delayed branch
      SD    8(R1),F4    ; changed because interchanged with SUB
    
```

Execution time has been reduced from 9 clock cycles to 6.

Notice that to create this schedule, the compiler had to determine that it could swap the SUB and SD by changing the address the SD stored to: The address was 0 (R1) and is now 8 (R1). This is not trivial, since most compilers would see that the SD instruction depends on the SUB and would refuse to interchange them. A smarter compiler could figure out the relationship and perform the interchange. The dependence among the LD, ADDD, and SD determines the clock cycle count for this loop.

In the above example, we complete one loop iteration and finish one vector element every 6 clock cycles, but the actual work of operating on the vector element takes just 3 of those 6 clock cycles. The remaining 3 clock cycles consist of loop overhead—the SUB and BNEZ—and a stall. To eliminate these 3 clock cycles we need to get more operations within the loop. A simple scheme for increasing the number of instructions between executions of the loop branch is *loop unrolling*. This is done by simply replicating the loop body multiple times, adjusting the loop termination code, and then scheduling the unrolled loop. To allow effective scheduling of the loop, we will want to use different registers for each iteration, thus increasing the register count.

Example

Show what our loop looks like unrolled three times (yielding four copies of the loop body), assuming R1 is initially a multiple of 4. Eliminate any obviously redundant computations, and do not reuse any of the registers.

Answer

Here is the result after dropping the unnecessary SUB and BNEZ operations duplicated during unrolling.

```

Loop: LD      F0, 0(R1)
      ADDD   F4, F0, F2
      SD     0(R1), F4 ;drop SUB & BNEZ
      LD     F6, -8(R1)
      ADDD   F8, F6, F2
      SD     -8(R1), F8 ;drop SUB & BNEZ
      LD     F10, -16(R1)
      ADDD   F12, F10, F2
      SD     -16(R1), F12 ;drop SUB & BNEZ
      LD     F14, -24(R1)
      ADDD   F16, F14, F2
      SD     -24(R1), F16
      SUB    R1, R1, #32
      BNEZ   R1, LOOP

```

We have eliminated three branches and three decrements of R1. The addresses on the loads and stores have been compensated for. Without scheduling, every operation is followed by a dependent operation, and thus will cause a stall. This loop will run in 27 clock cycles—each LD takes 2 clock cycles, each ADDD 3, the branch 2, and all other instructions 1—or 6.8 clock cycles for each of the four elements.

Although this unrolled version is currently slower than the scheduled version of the original loop, this will change when we schedule the unrolled loop. Loop unrolling is normally done early in the compilation process, so that redundant computations can be exposed and eliminated by the optimizer.

In real programs we do not normally know the upper bound on the loop. Suppose it is n , and we would like to unroll the loop k times. Instead of a single unrolled loop, we generate a pair of loops. The first executes $(n \bmod k)$ times and has a body that is the original loop. The unrolled version of the loop is surrounded by an outer loop that iterates $(n \operatorname{div} k)$ times. In the above example, unrolling improves the performance of this loop by eliminating overhead instructions, though it increases code size substantially. What will happen to the performance increase when the loop is scheduled on DLX?

Example

Show the unrolled loop in the previous example after it has been scheduled on DLX.

Answer

```

Loop: LD    F0, 0(R1)
      LD    F6, -8(R1)
      LD    F10, -16(R1)
      LD    F14, -24(R1)
      ADDD  F4, F0, F2
      ADDD  F8, F6, F2
      ADDD  F12, F10, F2
      ADDD  F16, F14, F2
      SD    0(R1), F4
      SD    -8(R1), F8
      SD    -16(R1), F12
      SUB   R1, R1, #32 ;branch dependence
      BNEZ  R1, LOOP
      SD    -24(R1), F16 ; 8-32 = -24

```

The execution time of the unrolled loop has dropped to a total of 14 clock cycles, or 3.5 clock cycles per element, compared to 6.8 per element before scheduling.

The gain from scheduling on the unrolled loop is even larger than on the original loop. This is because unrolling the loop exposes more computation that can be scheduled. Scheduling the loop in this fashion necessitates realizing that the loads and stores are independent and can be interchanged.

Loop unrolling is a simple but useful method for increasing the size of straightline code fragments that can be scheduled effectively. This compile-time transformation is similar to what Tomasulo's algorithm does with register renaming and out-of-order execution. As we will see, this is very important in attempts to lower the CPI by issuing instructions at a high rate.

A Superscalar Version of DLX

One method of decreasing the CPI of DLX is to **issue** more than one instruction per clock cycle. This would allow the instruction-execution rate to exceed the clock rate. Machines that issue multiple independent instructions per clock cycle when they are properly scheduled by the compiler have been called *superscalar machines*. In a superscalar machine, the hardware can issue a small number (say 2 to 4) of independent instructions in a single clock. However, if the instructions in the instruction stream are dependent or don't meet certain criteria, only the first instruction in sequence will be issued. A machine where the compiler has complete responsibility for creating a package of instructions that can be simultaneously issued, and the hardware does not dynamically make any decisions about multiple issue, should probably be regarded as a type of VLIW (very long instruction word), which we discuss in the next section.

What would the DLX machine look like as a superscalar? Let's assume two instructions issued per clock cycle. One of the instructions could be a load, store, branch, or integer ALU operation, and the other could be any floating-point operation. As we will see, issue of an integer operation in parallel with a floating-point operation is much simpler and less demanding than arbitrary dual issue.

Issuing two instructions per cycle will require fetching and decoding 64 bits of instructions. To keep the decoding simple, we could require that the instructions be paired and aligned on a 64-bit boundary, with the integer portion appearing first. Figure 6.46 shows how the instructions look as they go into the pipeline in pairs. This table does not address how the floating-point operations extend the EX cycle, but it is no different in the superscalar case than it was for the ordinary DLX pipeline; the concepts of Section 6.6 apply directly. With this pipeline, we have substantially boosted the rate at which we can issue floating-point instructions. To make this worthwhile, however, we need either pipelined floating-point units or multiple independent units. Otherwise, floating-point instructions can only be fetched, and not issued, since all the floating units will be busy.