

## Chapter 6

# Floating Point

### 6.1 Floating Point Representation

#### 6.1.1 Non-integral binary numbers

When number systems were discussed in the first chapter, only integer values were discussed. Obviously, it must be possible to represent non-integral numbers in other bases as well as decimal. In decimal, digits to the right of the decimal point have associated negative powers of ten:

$$0.123 = 1 \times 10^{-1} + 2 \times 10^{-2} + 3 \times 10^{-3}$$

Not surprisingly, binary numbers work similarly:

$$0.101_2 = 1 \times 2^{-1} + 0 \times 2^{-2} + 1 \times 2^{-3} = 0.625$$

This idea can be combined with the integer methods of Chapter 1 to convert a general number:

$$110.011_2 = 4 + 2 + 0.25 + 0.125 = 6.375$$

Converting from decimal to binary is not very difficult either. In general, divide the decimal number into two parts: integer and fraction. Convert the integer part to binary using the methods from Chapter 1. The fractional part is converted using the method described below.

Consider a binary fraction with the bits labeled  $a, b, c, \dots$ . The number in binary then looks like:

$$0.abcdef \dots$$

Multiply the number by two. The binary representation of the new number will be:

$$a.bcd ef \dots$$

$0.5625 \times 2 = 1.125$	first bit = 1
$0.125 \times 2 = 0.25$	second bit = 0
$0.25 \times 2 = 0.5$	third bit = 0
$0.5 \times 2 = 1.0$	fourth bit = 1

Figure 6.1: Converting 0.5625 to binary

$0.85 \times 2 = 1.7$
$0.7 \times 2 = 1.4$
$0.4 \times 2 = 0.8$
$0.8 \times 2 = 1.6$
$0.6 \times 2 = 1.2$
$0.2 \times 2 = 0.4$
$0.4 \times 2 = 0.8$
$0.8 \times 2 = 1.6$

Figure 6.2: Converting 0.85 to binary

Note that the first bit is now in the one's place. Replace the  $a$  with 0 to get:

$$0.bcdf \dots$$

and multiply by two again to get:

$$b.cdef \dots$$

Now the second bit ( $b$ ) is in the one's position. This procedure can be repeated until as many bits needed are found. Figure 6.1 shows a real example that converts 0.5625 to binary. The method stops when a fractional part of zero is reached.

As another example, consider converting 23.85 to binary. It is easy to convert the integral part ( $23 = 10111_2$ ), but what about the fractional part (0.85)? Figure 6.2 shows the beginning of this calculation. If one looks at

the numbers carefully, an infinite loop is found! This means that 0.85 is a repeating binary (as opposed to a repeating decimal in base 10)<sup>1</sup>. There is a pattern to the numbers in the calculation. Looking at the pattern, one can see that  $0.85 = 0.11\overline{0110}_2$ . Thus,  $23.85 = 10111.11\overline{0110}_2$ .

One important consequence of the above calculation is that 23.85 can not be represented *exactly* in binary using a finite number of bits. (Just as  $\frac{1}{3}$  can not be represented in decimal with a finite number of digits.) As this chapter shows, `float` and `double` variables in C are stored in binary. Thus, values like 23.85 can not be stored exactly into these variables. Only an approximation of 23.85 can be stored.

To simplify the hardware, floating point numbers are stored in a consistent format. This format uses scientific notation (but in binary, using powers of two, not ten). For example, 23.85 or  $10111.11011001100110\dots_2$  would be stored as:

$$1.011111011001100110\dots \times 2^{100}$$

(where the exponent (100) is in binary). A *normalized* floating point number has the form:

$$1.sssssssssssssss \times 2^{eeeeeee}$$

where *1.sssssssssssss* is the *significand* and *eeeeeee* is the *exponent*.

### 6.1.2 IEEE floating point representation

The IEEE (Institute of Electrical and Electronic Engineers) is an international organization that has designed specific binary formats for storing floating point numbers. This format is used on most (but not all!) computers made today. Often it is supported by the hardware of the computer itself. For example, Intel's numeric (or math) coprocessors (which are built into all its CPU's since the Pentium) use it. The IEEE defines two different formats with different precisions: single and double precision. Single precision is used by `float` variables in C and double precision is used by `double` variables.

Intel's math coprocessor also uses a third, higher precision called *extended precision*. In fact, all data in the coprocessor itself is in this precision. When it is stored in memory from the coprocessor it is converted to either single or double precision automatically.<sup>2</sup> Extended precision uses a slightly different general format than the IEEE float and double formats and so will not be discussed here.

---

<sup>1</sup>It should not be so surprising that a number might repeat in one base, but not another. Think about  $\frac{1}{3}$ ; it repeats in decimal, but in ternary (base 3) it would be  $0.1_3$ .

<sup>2</sup> Some compiler's (such as Borland) `long double` type uses this extended precision. However, other compilers use double precision for both `double` and `long double`. (This is allowed by ANSI C.)

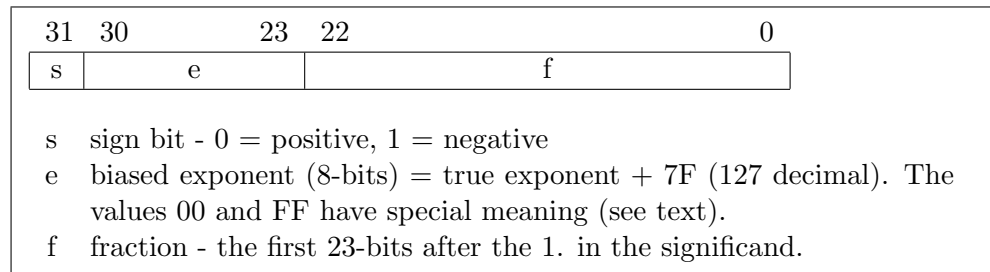


Figure 6.3: IEEE single precision

### IEEE single precision

Single precision floating point uses 32 bits to encode the number. It is usually accurate to 7 significant decimal digits. Floating point numbers are stored in a much more complicated format than integers. Figure 6.3 shows the basic format of a IEEE single precision number. There are several quirks to the format. Floating point numbers do not use the two's complement representation for negative numbers. They use a signed magnitude representation. Bit 31 determines the sign of the number as shown.

The binary exponent is not stored directly. Instead, the sum of the exponent and 7F is stored from bit 23 to 30. This *biased exponent* is always non-negative.

The fraction part assumes a normalized significand (in the form 1.*ssssssss*). Since the first bit is always an one, the leading one is *not stored*! This allows the storage of an additional bit at the end and so increases the precision slightly. This idea is known as the *hidden one representation*.

*One should always keep in mind that the bytes 41 BE CC CD can be interpreted different ways depending on what a program does with them! As a single precision floating point number, they represent 23.850000381, but as a double word integer, they represent 1,103,023,309! The CPU does not know which is the correct interpretation!*

How would 23.85 be stored? First, it is positive so the sign bit is 0. Next the true exponent is 4, so the biased exponent is  $7F + 4 = 83_{16}$ . Finally, the fraction is 0111101100110011001100 (remember the leading one is hidden). Putting this all together (to help clarify the different sections of the floating point format, the sign bit and the fraction have been underlined and the bits have been grouped into 4-bit nibbles):

$$\underline{0} \ 100 \ 0001 \ 1 \underline{011 \ 1110 \ 1100 \ 1100 \ 1100 \ 1100}_2 = 41BECCCC_{16}$$

This is not exactly 23.85 (since it is a repeating binary). If one converts the above back to decimal, one finds that it is approximately 23.849998474. This number is very close to 23.85, but it is not exact. Actually, in C, 23.85 would not be represented exactly as above. Since the left-most bit that was truncated from the exact representation is 1, the last bit is rounded up to 1. So 23.85 would be represented as 41 BE CC CD in hex using single precision. Converting this to decimal results in 23.850000381 which is a slightly better approximation of 23.85.

$e = 0$	and	$f = 0$	denotes the number zero (which can not be normalized) Note that there is a +0 and -0.
$e = 0$	and	$f \neq 0$	denotes a <i>denormalized number</i> . These are discussed in the next section.
$e = FF$	and	$f = 0$	denotes infinity ( $\infty$ ). There are both positive and negative infinities.
$e = FF$	and	$f \neq 0$	denotes an undefined result, known as <i>NaN</i> (Not a Number).

Table 6.1: Special values of  $f$  and  $e$



Figure 6.4: IEEE double precision

How would -23.85 be represented? Just change the sign bit: C1 BE CC CD. Do *not* take the two's complement!

Certain combinations of  $e$  and  $f$  have special meanings for IEEE floats. Table 6.1 describes these special values. An infinity is produced by an overflow or by division by zero. An undefined result is produced by an invalid operation such as trying to find the square root of a negative number, adding two infinities, *etc.*

Normalized single precision numbers can range in magnitude from  $1.0 \times 2^{-126}$  ( $\approx 1.1755 \times 10^{-35}$ ) to  $1.11111 \dots \times 2^{127}$  ( $\approx 3.4028 \times 10^{35}$ ).

**Denormalized numbers**

Denormalized numbers can be used to represent numbers with magnitudes too small to normalize (*i.e.* below  $1.0 \times 2^{-126}$ ). For example, consider the number  $1.001_2 \times 2^{-129}$  ( $\approx 1.6530 \times 10^{-39}$ ). In the given normalized form, the exponent is too small. However, it can be represented in the unnormalized form:  $0.01001_2 \times 2^{-127}$ . To store this number, the biased exponent is set to 0 (see Table 6.1) and the fraction is the complete significand of the number written as a product with  $2^{-127}$  (*i.e.* all bits are stored including the one to the left of the decimal point). The representation of  $1.001 \times 2^{-129}$  is then:

0 000 0000 0 001 0010 0000 0000 0000 0000

### IEEE double precision

IEEE double precision uses 64 bits to represent numbers and is usually accurate to about 15 significant decimal digits. As Figure 6.4 shows, the basic format is very similar to single precision. More bits are used for the biased exponent (11) and the fraction (52) than for single precision.

The larger range for the biased exponent has two consequences. The first is that it is calculated as the sum of the true exponent and 3FF (1023) (not 7F as for single precision). Secondly, a large range of true exponents (and thus a larger range of magnitudes) is allowed. Double precision magnitudes can range from approximately  $10^{-308}$  to  $10^{308}$ .

It is the larger field of the fraction that is responsible for the increase in the number of significant digits for double values.

As an example, consider 23.85 again. The biased exponent will be  $4 + 3FF = 403$  in hex. Thus, the double representation would be:

0 100 0000 0011 0111 1101 1001 1001 1001 1001 1001 1001 1001 1001 1010

or 40 37 D9 99 99 99 9A in hex. If one converts this back to decimal, one finds 23.8500000000000014 (there are 12 zeros!) which is a much better approximation of 23.85.

The double precision has the same special values as single precision<sup>3</sup>. Denormalized numbers are also very similar. The only main difference is that double denormalized numbers use  $2^{-1023}$  instead of  $2^{-127}$ .

## 6.2 Floating Point Arithmetic

Floating point arithmetic on a computer is different than in continuous mathematics. In mathematics, all numbers can be considered exact. As shown in the previous section, on a computer many numbers can not be represented exactly with a finite number of bits. All calculations are performed with limited precision. In the examples of this section, numbers with an 8-bit significand will be used for simplicity.

### 6.2.1 Addition

To add two floating point numbers, the exponents must be equal. If they are not already equal, then they must be made equal by shifting the significand of the number with the smaller exponent. For example, consider  $10.375 + 6.34375 = 16.71875$  or in binary:

$$\begin{array}{r} 1.0100110 \times 2^3 \\ + 1.1001011 \times 2^2 \\ \hline \end{array}$$

<sup>3</sup>The only difference is that for the infinity and undefined values, the biased exponent is 7FF not FF.

These two numbers do not have the same exponent so shift the significand to make the exponents the same and then add:

$$\begin{array}{r} 1.0100110 \times 2^3 \\ + \quad 0.1100110 \times 2^3 \\ \hline 10.0001100 \times 2^3 \end{array}$$

Note that the shifting of  $1.1001011 \times 2^2$  drops off the trailing one and after rounding results in  $0.1100110 \times 2^3$ . The result of the addition,  $10.0001100 \times 2^3$  (or  $1.00001100 \times 2^4$ ) is equal to  $10000.110_2$  or 16.75. This is *not* equal to the exact answer (16.71875)! It is only an approximation due to the round off errors of the addition process.

It is important to realize that floating point arithmetic on a computer (or calculator) is always an approximation. The laws of mathematics do not always work with floating point numbers on a computer. Mathematics assumes infinite precision which no computer can match. For example, mathematics teaches that  $(a + b) - b = a$ ; however, this may not hold true exactly on a computer!

### 6.2.2 Subtraction

Subtraction works very similarly and has the same problems as addition. As an example, consider  $16.75 - 15.9375 = 0.8125$ :

$$\begin{array}{r} 1.0000110 \times 2^4 \\ - \quad 1.1111111 \times 2^3 \\ \hline \end{array}$$

Shifting  $1.1111111 \times 2^3$  gives (rounding up)  $1.0000000 \times 2^4$

$$\begin{array}{r} 1.0000110 \times 2^4 \\ - \quad 1.0000000 \times 2^4 \\ \hline 0.0000110 \times 2^4 \end{array}$$

$0.0000110 \times 2^4 = 0.11_2 = 0.75$  which is not exactly correct.

### 6.2.3 Multiplication and division

For multiplication, the significands are multiplied and the exponents are added. Consider  $10.375 \times 2.5 = 25.9375$ :

$$\begin{array}{r} 1.0100110 \times 2^3 \\ \times \quad 1.0100000 \times 2^1 \\ \hline 10100110 \\ + \quad 10100110 \\ \hline 1.10011111000000 \times 2^4 \end{array}$$

Of course, the real result would be rounded to 8-bits to give:

$$1.1010000 \times 2^4 = 11010.000_2 = 26$$

Division is more complicated, but has similar problems with round off errors.

### 6.2.4 Ramifications for programming

The main point of this section is that floating point calculations are not exact. The programmer needs to be aware of this. A common mistake that programmers make with floating point numbers is to compare them assuming that a calculation is exact. For example, consider a function named  $f(x)$  that makes a complex calculation and a program is trying to find the function's roots<sup>4</sup>. One might be tempted to use the following statement to check to see if  $x$  is a root:

```
if ( f(x) == 0.0 )
```

But, what if  $f(x)$  returns  $1 \times 10^{-30}$ ? This very likely means that  $x$  is a *very* good approximation of a true root; however, the equality will be false. There may not be any IEEE floating point value of  $x$  that returns exactly zero, due to round off errors in  $f(x)$ .

A much better method would be to use:

```
if ( fabs(f(x)) < EPS )
```

where `EPS` is a macro defined to be a very small positive value (like  $1 \times 10^{-10}$ ). This is true whenever  $f(x)$  is very close to zero. In general, to compare a floating point value (say  $x$ ) to another ( $y$ ) use:

```
if ( fabs(x - y)/fabs(y) < EPS )
```

## 6.3 The Numeric Coprocessor

### 6.3.1 Hardware

The earliest Intel processors had no hardware support for floating point operations. This does not mean that they could not perform float operations. It just means that they had to be performed by procedures composed of many non-floating point instructions. For these early systems, Intel did provide an additional chip called a *math coprocessor*. A math coprocessor has machine instructions that perform many floating point operations much faster than using a software procedure (on early processors, at least 10 times

---

<sup>4</sup>A root of a function is a value  $x$  such that  $f(x) = 0$



faster!). The coprocessor for the 8086/8088 was called the 8087. For the 80286, there was a 80287 and for the 80386, a 80387. The 80486DX processor integrated the math coprocessor into the 80486 itself.<sup>5</sup> Since the Pentium, all generations of 80x86 processors have a builtin math coprocessor; however, it is still programmed as if it was a separate unit. Even earlier systems without a coprocessor can install software that emulates a math coprocessor. These emulator packages are automatically activated when a program executes a coprocessor instruction and run a software procedure that produces the same result as the coprocessor would have (though much slower, of course).

The numeric coprocessor has eight floating point registers. Each register holds 80 bits of data. Floating point numbers are *always* stored as 80-bit extended precision numbers in these registers. The registers are named ST0, ST1, ST2, . . . ST7. The floating point registers are used differently than the integer registers of the main CPU. The floating point registers are organized as a *stack*. Recall that a stack is a *Last-In First-Out* (LIFO) list. ST0 always refers to the value at the top of the stack. All new numbers are added to the top of the stack. Existing numbers are pushed down on the stack to make room for the new number.

There is also a status register in the numeric coprocessor. It has several flags. Only the 4 flags used for comparisons will be covered: C<sub>0</sub>, C<sub>1</sub>, C<sub>2</sub> and C<sub>3</sub>. The use of these is discussed later.

### 6.3.2 Instructions

To make it easy to distinguish the normal CPU instructions from coprocessor ones, all the coprocessor mnemonics start with an F.

#### Loading and storing

There are several instructions that load data onto the top of the coprocessor register stack:

FLD <i>source</i>	loads a floating point number from memory onto the top of the stack. The <i>source</i> may be a single, double or extended precision number or a coprocessor register.
FILD <i>source</i>	reads an <i>integer</i> from memory, converts it to floating point and stores the result on top of the stack. The <i>source</i> may be either a word, double word or quad word.
FLD1	stores a one on the top of the stack.
FLDZ	stores a zero on the top of the stack.

There are also several instructions that store data from the stack into memory. Some of these instructions also *pop* (*i.e.* remove) the number from

---

<sup>5</sup>However, the 80486SX did *not* have an integrated coprocessor. There was a separate 80487SX chip for these machines.

the stack as it stores it.

- FST *dest*** stores the top of the stack (ST0) into memory. The *destination* may either be a single or double precision number or a coprocessor register.
- FSTP *dest*** stores the top of the stack into memory just as **FST**; however, after the number is stored, its value is popped from the stack. The *destination* may either a single, double or extended precision number or a coprocessor register.
- FIST *dest*** stores the value of the top of the stack converted to an integer into memory. The *destination* may either a word or a double word. The stack itself is unchanged. How the floating point number is converted to an integer depends on some bits in the coprocessor's *control word*. This is a special (non-floating point) word register that controls how the coprocessor works. By default, the control word is initialized so that it rounds to the nearest integer when it converts to integer. However, the **FSTCW** (Store Control Word) and **FLDCW** (Load Control Word) instructions can be used to change this behavior.
- FISTP *dest*** Same as **FIST** except for two things. The top of the stack is popped and the *destination* may also be a quad word.

There are two other instructions that can move or remove data on the stack itself.

- FXCH ST*n*** exchanges the values in ST0 and ST*n* on the stack (where *n* is register number from 1 to 7).
- FFREE ST*n*** frees up a register on the stack by marking the register as unused or empty.

### Addition and subtraction

Each of the addition instructions compute the sum of ST0 and another operand. The result is always stored in a coprocessor register.

- FADD *src*** ST0 += *src*. The *src* may be any coprocessor register or a single or double precision number in memory.
- FADD *dest*, ST0** *dest* += ST0. The *dest* may be any coprocessor register.
- FADDP *dest* or  
FADDP *dest*, ST0** *dest* += ST0 then pop stack. The *dest* may be any coprocessor register.
- FIADD *src*** ST0 += (float) *src*. Adds an integer to ST0. The *src* must be a word or double word in memory.

There are twice as many subtraction instructions than addition because the order of the operands is important for subtraction (*i.e.*  $a + b = b + a$ , but  $a - b \neq b - a$ !). For each instruction, there is an alternate one that subtracts in the reverse order. These reverse instructions all end in either

```

1 segment .bss
2 array      resq SIZE
3 sum        resq 1
4
5 segment .text
6     mov     ecx, SIZE
7     mov     esi, array
8     fldz                    ; ST0 = 0
9 lp:
10    fadd     qword [esi]     ; ST0 += *(esi)
11    add      esi, 8          ; move to next double
12    loop     lp
13    fstp     qword sum       ; store result into sum

```

Figure 6.5: Array sum example

R or RP. Figure 6.5 shows a short code snippet that adds up the elements of an array of doubles. On lines 10 and 13, one must specify the size of the memory operand. Otherwise the assembler would not know whether the memory operand was a float (dword) or a double (qword).

<code>FSUB <i>src</i></code>	<code>ST0 -= <i>src</i></code> . The <i>src</i> may be any coprocessor register or a single or double precision number in memory.
<code>FSUBR <i>src</i></code>	<code>ST0 = <i>src</i> - ST0</code> . The <i>src</i> may be any coprocessor register or a single or double precision number in memory.
<code>FSUB <i>dest</i>, ST0</code>	<code><i>dest</i> -= ST0</code> . The <i>dest</i> may be any coprocessor register.
<code>FSUBR <i>dest</i>, ST0</code>	<code><i>dest</i> = ST0 - <i>dest</i></code> . The <i>dest</i> may be any coprocessor register.
<code>FSUBP <i>dest</i> or FSUBP <i>dest</i>, ST0</code>	<code><i>dest</i> -= ST0</code> then pop stack. The <i>dest</i> may be any coprocessor register.
<code>FSUBRP <i>dest</i> or FSUBRP <i>dest</i>, ST0</code>	<code><i>dest</i> = ST0 - <i>dest</i></code> then pop stack. The <i>dest</i> may be any coprocessor register.
<code>FISUB <i>src</i></code>	<code>ST0 -= (float) <i>src</i></code> . Subtracts an integer from ST0. The <i>src</i> must be a word or double word in memory.
<code>FISUBR <i>src</i></code>	<code>ST0 = (float) <i>src</i> - ST0</code> . Subtracts ST0 from an integer. The <i>src</i> must be a word or double word in memory.

## Multiplication and division

The multiplication instructions are completely analogous to the addition instructions.

FMUL <i>src</i>	ST0 *= <i>src</i> . The <i>src</i> may be any coprocessor register or a single or double precision number in memory.
FMUL <i>dest</i> , ST0	<i>dest</i> *= ST0. The <i>dest</i> may be any coprocessor register.
FMULP <i>dest</i> or FMULP <i>dest</i> , ST0	<i>dest</i> *= ST0 then pop stack. The <i>dest</i> may be any coprocessor register.
FIMUL <i>src</i>	ST0 *= (float) <i>src</i> . Multiplies an integer to ST0. The <i>src</i> must be a word or double word in memory.

Not surprisingly, the division instructions are analogous to the subtraction instructions. Division by zero results in an infinity.

FDIV <i>src</i>	ST0 /= <i>src</i> . The <i>src</i> may be any coprocessor register or a single or double precision number in memory.
FDIVR <i>src</i>	ST0 = <i>src</i> / ST0. The <i>src</i> may be any coprocessor register or a single or double precision number in memory.
FDIV <i>dest</i> , ST0	<i>dest</i> /= ST0. The <i>dest</i> may be any coprocessor register.
FDIVR <i>dest</i> , ST0	<i>dest</i> = ST0 / <i>dest</i> . The <i>dest</i> may be any coprocessor register.
FDIVP <i>dest</i> or FDIVP <i>dest</i> , ST0	<i>dest</i> /= ST0 then pop stack. The <i>dest</i> may be any coprocessor register.
FDIVRP <i>dest</i> or FDIVRP <i>dest</i> , ST0	<i>dest</i> = ST0 / <i>dest</i> then pop stack. The <i>dest</i> may be any coprocessor register.
FIDIV <i>src</i>	ST0 /= (float) <i>src</i> . Divides ST0 by an integer. The <i>src</i> must be a word or double word in memory.
FIDIVR <i>src</i>	ST0 = (float) <i>src</i> / ST0. Divides an integer by ST0. The <i>src</i> must be a word or double word in memory.

## Comparisons

The coprocessor also performs comparisons of floating point numbers. The FCOM family of instructions does this operation.

```

1  ;    if ( x > y )
2  ;
3      fld    qword [x]          ; ST0 = x
4      fcomp  qword [y]          ; compare ST0 and y
5      fstsw  ax                  ; move C bits into FLAGS
6      sahf
7      jna    else_part          ; if x not above y, goto else_part
8  then_part:
9      ; code for then part
10     jmp    end_if
11  else_part:
12     ; code for else part
13  end_if:

```

Figure 6.6: Comparison example

<b>FCOM <i>src</i></b>	compares ST0 and <i>src</i> . The <i>src</i> can be a coprocessor register or a float or double in memory.
<b>FCOMP <i>src</i></b>	compares ST0 and <i>src</i> , then pops stack. The <i>src</i> can be a coprocessor register or a float or double in memory.
<b>FCOMPP</b>	compares ST0 and ST1, then pops stack twice.
<b>FICOM <i>src</i></b>	compares ST0 and (float) <i>src</i> . The <i>src</i> can be a word or dword integer in memory.
<b>FICOMP <i>src</i></b>	compares ST0 and (float) <i>src</i> , then pops stack. The <i>src</i> can be a word or dword integer in memory.
<b>FTST</b>	compares ST0 and 0.

These instructions change the C<sub>0</sub>, C<sub>1</sub>, C<sub>2</sub> and C<sub>3</sub> bits of the coprocessor status register. Unfortunately, it is not possible for the CPU to access these bits directly. The conditional branch instructions use the FLAGS register, not the coprocessor status register. However, it is relatively simple to transfer the bits of the status word into the corresponding bits of the FLAGS register using some new instructions:

<b>FSTSW <i>dest</i></b>	Stores the coprocessor status word into either a word in memory or the AX register.
<b>SAHF</b>	Stores the AH register into the FLAGS register.
<b>LAHF</b>	Loads the AH register with the bits of the FLAGS register.

Figure 6.6 shows a short example code snippet. Lines 5 and 6 transfer the C<sub>0</sub>, C<sub>1</sub>, C<sub>2</sub> and C<sub>3</sub> bits of the coprocessor status word into the FLAGS register. The bits are transferred so that they are analogous to the result of a comparison of two *unsigned* integers. This is why line 7 uses a JNA instruction.

The Pentium Pro (and later processors (Pentium II and III)) support two new comparison operators that directly modify the CPU's FLAGS register.

**FCOMI *src*** compares ST0 and *src*. The *src* must be a coprocessor register.

**FCOMIP *src*** compares ST0 and *src*, then pops stack. The *src* must be a coprocessor register.

Figure 6.7 shows an example subroutine that finds the maximum of two doubles using the FCOMIP instruction. Do not confuse these instructions with the integer comparison functions (FICOM and FICOMP).

### Miscellaneous instructions

This section covers some other miscellaneous instructions that the coprocessor provides.

**FCHS** ST0 = - ST0 Changes the sign of ST0

**FABS** ST0 = |ST0| Takes the absolute value of ST0

**FSQRT** ST0 =  $\sqrt{\text{ST0}}$  Takes the square root of ST0

**FSCALE** ST0 = ST0  $\times 2^{\text{ST1}}$  multiplies ST0 by a power of 2 quickly. ST1 is not removed from the coprocessor stack. Figure 6.8 shows an example of how to use this instruction.

### 6.3.3 Examples

#### 6.3.4 Quadratic formula

The first example shows how the quadratic formula can be encoded in assembly. Recall that the quadratic formula computes the solutions to the quadratic equation:

$$ax^2 + bx + c = 0$$

The formula itself gives two solutions for  $x$ :  $x_1$  and  $x_2$ .

$$x_1, x_2 = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

The expression inside the square root ( $b^2 - 4ac$ ) is called the *discriminant*. Its value is useful in determining which of the following three possibilities are true for the solutions.

1. There is only one real degenerate solution.  $b^2 - 4ac = 0$
2. There are two real solutions.  $b^2 - 4ac > 0$
3. There are two complex solutions.  $b^2 - 4ac < 0$

Here is a small C program that uses the assembly subroutine:

---

**quadt.c**

---

```

1  #include <stdio.h>
2
3  int quadratic( double, double, double, double *, double *);
4
5  int main()
6  {
7      double a,b,c, root1, root2;
8
9      printf("Enter a, b, c: ");
10     scanf("%lf %lf %lf", &a, &b, &c);
11     if (quadratic( a, b, c, &root1, &root2) )
12         printf("roots: %.10g %.10g\n", root1, root2);
13     else
14         printf("No real roots\n");
15     return 0;
16 }
```

---

**quadt.c**

---

Here is the assembly routine:

---

**quad.asm**

---

```

1  ; function quadratic
2  ; finds solutions to the quadratic equation:
3  ;      a*x^2 + b*x + c = 0
4  ; C prototype:
5  ;   int quadratic( double a, double b, double c,
6  ;                  double * root1, double *root2 )
7  ; Parameters:
8  ;   a, b, c - coefficients of powers of quadratic equation (see above)
9  ;   root1   - pointer to double to store first root in
10 ;   root2   - pointer to double to store second root in
11 ; Return value:
12 ;   returns 1 if real roots found, else 0
13
14 %define a          qword [ebp+8]
15 %define b          qword [ebp+16]
16 %define c          qword [ebp+24]
17 %define root1      dword [ebp+32]
18 %define root2      dword [ebp+36]
19 %define disc       qword [ebp-8]
```

```

20 %define one_over_2a      qword [ebp-16]
21
22 segment .data
23 MinusFour      dw      -4
24
25 segment .text
26     global      _quadratic
27 _quadratic:
28     push      ebp
29     mov       ebp, esp
30     sub       esp, 16      ; allocate 2 doubles (disc & one_over_2a)
31     push      ebx          ; must save original ebx
32
33     fld      word [MinusFour]; stack -4
34     fld      a              ; stack: a, -4
35     fld      c              ; stack: c, a, -4
36     fmulp    st1            ; stack: a*c, -4
37     fmulp    st1            ; stack: -4*a*c
38     fld      b
39     fld      b              ; stack: b, b, -4*a*c
40     fmulp    st1            ; stack: b*b, -4*a*c
41     faddp    st1            ; stack: b*b - 4*a*c
42     ftst
43     fstsw    ax             ; test with 0
44     sahf
45     jb      no_real_solutions ; if disc < 0, no real solutions
46     fsqrt
47     fstp     disc           ; stack: sqrt(b*b - 4*a*c)
48     fld1
49     fld      a              ; stack: a, 1.0
50     fscale
51     fdivp    st1            ; stack: 1/(2*a)
52     fst      one_over_2a    ; stack: 1/(2*a)
53     fld      b              ; stack: b, 1/(2*a)
54     fld      disc           ; stack: disc, b, 1/(2*a)
55     fsubrp   st1            ; stack: disc - b, 1/(2*a)
56     fmulp    st1            ; stack: (-b + disc)/(2*a)
57     mov      ebx, root1
58     fstp     qword [ebx]    ; store in *root1
59     fld      b              ; stack: b
60     fld      disc           ; stack: disc, b
61     fchs

```



```

62      fsubrp  st1          ; stack: -disc - b
63      fmul   one_over_2a  ; stack: (-b - disc)/(2*a)
64      mov    ebx, root2
65      fstp   qword [ebx]   ; store in *root2
66      mov    eax, 1        ; return value is 1
67      jmp    short quit
68
69 no_real_solutions:
70      mov    eax, 0        ; return value is 0
71
72 quit:
73      pop    ebx
74      mov    esp, ebp
75      pop    ebp
76      ret

```

---

quad.asm

### 6.3.5 Reading array from file

In this example, an assembly routine reads doubles from a file. Here is a short C test program:

---

readt.c

---

```

1  /*
2   * This program tests the 32-bit read_doubles() assembly procedure.
3   * It reads the doubles from stdin. (Use redirection to read from file.)
4   */
5  #include <stdio.h>
6  extern int read_doubles( FILE *, double *, int );
7  #define MAX 100
8
9  int main()
10 {
11     int i,n;
12     double a[MAX];
13
14     n = read_doubles(stdin, a, MAX);
15
16     for( i=0; i < n; i++ )
17         printf("%3d %g\n", i, a[i]);
18     return 0;
19 }

```

---

 readt.c
 

---

Here is the assembly routine

---

 read.asm
 

---

```

1 segment .data
2 format db      "%lf", 0          ; format for fscanf()
3
4 segment .text
5     global  _read_doubles
6     extern  _fscanf
7
8 %define SIZEOF_DOUBLE  8
9 %define FP             dword [ebp + 8]
10 %define ARRAYP         dword [ebp + 12]
11 %define ARRAY_SIZE     dword [ebp + 16]
12 %define TEMP_DOUBLE    [ebp - 8]
13
14 ;
15 ; function _read_doubles
16 ; C prototype:
17 ;   int read_doubles( FILE * fp, double * arrayp, int array_size );
18 ; This function reads doubles from a text file into an array, until
19 ; EOF or array is full.
20 ; Parameters:
21 ;   fp          - FILE pointer to read from (must be open for input)
22 ;   arrayp      - pointer to double array to read into
23 ;   array_size  - number of elements in array
24 ; Return value:
25 ;   number of doubles stored into array (in EAX)
26
27 _read_doubles:
28     push    ebp
29     mov     ebp, esp
30     sub     esp, SIZEOF_DOUBLE      ; define one double on stack
31
32     push    esi                    ; save esi
33     mov     esi, ARRAYP            ; esi = ARRAYP
34     xor     edx, edx                ; edx = array index (initially 0)
35
36 while_loop:
37     cmp     edx, ARRAY_SIZE        ; is edx < ARRAY_SIZE?
38     jnl     short quit              ; if not, quit loop

```

```

39 ;
40 ; call fscanf() to read a double into TEMP_DOUBLE
41 ; fscanf() might change edx so save it
42 ;
43     push    edx                ; save edx
44     lea     eax, TEMP_DOUBLE
45     push    eax                ; push &TEMP_DOUBLE
46     push    dword format      ; push &format
47     push    FP                ; push file pointer
48     call    _fscanf
49     add     esp, 12
50     pop     edx                ; restore edx
51     cmp     eax, 1             ; did fscanf return 1?
52     jne     short quit        ; if not, quit loop
53
54 ;
55 ; copy TEMP_DOUBLE into ARRAYP[edx]
56 ; (The 8-bytes of the double are copied by two 4-byte copies)
57 ;
58     mov     eax, [ebp - 8]
59     mov     [esi + 8*edx], eax ; first copy lowest 4 bytes
60     mov     eax, [ebp - 4]
61     mov     [esi + 8*edx + 4], eax ; next copy highest 4 bytes
62
63     inc     edx
64     jmp     while_loop
65
66 quit:
67     pop     esi                ; restore esi
68
69     mov     eax, edx           ; store return value into eax
70
71     mov     esp, ebp
72     pop     ebp
73     ret

```

---

read.asm

### 6.3.6 Finding primes

This final example looks at finding prime numbers again. This implementation is more efficient than the previous one. It stores the primes it has found in an array and only divides by the previous primes it has found instead of every odd number to find new primes.

One other difference is that it computes the square root of the guess for the next prime to determine at what point it can stop searching for factors. It alters the coprocessor control word so that when it stores the square root as an integer, it truncates instead of rounding. This is controlled by bits 10 and 11 of the control word. These bits are called the RC (Rounding Control) bits. If they are both 0 (the default), the coprocessor rounds when converting to integer. If they are both 1, the coprocessor truncates integer conversions. Notice that the routine is careful to save the original control word and restore it before it returns.

Here is the C driver program:

---

```

                                fprime.c
1  #include <stdio.h>
2  #include <stdlib.h>
3  /*
4   * function find_primes
5   * finds the indicated number of primes
6   * Parameters:
7   *   a — array to hold primes
8   *   n — how many primes to find
9   */
10 extern void find_primes ( int * a, unsigned n );
11
12 int main()
13 {
14     int status;
15     unsigned i;
16     unsigned max;
17     int * a;
18
19     printf ("How many primes do you wish to find? ");
20     scanf ("%u", &max);
21
22     a = calloc ( sizeof(int), max);
23
24     if ( a ) {
25
26         find_primes (a,max);
27
28         /* print out the last 20 primes found */
29         for(i = ( max > 20 ) ? max - 20 : 0; i < max; i++ )
30             printf ("%3d %d\n", i+1, a[i]);

```

```

31
32     free(a);
33     status = 0;
34 }
35 else {
36     fprintf (stderr , "Can not create array of %u ints\n", max);
37     status = 1;
38 }
39
40 return status;
41 }

```

---

fprime.c

---

Here is the assembly routine:

---

```

1  segment .text          prime2.asm
2      global _find_primes
3  ;
4  ; function find_primes
5  ; finds the indicated number of primes
6  ; Parameters:
7  ;   array - array to hold primes
8  ;   n_find - how many primes to find
9  ; C Prototype:
10 ;extern void find_primes( int * array, unsigned n_find )
11 ;
12 %define array          ebp + 8
13 %define n_find         ebp + 12
14 %define n              ebp - 4          ; number of primes found so far
15 %define isqrt          ebp - 8          ; floor of sqrt of guess
16 %define orig_cntl_wd   ebp - 10         ; original control word
17 %define new_cntl_wd    ebp - 12         ; new control word
18
19 _find_primes:
20     enter    12,0          ; make room for local variables
21
22     push    ebx            ; save possible register variables
23     push    esi
24
25     fstcw   word [orig_cntl_wd]        ; get current control word
26     mov     ax, [orig_cntl_wd]

```

```

27         or      ax, 0C00h                ; set rounding bits to 11 (truncate)
28         mov     [new_cntl_wd], ax
29         fldcw   word [new_cntl_wd]
30
31         mov     esi, [array]              ; esi points to array
32         mov     dword [esi], 2            ; array[0] = 2
33         mov     dword [esi + 4], 3        ; array[1] = 3
34         mov     ebx, 5                    ; ebx = guess = 5
35         mov     dword [n], 2              ; n = 2
36     ;
37     ; This outer loop finds a new prime each iteration, which it adds to the
38     ; end of the array. Unlike the earlier prime finding program, this function
39     ; does not determine primeness by dividing by all odd numbers. It only
40     ; divides by the prime numbers that it has already found. (That's why they
41     ; are stored in the array.)
42     ;
43     while_limit:
44         mov     eax, [n]
45         cmp     eax, [n_find]              ; while ( n < n_find )
46         jnb     short quit_limit
47
48         mov     ecx, 1                    ; ecx is used as array index
49         push    ebx                        ; store guess on stack
50         fld     dword [esp]                ; load guess onto coprocessor stack
51         pop     ebx                        ; get guess off stack
52         fsqrt
53         fistp   dword [isqrt]              ; isqrt = floor(sqrt(guess))
54     ;
55     ; This inner loop divides guess (ebx) by earlier computed prime numbers
56     ; until it finds a prime factor of guess (which means guess is not prime)
57     ; or until the prime number to divide is greater than floor(sqrt(guess))
58     ;
59     while_factor:
60         mov     eax, dword [esi + 4*ecx]    ; eax = array[ecx]
61         cmp     eax, [isqrt]                ; while ( isqrt < array[ecx]
62         jnbe    short quit_factor_prime
63         mov     eax, ebx
64         xor     edx, edx
65         div     dword [esi + 4*ecx]
66         or      edx, edx                    ; && guess % array[ecx] != 0 )
67         jz      short quit_factor_not_prime
68         inc     ecx                        ; try next prime

```

```
69         jmp     short while_factor
70
71     ;
72     ; found a new prime !
73     ;
74     quit_factor_prime:
75         mov     eax, [n]
76         mov     dword [esi + 4*eax], ebx        ; add guess to end of array
77         inc     eax
78         mov     [n], eax                        ; inc n
79
80     quit_factor_not_prime:
81         add     ebx, 2                          ; try next odd number
82         jmp     short while_limit
83
84     quit_limit:
85
86         fldcw   word [orig_cntl_wd]            ; restore control word
87         pop     esi                            ; restore register variables
88         pop     ebx
89
90         leave
91         ret
```

---

```

1  global _dmax
2
3  segment .text
4  ; function _dmax
5  ; returns the larger of its two double arguments
6  ; C prototype
7  ; double dmax( double d1, double d2 )
8  ; Parameters:
9  ;   d1   - first double
10 ;   d2   - second double
11 ; Return value:
12 ;   larger of d1 and d2 (in ST0)
13 %define d1   ebp+8
14 %define d2   ebp+16
15 _dmax:
16     enter    0, 0
17
18     fld     qword [d2]
19     fld     qword [d1]           ; ST0 = d1, ST1 = d2
20     fcomip  st1                 ; ST0 = d2
21     jna     short d2_bigger
22     fcomp   st0                 ; pop d2 from stack
23     fld     qword [d1]           ; ST0 = d1
24     jmp     short exit
25 d2_bigger:                       ; if d2 is max, nothing to do
26 exit:
27     leave
28     ret

```

Figure 6.7: FCOMIP example



```
1 segment .data
2 x          dq  2.75          ; converted to double format
3 five       dw  5
4
5 segment .text
6     fild    dword [five]     ; ST0 = 5
7     fld     qword [x]        ; ST0 = 2.75, ST1 = 5
8     fscale                     ; ST0 = 2.75 * 32, ST1 = 5
```

Figure 6.8: FSCALE example