
D.1	Introduction	D-2
D.2	80x86 Registers and Data Addressing Modes	D-3
D.3	80x86 Integer Operations	D-6
D.4	80x86 Floating-Point Operations	D-10
D.5	80x86 Instruction Encoding	D-12
D.6	Putting It All Together: Measurements of Instruction Set Usage	D-14
D.7	Concluding Remarks	D-20
D.8	Historical Perspective and References	D-21

D

An Alternative to RISC: The Intel 80x86

The x86 isn't all that complex—it just doesn't make a lot of sense.

Mike Johnson

*Leader of 80x86 Design at AMD,
Microprocessor Report (1994)*

D.1 Introduction

MIPS was the vision of a single architect. The pieces of this architecture fit nicely together and the whole architecture can be described succinctly. Such is not the case of the 80x86: It is the product of several independent groups who evolved the architecture over 20 years, adding new features to the original instruction set as you might add clothing to a packed bag. Here are important 80x86 milestones:

- 1978—The Intel 8086 architecture was announced as an assembly language-compatible extension of the then-successful Intel 8080, an 8-bit microprocessor. The 8086 is a 16-bit architecture, with all internal registers 16 bits wide. Whereas the 8080 was a straightforward accumulator machine, the 8086 extended the architecture with additional registers. Because nearly every register has a dedicated use, the 8086 falls somewhere between an accumulator machine and a general-purpose register machine, and can fairly be called an *extended accumulator* machine.
- 1980—The Intel 8087 floating-point coprocessor is announced. This architecture extends the 8086 with about 60 floating-point instructions. Its architects rejected extended accumulators to go with a hybrid of stacks and registers, essentially an *extended stack* architecture: A complete stack instruction set is supplemented by a limited set of register-memory instructions.
- 1982—The 80286 extended the 8086 architecture by increasing the address space to 24 bits, by creating an elaborate memory mapping and protection model, and by adding a few instructions to round out the instruction set and to manipulate the protection model. Because it was important to run 8086 programs without change, the 80286 offered a *real addressing mode* to make the machine look just like an 8086.
- 1985—The 80386 extended the 80286 architecture to 32 bits. In addition to a 32-bit architecture with 32-bit registers and a 32-bit address space, the 80386 added new addressing modes and additional operations. The added instructions make the 80386 nearly a general-purpose register machine. The 80386 also added paging support in addition to segmented addressing (see Chapter 5). Like the 80286, the 80386 has a mode to execute 8086 programs without change.

This history illustrates the impact of the “golden handcuffs” of compatibility on the 80x86, as the existing software base at each step was too important to jeopardize with significant architectural changes. Fortunately, the subsequent 80486 in 1989, Pentium in 1992, and P6 in 1995 were aimed at higher performance, with only four instructions added to the user-visible instruction set: three to help with multiprocessing plus a conditional move instruction.

Since 1997 Intel has added hundreds of instructions to support multimedia by operating on many narrower data types within a single clock (see Chapter 2). These SIMD or vector instructions are primarily used in handcoded libraries or drivers and rarely generated by compilers. The first extension, called MMX,

appeared in 1997. It consists of 57 instructions that pack and unpack multiple bytes, 16-bit words, or 32-bit double words into 64-bit registers and performs shift, logical, and integer arithmetic on the narrow data items in parallel. It supports both saturating and nonsaturating arithmetic. MMX uses the registers comprising the floating-point stack and hence there is no new state for operating systems to save.

In 1999 Intel added another 70 instructions, labeled SSE as part of Pentium III. The primary changes were to add eight separate registers, double their width to 128 bits, and add a single-precision floating-point data type. Hence four 32-bit floating-point operations can be performed in parallel. To improve memory performance, SSE included cache prefetch instructions plus streaming store instructions that bypass the caches and write directly to memory.

In 2001 Intel added yet another 144 instructions, this time labeled SSE2. The new data type is double-precision arithmetic, which allows pairs of 64-bit floating-point operations in parallel. Almost all of these 144 instructions are versions of existing MMX and SSE instructions that operate on 64 bits of data in parallel. Not only does this change enable multimedia operations, it gives the compiler a different target for floating-point operations than the unique stack architecture. Compilers can choose to use the eight SSE registers as floating-point registers as found in the RISC machines. This change has boosted performance on the Pentium 4, the first microprocessor to include SSE2 instructions. At the time of announcement, a 1.5 GHz Pentium 4 was 1.24 times faster than a 1 GHz Pentium III for SPECint2000(base), but it was 1.88 times faster for SPECfp2000(base).

Whatever the artistic failures of the 80x86, keep in mind that there are more instances of this architectural family than of any other server or desktop processor in the world, perhaps 500 million in 2001. Nevertheless, its checkered ancestry has led to an architecture that is difficult to explain and impossible to love.

We start our explanation with the registers and addressing modes, move on to the integer operations, then cover the floating-point operations, and conclude with an examination of instruction encoding.

D.2

80x86 Registers and Data Addressing Modes

The evolution of the instruction set can be seen in the registers of the 80x86 (Figure D.1). Original registers are shown in black type, with the extensions of the 80386 shown in a lighter shade, a coloring scheme followed in subsequent figures. The 80386 basically extended all 16-bit registers (except the segment registers) to 32 bits, prefixing an “E” to their name to indicate the 32-bit version. The arithmetic, logical, and data transfer instructions are two-operand instructions that allow the combinations shown in Figure D.2.

To explain the addressing modes we need to keep in mind whether we are talking about the 16-bit mode used by both the 8086 and 80286 or the 32-bit mode available on the 80386 and its successors. The seven data memory addressing modes supported are

D-4 ■ Appendix D *An Alternative to RISC: The Intel 80x86*

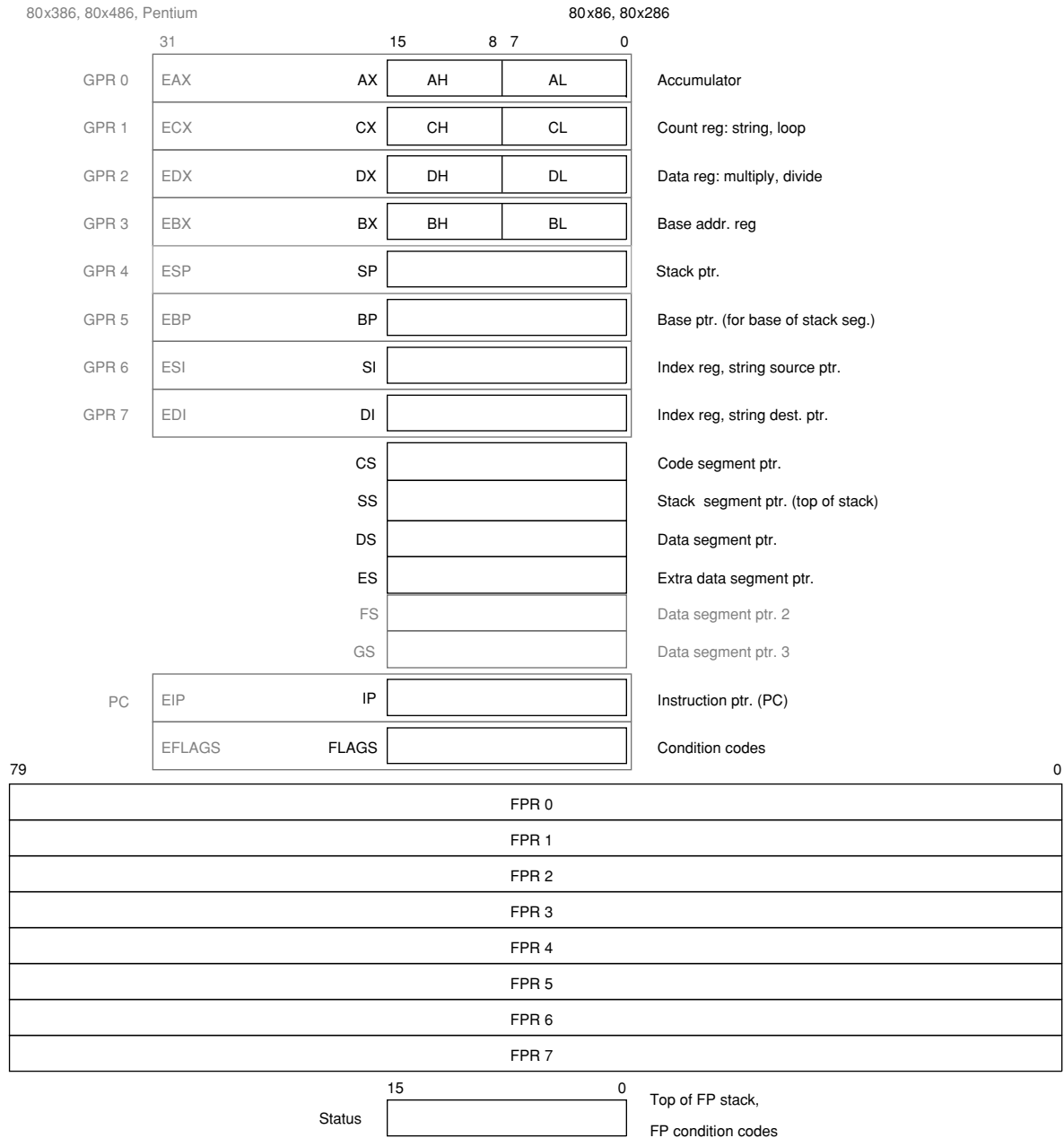


Figure D.1 The 80x86 has evolved over time, and so has its register set. The original set is shown in black, and the extended set in gray. The 8086 divided the first four registers in half so that they could be used either as one 16-bit register or as two 8-bit registers. Starting with the 80386, the top eight registers were extended to 32 bits and could also be used as general-purpose registers. The floating-point registers on the bottom are 80 bits wide, and although they look like regular registers they are not. They implement a stack, with the top of stack pointed to by the status register. One operand must be the top of stack, and the other can be any of the other seven registers below the top of stack.

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

Figure D.2 Instruction types for the arithmetic, logical, and data transfer instructions. The 80x86 allows the combinations shown. The only restriction is the absence of a memory-memory mode. Immediates may be 8, 16, or 32 bits in length; a register is any one of the 14 major registers in Figure D.1 (not IP or FLAGS).

- absolute
- register indirect
- based
- indexed
- based indexed with displacement
- based with scaled indexed
- based with scaled indexed and displacement

Displacements can be 8 or 32 bits in 32-bit mode, and 8 or 16 bits in 16-bit mode. If we count the size of the address as a separate addressing mode, the total is 11 addressing modes.

Although a memory operand can use any addressing mode, there are restrictions on what registers can be used in a mode. Section D.5 on 80x86 instruction encodings gives the full set of restrictions on registers, but the following description of addressing modes gives the basic register options:

- *Absolute*—With 16-bit or 32-bit displacement, depending on the mode.
- *Register indirect*—BX, SI, DI in 16-bit mode and EAX, ECX, EDX, EBX, ESI, and EDI in 32-bit mode.
- *Based mode with 8-bit or 16-bit/32-bit displacement*—BP, BX, SI, DI in 16-bit mode and EAX, ECX, EDX, EBX, ESI, and EDI in 32-bit mode. The displacement is either 8 bits or the size of the address mode: 16 or 32 bits. (Intel gives two different names to this single addressing mode, *based* and *indexed*, but they are essentially identical and we combine them. This book uses indexed addressing to mean something different and is explained next.)
- *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. (The 32-bit mode uses a different addressing mode to get the same effect.)

- *Based indexed with 8- or 16-bit displacement*—The address is the sum of displacement and contents of two registers. The same restrictions on registers apply as in indexed mode.
- *Base plus scaled indexed*—This addressing mode and the next were added in the 80386, and are only available in 32-bit mode. The address calculation is

$$\text{Base register} + 2^{\text{Scale}} \times \text{Index register}$$
 where *Scale* has the value 0, 1, 2, or 3, *Index register* can be any of the eight 32-bit general registers except ESP, and *Base register* can be any of the eight 32-bit general registers.
- *Base plus scaled index with 8- or 32-bit displacement*—The address is the sum of the displacement and the address calculated by the scaled mode immediately above. The same restrictions on registers apply.

The 80x86 uses Little Endian addressing.

Ideally, we would refer discussion of 80x86 logical and physical addresses to Chapter 5, but the segmented address space prevents us from hiding that information. Figure D.3 shows the memory mapping options on the generations of 80x86 machines; Chapter 5 describes the segmented protection scheme in greater detail.

The assembly language programmer clearly must specify which segment register should be used with an address, no matter which address mode is used. To save space in the instructions, segment registers are selected automatically depending on which address register is used. The rules are simple: References to instructions (IP) use the code segment register (CS), references to the stack (BP or SP) use the stack segment register (SS), and the default segment register for the other registers is the data segment register (DS). The next section explains how they can be overridden.

D.3

80x86 Integer Operations

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 80386 adds 32-bit addresses and data, called double words. Almost every operation works on both 8-bit data and one longer data size. That size is determined by the mode and is either 16 or 32 bits.

Clearly some programs want to operate on data of all three sizes, so the 80x86 architects provide a convenient way to specify each version without expanding code size significantly. They decided that most programs would be dominated by either 16- or 32-bit data, and so it made sense to be able to set a default large size. This default size is set by a bit in the code segment register. To override the default size, an 8-bit *prefix* is attached to the instruction to tell the machine to use the other large size for this instruction.

The prefix solution was borrowed from the 8086, which allows multiple prefixes to modify instruction behavior. The three original prefixes override the

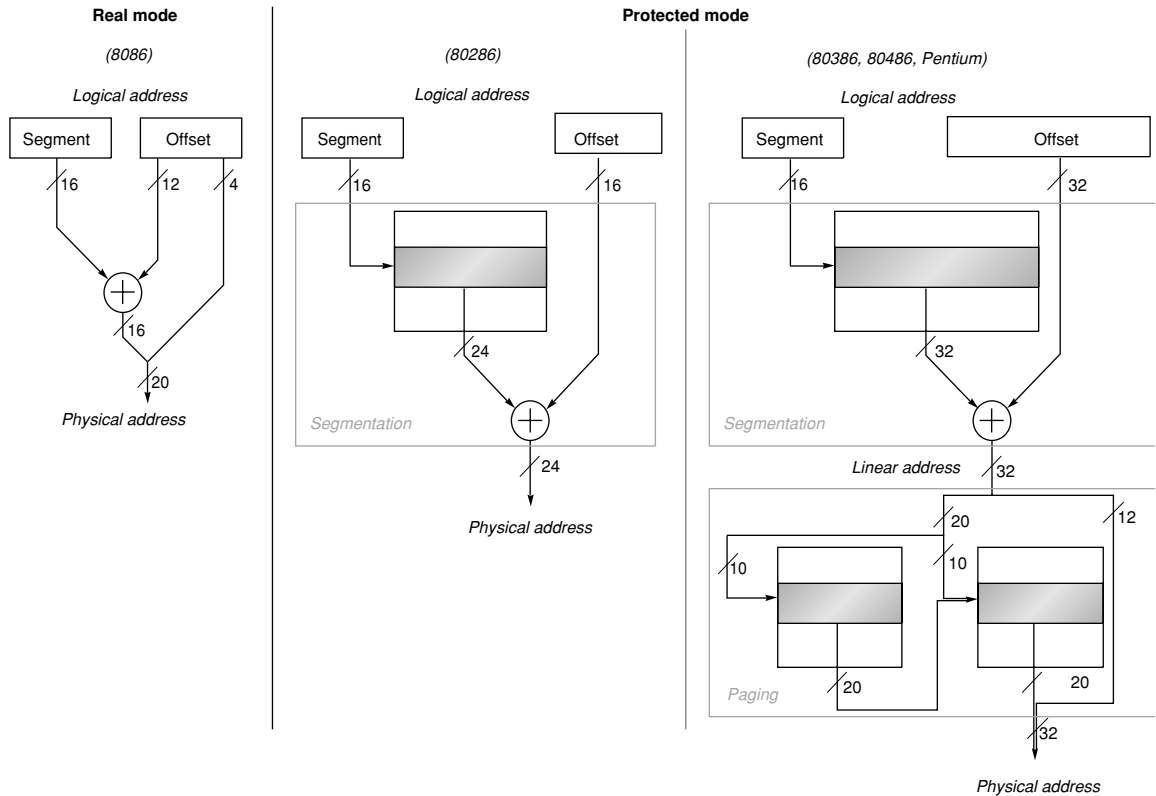


Figure D.3 The original segmented scheme of the 8086 is shown on the left. All 80x86 processors support this style of addressing, called *real mode*. It simply takes the contents of a segment register, shifts it left 4 bits, and adds it to the 16-bit offset, forming a 20-bit physical address. The 80286 (center) used the contents of the segment register to select a segment descriptor, which includes a 24-bit base address among other items. It is added to the 16-bit offset to form the 24-bit physical address. The 80386 and successors (right) expand this base address in the segment descriptor to 32 bits and also add an optional paging layer below segmentation. A 32-bit linear address is first formed from the segment and offset, and then this address is divided into two 10-bit fields and a 12-bit page offset. The first 10-bit field selects the entry in the first-level page table, and then this entry is used in combination with the second 10-bit field to access the second-level page table to select the upper 20 bits of the physical address. Prepending this 20-bit address to the final 12-bit field gives the 32-bit physical address. Paging can be turned off, redefining the 32-bit linear address as the physical address. Note that a “flat” 80x86 address space comes simply by loading the same value in all the segment registers; that is, it doesn’t matter which segment register is selected.

default segment register, lock the bus so as to perform a semaphore (see Chapter 6), or repeat the following instruction until CX counts down to zero. This last prefix was intended to be paired with a byte move instruction to move a variable number of bytes. The 80386 also added a prefix to override the default address size.

The 80x86 integer 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 branches and unconditional jumps, calls, and returns.
4. String instructions, including string move and string compare.

Figure D.4 shows some typical 80x86 instructions and their functions.

The data transfer, arithmetic, and logic instructions are unremarkable, except that the arithmetic and logic instruction operations allow the destination to be either a register or a memory location.

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 in 16-bit mode. One of these is used as the instruction pointer, while the other is loaded into CS and becomes the new code segment. In 32-bit mode the first field is expanded to 32 bits to match the 32-bit program counter (EIP).

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 the 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.

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]←IP+5; SP←SP−2; M[SS:SP]←CS; 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 D.4 Some typical 80x86 instructions and their functions. A list of frequent operations appears in Figure D.5. 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. The hardware description language is described on the back inside cover of this book.

String instructions are part of the 8080 ancestry of the 80x86 and are not commonly executed in most programs.

Figure D.5 lists some of the integer 80x86 instructions. Many of the instructions are available in both byte and word formats.

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- 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
RCR	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 D.5 Some typical operations on the 80x86. 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.

D.4 80x86 Floating-Point Operations

Intel provided a stack architecture with its floating-point instructions: loads push numbers onto the stack, operations find operands in the top two elements of the stacks, and stores can pop elements off the stack, just as the stack example in Figure 2.2 on page 93 suggests.

Intel supplemented this stack architecture with instructions and addressing modes that allow the architecture to have some of the benefits of a register-memory model. In addition to finding operands in the top two elements of the stack, one operand can be in memory or in one of the seven registers below the top of the stack.

This hybrid is still a restricted register-memory model, however, in that loads always move data to the top of the stack while incrementing the top of stack pointer and stores can only move the top of stack to memory. Intel uses the notation *ST* to indicate the top of stack, and *ST(i)* to represent the *i*th register below the top of stack.

One novel feature of this architecture is that the operands are wider in the register stack than they are stored in memory, and all operations are performed at this wide internal precision. Numbers are automatically converted to the internal 80-bit format on a load and converted back to the appropriate size on a store. Memory data can be 32-bit (single-precision) or 64-bit (double-precision) floating-point numbers, called *real* by Intel. The register-memory version of these instructions will then convert the memory operand to this Intel 80-bit format before performing the operation. The data transfer instructions also will automatically convert 16- and 32-bit integers to reals, and vice versa, for integer loads and stores.

The 80x86 floating-point operations can be divided into four major classes:

1. Data movement instructions, including load, load constant, and store.
2. Arithmetic instructions, including add, subtract, multiply, divide, square root, and absolute value.
3. Comparison, including instructions to send the result to the integer CPU so that it can branch.
4. Transcendental instructions, including sine, cosine, log, and exponentiation.

Figure D.6 shows some of the 60 floating-point operations. We use the curly brackets $\{ \}$ to show optional variations of the basic operations: $\{I\}$ means there is an integer version of the instruction, $\{P\}$ means this variation will pop one operand off the stack after the operation, and $\{R\}$ means reverse the sense of the operands in this operation.

Not all combinations are provided. Hence

$$F\{I\}SUB\{R\}\{P\}$$

represents these instructions found in the 80x86:

Data transfer	Arithmetic	Compare	Transcendental
F{I}LD mem/ST(i)	F{I}ADD{P} mem/ST(i)	F{I}COM{P}{P}	FPATAN
F{I}ST{P} mem/ST(i)	F{I}SUB{R}{P} mem/ST(i)	F{I}UCOM{P}{P}	F2XM1
FLDPI	F{I}MUL{P} mem/ST(i)	FSTSW AX/mem	FCOS
FLD1	F{I}DIV{R}{P} mem/ST(i)		FPTAN
FLDZ	FSQRT		FPREM
	FABS		FSIN
	FRNDINT		FYL2X

Figure D.6 The floating-point instructions of the 80x86. The first column shows the data transfer instructions, which move data to memory or to one of the registers below the top of the stack. The last three operations push constants on the stack: pi, 1.0, and 0.0. The second column contains the arithmetic operations described above. Note that the last three operate only on the top of stack. The third column is the compare instructions. Since there are no special floating-point branch instructions, the result of the compare must be transferred to the integer CPU via the FSTSW instruction, either into the AX register or into memory, followed by an SAHF instruction to set the condition codes. The floating-point comparison can then be tested using integer branch instructions. The final column gives the higher-level floating-point operations.

FSUB
FISUB
FSUBR
FISUBR
FSUBP
FSUBRP

There are no pop or reverse pop versions of the integer subtract instructions.

Note that we get even more combinations when including the operand modes for these operations. The floating-point add has these options, ignoring the integer and pop versions of the instruction:

FADD		Both operands in stack, result replaces top of stack.
FADD	ST(i)	One source operand is <i>i</i> th register below the top of stack, and the result replaces the top of stack.
FADD	ST(i),ST	One source operand is the top of stack, and the result replaces <i>i</i> th register below the top of stack.
FADD	mem32	One source operand is a 32-bit location in memory, and the result replaces the top of stack.
FADD	mem64	One source operand is a 64-bit location in memory, and the result replaces the top of stack.

As mentioned above SSE2 presents yet another model of IEEE floating-point registers.

D.5 80x86 Instruction Encoding

Saving the worst for last, the encoding of instructions in the 8086 is complex, with many different instruction formats. Instructions may vary from one byte, when there are no operands, to up to six bytes, when the instruction contains a 16-bit immediate and uses 16-bit displacement addressing. Prefix instructions increase 8086 instruction length beyond the obvious sizes.

The 80386 additions expand the instruction size even further, as Figure D.7 shows. Both the displacement and immediate fields can be 32 bits long, two more prefixes are possible, the opcode can be 16 bits long, and the scaled index mode specifier adds another 8 bits. The maximum possible 80386 instruction is 17 bytes long.

Figure D.8 shows the instruction format for several of the example instructions in Figure D.4. The opcode byte usually contains a bit saying whether the operand is a byte wide or the larger size, 16 bits or 32 bits depending on the mode. 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. Other instructions use a “postbyte” or extra opcode

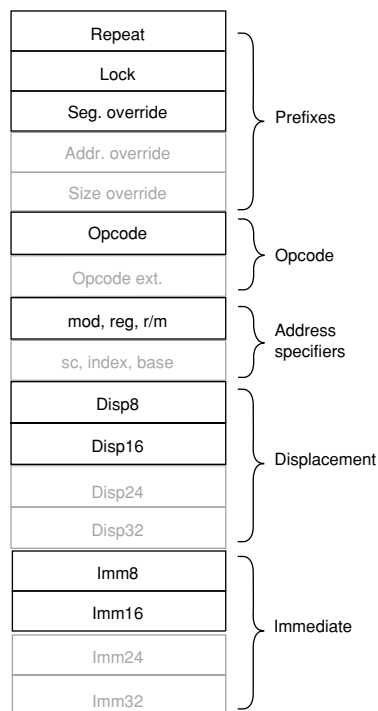


Figure D.7 The instruction format of the 8086 (black type) and the extensions for the 80386 (shaded type). Every field is optional except the opcode.

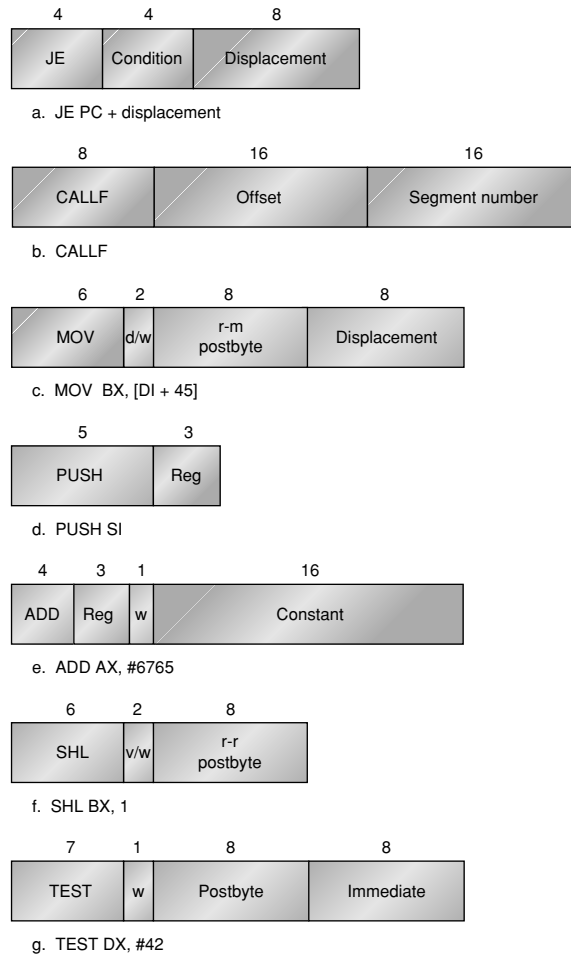


Figure D.8 Typical 8086 instruction formats. The encoding of the postbyte is shown in Figure D.9. Many instructions contain the 1-bit field *w*, which says whether the operation is a byte or a word. Fields of the form *v/w* or *d/w* are a *d*-field or *v*-field followed by the *w*-field. The *d*-field 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.

byte, labeled “mod, reg, r/m” in Figure D.7, which contains the addressing mode information. This postbyte is used for many of the instructions that address memory. The based with scaled index uses a second postbyte, labeled “sc, index, base” in Figure D.7.

The floating-point instructions are encoded in the escape opcode of the 8086 and the postbyte address specifier. The memory operations reserve 2 bits to decide whether the operand is a 32- or 64-bit real or a 16- or 32-bit integer. Those same 2 bits are used in versions that do not access memory to decide whether the stack should be popped after the operation and whether the top of stack or a lower register should get the result.

Alas, you cannot separate the restrictions on registers from the encoding of the addressing modes in the 80x86. Hence Figures D.9 and D.10 show the encoding of the two postbyte address specifiers for both 16- and 32-bit mode.

D.6

Putting It All Together: Measurements of Instruction Set Usage

In this section we present detailed measurements for the 80x86, and then compare the measurements to MIPS for the same programs. To facilitate comparisons among dynamic instruction set measurements, we use a subset of the SPEC92 programs. The 80x86 results were taken in 1994 using the Sun Solaris FORTRAN and C compilers V2.0 and executed in 32-bit mode. These compilers were comparable in quality to the compilers used for MIPS.

reg	w = 0	w = 1		r/m	mod = 0		mod = 1		mod = 2		mod = 3
		16b	32b		16b	32b	16b	32b	16b	32b	
0	AL	AX	EAX	0	addr=BX+SI	=EAX	<i>same</i>	<i>same</i>	<i>same</i>	<i>same</i>	<i>same</i>
1	CL	CX	ECX	1	addr=BX+DI	=ECX	<i>addr as</i>	<i>addr as</i>	<i>addr as</i>	<i>addr as</i>	<i>as</i>
2	DL	DX	EDX	2	addr=BP+SI	=EDX	<i>mod=0</i>	<i>mod=0</i>	<i>mod=0</i>	<i>mod=0</i>	<i>reg</i>
3	BL	BX	EBX	3	addr=BP+SI	=EBX	<i>+ disp8</i>	<i>+ disp8</i>	<i>+ disp16</i>	<i>+ disp32</i>	<i>field</i>
4	AH	SP	ESP	4	addr=SI	=(sib)	SI+disp16	(sib)+disp8	SI+disp8	(sib)+disp32	"
5	CH	BP	EBP	5	addr=DI	=disp32	DI+disp8	EBP+disp8	DI+disp16	EBP+disp32	"
6	DH	SI	ESI	6	addr=disp16	=ESI	BP+disp8	ESI+disp8	BP+disp16	ESI+disp32	"
7	BH	DI	EDI	7	addr=BX	=EDI	BX+disp8	EDI+disp8	BX+disp16	EDI+disp32	"

Figure D.9 The encoding of the first address specifier of the 80x86, "mod, reg, r/m." The first four columns show the encoding of the 3-bit reg field, which depends on the w bit from the opcode and whether the machine is in 16- or 32-bit mode. The remaining columns explain the mod and r/m fields. The meaning of the 3-bit r/m field depends on the value in the 2-bit mod field and the address size. Basically, the registers used in the address calculation are listed in the sixth and seventh columns, under mod = 0, with mod = 1 adding an 8-bit displacement and mod = 2 adding a 16- or 32-bit displacement, depending on the address mode. The exceptions are r/m = 6 when mod = 1 or mod = 2 in 16-bit mode selects BP plus the displacement; r/m = 5 when mod = 1 or mod = 2 in 32-bit mode selects EBP plus displacement; and r/m = 4 in 32-bit mode when mod ≠ 3 (sib) means use the scaled index mode shown in Figure D.10. When mod = 3, the r/m field indicates a register, using the same encoding as the reg field combined with the w bit.

Index		Base
0	EAX	EAX
1	ECX	ECX
2	EDX	EDX
3	EBX	EBX
4	no index	ESP
5	EBP	if mod = 0, disp32 if mod ≠ 0, EBP
6	ESI	ESI
7	EDI	EDI

Figure D.10 Based plus scaled index mode address specifier found in the 80386. This mode is indicated by the (*sib*) notation in Figure D.9. Note that this mode expands the list of registers to be used in other modes: register indirect using ESP comes from Scale = 0, Index = 4, and Base = 4, and base displacement with EBP comes from Scale = 0, Index = 5, and mod = 0. The two-bit scale field is used in this formula of the effective address: Base register + $2^{\text{Scale}} \times$ Index register.

Remember that these measurements depend on the benchmarks chosen and the compiler technology used. Although we feel that the measurements in this section are reasonably indicative of the usage of these 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, and 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.

We start with an evaluation of the features of the 80x86 in isolation, and later compare instruction counts with those of DLX.

Measurements of 80x86 Operand Addressing

We start with addressing modes. Figure D.11 shows the distribution of the operand types in the 80x86. These measurements cover the “second” operand of the operation; for example,

```
mov EAX, [45]
```

counts as a single memory operand. If the types of the first operand were counted, the percentage of register usage would increase by about a factor of 1.5.

The 80x86 memory operands are divided into their respective addressing modes in Figure D.12. Probably the biggest surprise is the popularity of the addressing modes added by the 80386, the last four rows of the figure. They

	Integer average	FP average
Register	45%	22%
Immediate	16%	6%
Memory	39%	72%

Figure D.11 Operand type distribution for the average of five SPECint92 programs (compress, eqntott, espresso, gcc, li) and the average of five SPECfp92 programs (doduc, ear, hydro2d, mdljdp2, su2cor).

Addressing mode	Integer average	FP average
Register indirect	13%	3%
Base + 8-bit disp.	31%	15%
Base + 32-bit disp.	9%	25%
Indexed	0%	0%
Based indexed + 8-bit disp.	0%	0%
Based indexed + 32-bit disp.	0%	1%
Base + scaled indexed	22%	7%
Base + scaled indexed + 8-bit disp.	0%	8%
Base + scaled indexed + 32-bit disp.	4%	4%
32-bit direct	20%	37%

Figure D.12 Operand addressing mode distribution by program. This chart does not include addressing modes used by branches or control instructions.

account for about half of all the memory accesses. Another surprise is the popularity of direct addressing. On most other machines, the equivalent of the direct addressing mode is rare. Perhaps the segmented address space of the 80x86 makes direct addressing more useful, since the address is relative to a base address from the segment register.

These addressing modes largely determine the size of the Intel instructions. Figure D.13 shows the distribution of instruction sizes. The average number of bytes per instruction for integer programs is 2.8, with a standard deviation of 1.5, and 4.1 with a standard deviation of 1.9 for floating-point programs. The difference in length arises partly from the differences in the addressing modes: Integer programs rely more on the shorter register indirect and 8-bit displacement addressing modes, while floating-point programs more frequently use the 80386 addressing modes with the longer 32-bit displacements.

Given that the floating-point instructions have aspects of both stacks and registers, how are they used? Figure D.14 shows that, at least for the compilers used

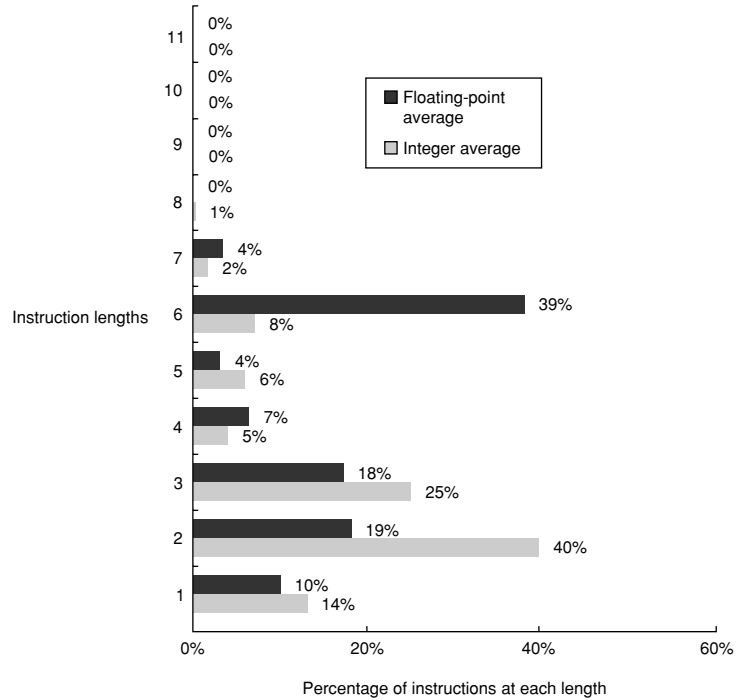


Figure D.13 Averages of the histograms of 80x86 instruction lengths for five SPECint92 programs and for five SPECfp92 programs, all running in 32-bit mode.

Option	doduc	ear	hydro2d	mdljdp2	su2cor	FP average
Stack (2nd operand ST (1))	1.1%	0.0%	0.0%	0.2%	0.6%	0.4%
Register (2nd operand ST(i), i > 1)	17.3%	63.4%	14.2%	7.1%	30.7%	26.5%
Memory	81.6%	36.6%	85.8%	92.7%	68.7%	73.1%

Figure D.14 The percentage of instructions for the floating-point operations (add, sub, mul, div) that use each of the three options for specifying a floating-point operand on the 80x86. The three options are 1) the strict stack model of implicit operands on the stack, 2) register version naming an explicit operand that is not one of the top two elements of the stack, and 3) memory operand.

in this measurement, the stack model of execution is rarely followed. (See Section D.8 for a historical explanation of this observation.)

Finally, Figures D.15 and D.16 show the instruction mixes for 10 SPEC92 programs.

Instruction	compress	eqntott	espresso	gcc (cc1)	li	Int. average
load	20.8%	18.5%	21.9%	24.9%	23.3%	22%
store	13.8%	3.2%	8.3%	16.6%	18.7%	12%
add	10.3%	8.8%	8.15%	7.6%	6.1%	8%
sub	7.0%	10.6%	3.5%	2.9%	3.6%	5%
mul				0.1%		0%
div						0%
compare	8.2%	27.7%	15.3%	13.5%	7.7%	16%
mov reg-reg	7.9%	0.6%	5.0%	4.2%	7.8%	4%
load imm	0.5%	0.2%	0.6%	0.4%		0%
cond. branch	15.5%	28.6%	18.9%	17.4%	15.4%	20%
uncond. branch	1.2%	0.2%	0.9%	2.2%	2.2%	1%
call	0.5%	0.4%	0.7%	1.5%	3.2%	1%
return, jmp indirect	0.5%	0.4%	0.7%	1.5%	3.2%	1%
shift	3.8%		2.5%	1.7%		1%
and	8.4%	1.0%	8.7%	4.5%	8.4%	6%
or	0.6%		2.7%	0.4%	0.4%	1%
other (xor, not, . . .)	0.9%		2.2%	0.1%		1%
load FP						0%
store FP						0%
add FP						0%
sub FP						0%
mul FP						0%
div FP						0%
compare FP						0%
mov reg-reg FP						0%
other (abs, sqrt, . . .)						0%

Figure D.15 80x86 instruction mix for five SPECint92 programs.

Comparative Operation Measurements

Figures D.17 and D.18 show the number of instructions executed for each of the 10 programs on the 80x86 and the ratio of instruction execution compared with that for DLX: Numbers less than 1.0 mean the 80x86 executes fewer instructions than DLX. The instruction count is surprisingly close to DLX for many integer programs, as you would expect a load-store instruction set architecture like DLX to execute more instructions than a register-memory architecture like the 80x86. The floating-point programs always have higher counts for the 80x86, presumably due to the lack of floating-point registers and the use of a stack architecture.

Instruction	doduc	ear	hydro2d	mdljdp2	su2cor	FP average
load	8.9%	6.5%	18.0%	27.6%	27.6%	20%
store	12.4%	3.1%	11.5%	7.8%	7.8%	8%
add	5.4%	6.6%	14.6%	8.8%	8.8%	10%
sub	1.0%	2.4%	3.3%	2.4%	2.4%	3%
mul						0%
div						0%
compare	1.8%	5.1%	0.8%	1.0%	1.0%	2%
mov reg-reg	3.2%	0.1%	1.8%	2.3%	2.3%	2%
load imm	0.4%	1.5%				0%
cond. branch	5.4%	8.2%	5.1%	2.7%	2.7%	5%
uncond branch	0.8%	0.4%	1.3%	0.3%	0.3%	1%
call	0.5%	1.6%		0.1%	0.1%	0%
return, jmp indirect	0.5%	1.6%		0.1%	0.1%	0%
shift	1.1%		4.5%	2.5%	2.5%	2%
and	0.8%	0.8%	0.7%	1.3%	1.3%	1%
or	0.1%			0.1%	0.1%	0%
other (xor, not, . . .)						0%
load FP	14.1%	22.5%	9.1%	12.6%	12.6%	14%
store FP	8.6%	11.4%	4.1%	6.6%	6.6%	7%
add FP	5.8%	6.1%	1.4%	6.6%	6.6%	5%
sub FP	2.2%	2.7%	3.1%	2.9%	2.9%	3%
mul FP	8.9%	8.0%	4.1%	12.0%	12.0%	9%
div FP	2.1%		0.8%	0.2%	0.2%	0%
compare FP	9.4%	6.9%	10.8%	0.5%	0.5%	5%
mov reg-reg FP	2.5%	0.8%	0.3%	0.8%	0.8%	1%
other (abs, sqrt, . . .)	3.9%	3.8%	4.1%	0.8%	0.8%	2%

Figure D.16 80x86 instruction mix for five SPECfp92 programs.

Another question is the total amount of data traffic for the 80x86 versus DLX, since the 80x86 can specify memory operands as part of operations while DLX can only access via loads and stores. Figures D.17 and D.18 also show the data reads, data writes, and data read-modify-writes for these 10 programs. The total accesses ratio to DLX of each memory access type is shown in the bottom rows, with the read-modify-write counting as one read and one write. The 80x86 performs about two to four times as many data accesses as DLX for floating-point programs, and 1.25 times as many for integer programs. Finally, Figure D.19 shows the percentage of instructions in each category for 80x86 and DLX.

	compress	eqntott	espresso	gcc (cc1)	li	Int. avg.
Instructions executed on 80x86 (millions)	2226	1203	2216	3770	5020	
Instructions executed ratio to DLX	0.61	1.74	0.85	0.96	0.98	1.03
Data reads on 80x86 (millions)	589	229	622	1079	1459	
Data writes on 80x86 (millions)	311	39	191	661	981	
Data read-modify-writes on 80x86 (millions)	26	1	129	48	48	
Total data reads on 80x86 (millions)	615	230	751	1127	1507	
Data read ratio to DLX	0.85	1.09	1.38	1.25	0.94	1.10
Total data writes on 80x86 (millions)	338	40	319	709	1029	
Data write ratio to DLX	1.67	9.26	2.39	1.25	1.20	3.15
Total data accesses on 80x86 (millions)	953	269	1070	1836	2536	
Data access ratio to DLX	1.03	1.25	1.58	1.25	1.03	1.23

Figure D.17 Instructions executed and data accesses on 80x86 and ratios compared to DLX for five SPECint92 programs.

	doduc	ear	hydro2d	mdljdp2	su2cor	FP average
Instructions executed on 80x86 (millions)	1223	15,220	13,342	6197	6197	
Instructions executed ratio to DLX	1.19	1.19	2.53	2.09	1.62	1.73
Data reads on 80x86 (millions)	515	6007	5501	3696	3643	
Data writes on 80x86 (millions)	260	2205	2085	892	892	
Data read-modify-writes on 80x86 (millions)	1	0	189	124	124	
Total data reads on 80x86 (millions)	517	6007	5690	3820	3767	
Data read ratio to DLX	2.04	2.36	4.48	4.77	3.91	3.51
Total data writes on 80x86 (millions)	261	2205	2274	1015	1015	
Data write ratio to DLX	3.68	33.25	38.74	16.74	9.35	20.35
Total data accesses on 80x86 (millions)	778	8212	7965	4835	4782	
Data access ratio to DLX	2.40	3.14	5.99	5.73	4.47	4.35

Figure D.18 Instructions executed and data accesses for five SPECfp92 programs on 80x86 and ratio to DLX.

D.7 Concluding Remarks

Beauty is in the eye of the beholder.

Old Adage

As we have seen, “orthogonal” is not a term found in the Intel architectural dictionary. To fully understand which registers and which addressing modes are

Category	Integer average		FP average	
	x86	DLX	x86	DLX
Total data transfer	34%	36%	28%	2%
Total integer arithmetic	34%	31%	16%	12%
Total control	24%	20%	6%	10%
Total logical	8%	13%	3%	2%
Total FP data transfer	0%	0%	22%	33%
Total FP arithmetic	0%	0%	25%	41%

Figure D.19 Percentage of instructions executed by category for 80x86 and DLX for the averages of five SPECint92 and SPECfp92 programs of Figures D.17 and D.18.

available, you need to see the encoding of all addressing modes and sometimes the encoding of the instructions.

Some argue that the inelegance of the 80x86 instruction set is unavoidable, the price that must be paid for rampant success by any architecture. We reject that notion. Obviously no successful architecture can jettison features that were added in previous implementations, and over time some features may be seen as undesirable. The awkwardness of the 80x86 began at its core with the 8086 instruction set and was exacerbated by the architecturally inconsistent expansions of the 8087, 80286, and 80386.

A counterexample is the IBM 360/370 architecture, which is much older than the 80x86. It dominates the mainframe market just as the 80x86 dominates the PC market. Due undoubtedly to a better base and more compatible enhancements, this instruction set makes much more sense than the 80x86 more than 30 years after its first implementation.

For better or worse, Intel had a 16-bit microprocessor years before its competitors' more elegant architectures, and this head start led to the selection of the 8086 as the CPU for the IBM PC. What it lacks in style is made up in quantity, making the 80x86 beautiful from the right perspective.

The saving grace of the 80x86 is that its architectural components are not too difficult to implement, as Intel has demonstrated by rapidly improving performance of integer programs since 1978. High floating-point performance is a larger challenge in this architecture.

D.8

Historical Perspective and References

The complexity of the x86 is not an impassable barrier. ... The biggest weakness in the x86 instruction set is the lack of registers coupled with an extremely painful addressing scheme.

Mike Johnson, Leader of 80x86 Design at AMD
Microprocessor Report (1994)

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 ancestors of the 80x86 were the first microprocessors, produced late in the first half of the 1970s. The Intel 4004 and 8008 were extremely simple 4- and 8-bit accumulator-style machines. Morse et al. [1980] describe the evolution of the 8086 from the 8080 in the late 1970s as 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, has made the 8086 architecture ubiquitous. The success of the IBM PC was due in part because IBM opened the architecture of the PC and enabled the PC-clone industry to flourish. As discussed in the introduction of this appendix, the 80286, 80386, 80486, Pentium, and P6 have extended the architecture and provided a series of performance enhancements.

Although the 68000 was chosen for the popular Macintosh, the Macintosh was never as pervasive as the PC, partly because Apple did not allow clones until recently, 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.

Kahan’s history [1990] of the stack architecture selection for the 8086 is entertaining. The floating-point architecture of the companion 8087 had to be retrofitted into the 8086 opcode space, making it inconvenient to offer two operands per instruction as found in the rest of the 8086. Hence the decision for one operand per instruction using a stack: “The designer’s task was to make a Virtue of this Necessity.” Rather than the classical stack architecture, which has no provision for avoiding common subexpressions from being pushed and popped from memory into the top of the stack found in registers, Intel tried to combine a flat register file with a stack. The reasoning was the restriction of the top of stack as one operand was not so bad since it only required the execution of an FXCH instruction (which swapped registers) to get the same result as a two-operand instruction, and FXCH was much faster than the floating-point operations of the 8087.

Since floating-point expressions are not that complex, Kahan reasoned that eight registers meant that the stack would rarely overflow. Hence he urged that the 8087 use this hybrid scheme with the provision that stack overflow or stack underflow would interrupt the 8086 so that interrupt software could give the illusion to the compiler writer of an unlimited stack for floating-point data. The Intel 8087 was implemented in Israel, and 7500 miles and 10 time zones made communication difficult from California. According to Palmer and Morse [1984]:

Unfortunately, nobody tried to write a software stack manager until after the 8087 was built, and by then it was too late; what was too complicated to perform in hardware turned out to be even worse in software. One thing found lacking is the ability to conveniently determine if an invalid-operation exception is indeed due to a stack overflow. . . . Also lacking is the ability to restart the instruction that caused the stack overflow . . . [p. 93]

The result is that the stack exceptions are too slow to handle in software. As Kahan [1990] says:

Consequently, almost all higher-level languages' compilers emit inefficient code for the 80x87 family, degrading the chip's performance by typically 50% with spurious stores and loads necessary simply to preclude stack over/underflow. . . .

I still regret that the 8087's stack implementation was not quite so neat as my original intention. . . . If the original design had been realized, compilers today would use the 80x87 and its descendents more efficiently, and Intel's competitors could more easily market faster but compatible 80x87 imitations.

The P6 renames the floating-point registers (see Chapter 3), effectively providing up to 40 floating-point registers at any given instant. The main effect of the stack organization is to force design teams to use transistors for dereferencing the stack before doing the renaming.

Hewlett-Packard and Intel have announced a new, common instruction set architecture. It is also upward compatible with the 80x86, and thus the 80x86 instruction set will be available in some form in computers of this century. Instruction set anthropologists will peel off layer by layer from such machines until they uncover artifacts from the first microprocessor. Given such a find, how will they judge 20th-century computer architecture?

References

- Crawford, J., and P. Gelsinger [1988]. *Programming the 80386*, Sybex Books, Alameda, Calif.
- Kahan, J. [1990]. "On the advantage of the 8087's stack," unpublished course notes, Computer Science Division, University of California at Berkeley.
- Morse, S., B. Ravenal, S. Mazor, and W. Pohlman [1980]. "Intel microprocessors—8080 to 8086," *Computer* 13:10 (October).
- Palmer, J., and S. Morse [1984]. *The 8087 Primer*, J. Wiley, New York, p. 93.
- Wakerly, J. [1989]. *Microcomputer Architecture and Programming*, J. Wiley, New York.