recommend the one coauthored by James Gosling, the creator of the language [5]) describe the data formats and arithmetic operations supported by Java.

Most books on logic design [58, 116] have a section on encodings and arithmetic operations. Such books describe different ways of implementing arithmetic circuits. Overton's book on IEEE floating point [82] provides a detailed description of the format as well as the properties from the perspective of a numerical applications programmer.

## Homework Problems

### 2.55 ◆
Compile and run the sample code that uses show_bytes (file show-bytes.c) on different machines to which you have access. Determine the byte orderings used by these machines.

### 2.56 ◆
Try running the code for show_bytes for different sample values.

### 2.57 ◆
Write procedures show_short, show_long, and show_double that print the byte representations of C objects of types short, long, and double, respectively. Try these out on several machines.

### 2.58 ◆◆
Write a procedure is_little_endian that will return 1 when compiled and run on a little-endian machine, and will return 0 when compiled and run on a big-endian machine. This program should run on any machine, regardless of its word size.

### 2.59 ◆◆
Write a C expression that will yield a word consisting of the least significant byte of x and the remaining bytes of y. For operands x = 0x89ABCDEF and y = 0x76543210, this would give 0x765432EF.

### 2.60 ◆◆
Suppose we number the bytes in a $w$-bit word from 0 (least significant) to $w/8 - 1$ (most significant). Write code for the following C function, which will return an unsigned value in which byte i of argument x has been replaced by byte b:

```
unsigned replace_byte (unsigned x, int i, unsigned char b);
```

Here are some examples showing how the function should work:

```
replace_byte(0x12345678, 2, 0xAB) --> 0x12AB5678
replace_byte(0x12345678, 0, 0xAB) --> 0x123456AB
```

### Bit-Level Integer Coding Rules

In several of the following problems, we will artificially restrict what programming constructs you can use to help you gain a better understanding of the bit-level,

logic, and arithmetic operations of C. In answering these problems, your code must follow these rules:

- Assumptions
  - Integers are represented in two's-complement form.
  - Right shifts of signed data are performed arithmetically.
  - Data type int is $w$ bits long. For some of the problems, you will be given a specific value for $w$, but otherwise your code should work as long as $w$ is a multiple of 8. You can use the expression sizeof(int)<<3 to compute $w$.
- Forbidden
  - Conditionals (if or ?:), loops, switch statements, function calls, and macro invocations.
  - Division, modulus, and multiplication.
  - Relative comparison operators (<, >, <=, and >=).
- Allowed operations
  - All bit-level and logic operations.
  - Left and right shifts, but only with shift amounts between 0 and $w-1$.
  - Addition and subtraction.
  - Equality (==) and inequality (!=) tests. (Some of the problems do not allow these.)
  - Integer constants INT_MIN and INT_MAX.
  - Casting between data types int and unsigned, either explicitly or implicitly.

Even with these rules, you should try to make your code readable by choosing descriptive variable names and using comments to describe the logic behind your solutions. As an example, the following code extracts the most significant byte from integer argument x:

```
/* Get most significant byte from x */
int get_msb(int x) {
    /* Shift by w-8 */
    int shift_val = (sizeof(int)-1)<<3;
    /* Arithmetic shift */
    int xright = x >> shift_val;
    /* Zero all but LSB */
    return xright & 0xFF;
}
```

**2.61** ◆◆
Write C expressions that evaluate to 1 when the following conditions are true and to 0 when they are false. Assume x is of type int.

A. Any bit of x equals 1.

B. Any bit of x equals 0.

    C. Any bit in the least significant byte of x equals 1.

    D. Any bit in the most significant byte of x equals 0.

Your code should follow the bit-level integer coding rules (page 164), with the additional restriction that you may not use equality (==) or inequality (!=) tests.

**2.63** ◆◆◆

Write a function int_shifts_are_arithmetic() that yields 1 when run on a machine that uses arithmetic right shifts for data type int and yields 0 otherwise. Your code should work on a machine with any word size. Test your code on several machines.

**2.63** ◆◆◆

Fill in code for the following C functions. Function srl performs a logical right shift using an arithmetic right shift (given by value xsra), followed by other operations not including right shifts or division. Function sra performs an arithmetic right shift using a logical right shift (given by value xsrl), followed by other operations not including right shifts or division. You may use the computation 8*sizeof(int) to determine $w$, the number of bits in data type int. The shift amount k can range from 0 to $w - 1$.

```
unsigned srl(unsigned x, int k) {
    /* Perform shift arithmetically */
    unsigned xsra = (int) x >> k;
.
.
.
.
.
.
}

int sra(int x, int k) {
    /* Perform shift logically */
    int xsrl = (unsigned) x >> k;
.
.
.
.
.
.
}
```

**2.64** ◆

Write code to implement the following function:

```
/* Return 1 when any odd bit of x equals 1; 0 otherwise.
   Assume w=32 */
int any_odd_one(unsigned x);
```

    Your function should follow the bit-level integer coding rules (page 164), except that you may assume that data type int has $w = 32$ bits.

**2.65** ◆◆◆◆

Write code to implement the following function:

```
/* Return 1 when x contains an odd number of 1s; 0 otherwise.
   Assume w=32 */
int odd_ones(unsigned x);
```

Your function should follow the bit-level integer coding rules (page 164), except that you may assume that data type int has $w = 32$ bits.

Your code should contain a total of at most 12 arithmetic, bitwise, and logical operations.

**2.66** ◆◆◆

Write code to implement the following function:

```
/*
 * Generate mask indicating leftmost 1 in x.  Assume w=32.
 * For example, 0xFF00 -> 0x8000, and 0x6600 --> 0x4000.
 * If x = 0, then return 0.
 */
int leftmost_one(unsigned x);
```

Your function should follow the bit-level integer coding rules (page 164), except that you may assume that data type int has $w = 32$ bits.

Your code should contain a total of at most 15 arithmetic, bitwise, and logical operations.

*Hint:* First transform x into a bit vector of the form $[0 \cdots 011 \cdots 1]$.

**2.67** ◆◆

You are given the task of writing a procedure int_size_is_32() that yields 1 when run on a machine for which an int is 32 bits, and yields 0 otherwise. You are not allowed to use the sizeof operator. Here is a first attempt:

```
1    /* The following code does not run properly on some machines */
2    int bad_int_size_is_32() {
3        /* Set most significant bit (msb) of 32-bit machine */
4        int set_msb = 1 << 31;
5        /* Shift past msb of 32-bit word */
6        int beyond_msb = 1 << 32;
7
8        /* set_msb is nonzero when word size >= 32
9           beyond_msb is zero when word size <= 32  */
10       return set_msb && !beyond_msb;
11   }
```

When compiled and run on a 32-bit SUN SPARC, however, this procedure returns 0. The following compiler message gives us an indication of the problem:

```
warning: left shift count >= width of type
```

    A. In what way does our code fail to comply with the C standard?

    B. Modify the code to run properly on any machine for which data type `int` is at least 32 bits.

    C. Modify the code to run properly on any machine for which data type `int` is at least 16 bits.

**2.68** ◆◆

Write code for a function with the following prototype:

```
/*
 * Mask with least signficant n bits set to 1
 * Examples: n = 6 --> 0x3F, n = 17 --> 0x1FFFF
 * Assume 1 <= n <= w
 */
int lower_one_mask(int n);
```

    Your function should follow the bit-level integer coding rules (page 164). Be careful of the case $n = w$.

**2.69** ◆◆◆

Write code for a function with the following prototype:

```
/*
 * Do rotating left shift.  Assume 0 <= n < w
 * Examples when x = 0x12345678 and w = 32:
 *    n=4 -> 0x23456781, n=20 -> 0x67812345
 */
unsigned rotate_left(unsigned x, int n);
```

    Your function should follow the bit-level integer coding rules (page 164). Be careful of the case $n = 0$.

**2.70** ◆◆

Write code for the function with the following prototype:

```
/*
 * Return 1 when x can be represented as an n-bit, 2's-complement
 * number; 0 otherwise
 * Assume 1 <= n <= w
 */
int fits_bits(int x, int n);
```

    Your function should follow the bit-level integer coding rules (page 164).

**2.71** ◆

You just started working for a company that is implementing a set of procedures to operate on a data structure where 4 signed bytes are packed into a 32-bit unsigned. Bytes within the word are numbered from 0 (least significant) to 3

(most significant). You have been assigned the task of implementing a function for a machine using two's-complement arithmetic and arithmetic right shifts with the following prototype:

```
/* Declaration of data type where 4 bytes are packed
   into an unsigned */
typedef unsigned packed_t;

/* Extract byte from word.  Return as signed integer */
int xbyte(packed_t word, int bytenum);
```

That is, the function will extract the designated byte and sign extend it to be a 32-bit `int`.

Your predecessor (who was fired for incompetence) wrote the following code:

```
/* Failed attempt at xbyte */
int xbyte(packed_t word, int bytenum)
{
    return (word >> (bytenum << 3)) & 0xFF;
}
```

A. What is wrong with this code?

B. Give a correct implementation of the function that uses only left and right shifts, along with one subtraction.

## 2.72 ◆◆

You are given the task of writing a function that will copy an integer `val` into a buffer `buf`, but it should do so only if enough space is available in the buffer.

Here is the code you write:

```
/* Copy integer into buffer if space is available */
/* WARNING: The following code is buggy */
void copy_int(int val, void *buf, int maxbytes) {
    if (maxbytes-sizeof(val) >= 0)
            memcpy(buf, (void *) &val, sizeof(val));
}
```

This code makes use of the library function `memcpy`. Although its use is a bit artificial here, where we simply want to copy an `int`, it illustrates an approach commonly used to copy larger data structures.

You carefully test the code and discover that it *always* copies the value to the buffer, even when `maxbytes` is too small.

A. Explain why the conditional test in the code always succeeds. *Hint:* The `sizeof` operator returns a value of type `size_t`.

B. Show how you can rewrite the conditional test to make it work properly.

**2.73** ◆◆

Write code for a function with the following prototype:

```
/* Addition that saturates to TMin or TMax */
int saturating_add(int x, int y);
```

Instead of overflowing the way normal two's-complement addition does, saturating addition returns *TMax* when there would be positive overflow, and *TMin* when there would be negative overflow. Saturating arithmetic is commonly used in programs that perform digital signal processing.

Your function should follow the bit-level integer coding rules (page 164).

**2.74** ◆◆

Write a function with the following prototype:

```
/* Determine whether arguments can be subtracted without overflow */
int tsub_ok(int x, int y);
```

This function should return 1 if the computation x-y does not overflow.

**2.75** ◆◆◆

Suppose we want to compute the complete $2w$-bit representation of $x \cdot y$, where both $x$ and $y$ are unsigned, on a machine for which data type unsigned is $w$ bits. The low-order $w$ bits of the product can be computed with the expression x*y, so we only require a procedure with prototype

```
unsigned unsigned_high_prod(unsigned x, unsigned y);
```

that computes the high-order $w$ bits of $x \cdot y$ for unsigned variables.

We have access to a library function with prototype

```
int signed_high_prod(int x, int y);
```

that computes the high-order $w$ bits of $x \cdot y$ for the case where $x$ and $y$ are in two's-complement form. Write code calling this procedure to implement the function for unsigned arguments. Justify the correctness of your solution.

*Hint:* Look at the relationship between the signed product $x \cdot y$ and the unsigned product $x' \cdot y'$ in the derivation of Equation 2.18.

**2.76** ◆

The library function calloc has the following declaration:

```
void *calloc(size_t nmemb, size_t size);
```

According to the library documentation, "The calloc function allocates memory for an array of nmemb elements of size bytes each. The memory is set to zero. If nmemb or size is zero, then calloc returns NULL."

Write an implementation of calloc that performs the allocation by a call to malloc and sets the memory to zero via memset. Your code should not have any vulnerabilities due to arithmetic overflow, and it should work correctly regardless of the number of bits used to represent data of type size_t.

As a reference, functions malloc and memset have the following declarations:

```
void *malloc(size_t size);
void *memset(void *s, int c, size_t n);
```

**2.77** ◆◆

Suppose we are given the task of generating code to multiply integer variable x by various different constant factors $K$. To be efficient, we want to use only the operations +, -, and <<. For the following values of $K$, write C expressions to perform the multiplication using at most three operations per expression.

A.  $K = 17$

B.  $K = -7$

C.  $K = 60$

D.  $K = -112$

**2.78** ◆◆

Write code for a function with the following prototype:

```
/* Divide by power of 2. Assume 0 <= k < w-1 */
int divide_power2(int x, int k);
```

The function should compute $x/2^k$ with correct rounding, and it should follow the bit-level integer coding rules (page 164).

**2.79** ◆◆

Write code for a function mul3div4 that, for integer argument x, computes $3 * x/4$ but follows the bit-level integer coding rules (page 164). Your code should replicate the fact that the computation 3*x can cause overflow.

**2.80** ◆◆◆

Write code for a function threefourths that, for integer argument x, computes the value of $\frac{3}{4}x$, rounded toward zero. It should not overflow. Your function should follow the bit-level integer coding rules (page 164).

**2.81** ◆◆

Write C expressions to generate the bit patterns that follow, where $a^k$ represents $k$ repetitions of symbol $a$. Assume a $w$-bit data type. Your code may contain references to parameters j and k, representing the values of $j$ and $k$, but not a parameter representing $w$.

A.  $1^{w-k}0^k$

B.  $0^{w-k-j}1^k0^j$

**2.82** ◆

We are running programs where values of type int are 32 bits. They are represented in two's complement, and they are right shifted arithmetically. Values of type unsigned are also 32 bits.

We generate arbitrary values x and y, and convert them to unsigned values as follows:

```
/* Create some arbitrary values */
int x = random();
int y = random();
/* Convert to unsigned */
unsigned ux = (unsigned) x;
unsigned uy = (unsigned) y;
```

For each of the following C expressions, you are to indicate whether or not the expression *always* yields 1. If it always yields 1, describe the underlying mathematical principles. Otherwise, give an example of arguments that make it yield 0.

A. `(x<y) == (-x>-y)`

B. `((x+y)<<4) + y-x == 17*y+15*x`

C. `~x+~y+1 == ~(x+y)`

D. `(ux-uy) == -(unsigned)(y-x)`

E. `((x >> 2) << 2) <= x`

### 2.83 ◆◆

Consider numbers having a binary representation consisting of an infinite string of the form $0.y\,y\,y\,y\,y\,y\cdots$, where $y$ is a $k$-bit sequence. For example, the binary representation of $\frac{1}{3}$ is $0.01010101\cdots$ ($y = 01$), while the representation of $\frac{1}{5}$ is $0.001100110011\cdots$ ($y = 0011$).

A. Let $Y = B2U_k(y)$, that is, the number having binary representation $y$. Give a formula in terms of $Y$ and $k$ for the value represented by the infinite string. *Hint:* Consider the effect of shifting the binary point $k$ positions to the right.

B. What is the numeric value of the string for the following values of $y$?
   (a) 101
   (b) 0110
   (c) 010011

### 2.84 ◆

Fill in the return value for the following procedure, which tests whether its first argument is less than or equal to its second. Assume the function f2u returns an unsigned 32-bit number having the same bit representation as its floating-point argument. You can assume that neither argument is *NaN*. The two flavors of zero, $+0$ and $-0$, are considered equal.

```
int float_le(float x, float y) {
    unsigned ux = f2u(x);
    unsigned uy = f2u(y);
```

```
    /* Get the sign bits */
    unsigned sx = ux >> 31;
    unsigned sy = uy >> 31;

    /* Give an expression using only ux, uy, sx, and sy */
    return          ;
}
```

**2.85** ◆

Given a floating-point format with a $k$-bit exponent and an $n$-bit fraction, write formulas for the exponent $E$, the significand $M$, the fraction $f$, and the value $V$ for the quantities that follow. In addition, describe the bit representation.

   A. The number 7.0

   B. The largest odd integer that can be represented exactly

   C. The reciprocal of the smallest positive normalized value

**2.86** ◆

Intel-compatible processors also support an "extended-precision" floating-point format with an 80-bit word divided into a sign bit, $k = 15$ exponent bits, a single *integer* bit, and $n = 63$ fraction bits. The integer bit is an explicit copy of the implied bit in the IEEE floating-point representation. That is, it equals 1 for normalized values and 0 for denormalized values. Fill in the following table giving the approximate values of some "interesting" numbers in this format:

| Description | Extended precision | |
| --- | --- | --- |
| | Value | Decimal |
| Smallest positive denormalized | _____ | _____ |
| Smallest positive normalized | _____ | _____ |
| Largest normalized | _____ | _____ |

This format can be used in C programs compiled for Intel-compatible machines by declaring the data to be of type `long double`. However, it forces the compiler to generate code based on the legacy 8087 floating-point instructions. The resulting program will most likely run much slower than would be the case for data type `float` or `double`.

**2.87** ◆

The 2008 version of the IEEE floating-point standard, named IEEE 754-2008, includes a 16-bit "half-precision" floating-point format. It was originally devised by computer graphics companies for storing data in which a higher dynamic range is required than can be achieved with 16-bit integers. This format has 1 sign bit, 5 exponent bits ($k = 5$), and 10 fraction bits ($n = 10$). The exponent bias is $2^{5-1} - 1 = 15$.

Fill in the table that follows for each of the numbers given, with the following instructions for each column:

Hex: The four hexadecimal digits describing the encoded form.

$M$: The value of the significand. This should be a number of the form $x$ or $\frac{x}{y}$, where $x$ is an integer and $y$ is an integral power of 2. Examples include 0, $\frac{67}{64}$, and $\frac{1}{256}$.

$E$: The integer value of the exponent.

$V$: The numeric value represented. Use the notation $x$ or $x \times 2^z$, where $x$ and $z$ are integers.

$D$: The (possibly approximate) numerical value, as is printed using the %f formatting specification of printf.

As an example, to represent the number $\frac{7}{8}$, we would have $s = 0$, $M = \frac{7}{4}$, and $E = -1$. Our number would therefore have an exponent field of $01110_2$ (decimal value $15 - 1 = 14$) and a significand field of $1100000000_2$, giving a hex representation 3B00. The numerical value is 0.875.

You need not fill in entries marked —.

| Description | Hex | $M$ | $E$ | $V$ | $D$ |
|---|---|---|---|---|---|
| $-0$ | _____ | _____ | _____ | $-0$ | $-0.0$ |
| Smallest value $> 2$ | _____ | _____ | _____ | _____ | _____ |
| 512 | _____ | _____ | _____ | 512 | 512.0 |
| Largest denormalized | _____ | _____ | _____ | _____ | _____ |
| $-\infty$ | _____ | — | — | $-\infty$ | $-\infty$ |
| Number with hex representation 3BB0 | 3BB0 | _____ | _____ | _____ | _____ |

### 2.88 ◆◆

Consider the following two 9-bit floating-point representations based on the IEEE floating-point format.

1. **Format A**
   - There is 1 sign bit.
   - There are $k = 5$ exponent bits. The exponent bias is 15.
   - There are $n = 3$ fraction bits.

2. **Format B**
   - There is 1 sign bit.
   - There are $k = 4$ exponent bits. The exponent bias is 7.
   - There are $n = 4$ fraction bits.

In the following table, you are given some bit patterns in format A, and your task is to convert them to the closest value in format B. If rounding is necessary you should *round toward* $+\infty$. In addition, give the values of numbers given by the format A and format B bit patterns. Give these as whole numbers (e.g., 17) or as fractions (e.g., 17/64 or $17/2^6$).

| Format A | | Format B | |
|---|---|---|
| Bits | Value | Bits | Value |
| 1 01111 001 | $\frac{-9}{8}$ | 1 0111 0010 | $\frac{-9}{8}$ |
| 0 10110 011 | _____ | _____ | _____ |
| 1 00111 010 | _____ | _____ | _____ |
| 0 00000 111 | _____ | _____ | _____ |
| 1 11100 000 | _____ | _____ | _____ |
| 0 10111 100 | _____ | _____ | _____ |

**2.89** ◆

We are running programs on a machine where values of type int have a 32-bit two's-complement representation. Values of type float use the 32-bit IEEE format, and values of type double use the 64-bit IEEE format.

We generate arbitrary integer values x, y, and z, and convert them to values of type double as follows:

```
/* Create some arbitrary values */
int x = random();
int y = random();
int z = random();
/* Convert to double */
double   dx = (double) x;
double   dy = (double) y;
double   dz = (double) z;
```

For each of the following C expressions, you are to indicate whether or not the expression *always* yields 1. If it always yields 1, describe the underlying mathematical principles. Otherwise, give an example of arguments that make it yield 0. Note that you cannot use an IA32 machine running GCC to test your answers, since it would use the 80-bit extended-precision representation for both float and double.

  A.  (float) x == (float) dx

  B.  dx – dy == (double) (x-y)

  C.  (dx + dy) + dz == dx + (dy + dz)

  D.  (dx * dy) * dz == dx * (dy * dz)

  E.  dx / dx == dz / dz

**2.90** ◆

You have been assigned the task of writing a C function to compute a floating-point representation of $2^x$. You decide that the best way to do this is to directly construct the IEEE single-precision representation of the result. When $x$ is too small, your routine will return 0.0. When $x$ is too large, it will return $+\infty$. Fill in the blank portions of the code that follows to compute the correct result. Assume the

function u2f returns a floating-point value having an identical bit representation as its unsigned argument.

```
float fpwr2(int x)
{
    /* Result exponent and fraction */
    unsigned exp, frac;
    unsigned u;

    if (x < _____) {
        /* Too small.  Return 0.0 */
        exp = _____;
        frac = _____;
    } else if (x < _____) {
        /* Denormalized result */
        exp = _____;
        frac = _____;
    } else if (x < _____) {
        /* Normalized result. */
        exp = _____;
        frac = _____;
    } else {
        /* Too big.  Return +oo */
        exp = _____;
        frac = _____;
    }

    /* Pack exp and frac into 32 bits */
    u = exp << 23 | frac;
    /* Return as float */
    return u2f(u);
}
```

**2.91** ◆

Around 250 B.C., the Greek mathematician Archimedes proved that $\frac{223}{71} < \pi < \frac{22}{7}$. Had he had access to a computer and the standard library <math.h>, he would have been able to determine that the single-precision floating-point approximation of $\pi$ has the hexadecimal representation 0x40490FDB. Of course, all of these are just approximations, since $\pi$ is not rational.

A. What is the fractional binary number denoted by this floating-point value?

B. What is the fractional binary representation of $\frac{22}{7}$? *Hint:* See Problem 2.83.

C. At what bit position (relative to the binary point) do these two approximations to $\pi$ diverge?

## Bit-Level Floating-Point Coding Rules

In the following problems, you will write code to implement floating-point functions, operating directly on bit-level representations of floating-point numbers. Your code should exactly replicate the conventions for IEEE floating-point operations, including using round-to-even mode when rounding is required.

To this end, we define data type float_bits to be equivalent to unsigned:

```
/* Access bit-level representation floating-point number */
typedef unsigned float_bits;
```

Rather than using data type float in your code, you will use float_bits. You may use both int and unsigned data types, including unsigned and integer constants and operations. You may not use any unions, structs, or arrays. Most significantly, you may not use any floating-point data types, operations, or constants. Instead, your code should perform the bit manipulations that implement the specified floating-point operations.

The following function illustrates the use of these coding rules. For argument $f$, it returns $\pm 0$ if $f$ is denormalized (preserving the sign of $f$), and returns $f$ otherwise.

```
/* If f is denorm, return 0.  Otherwise, return f */
float_bits float_denorm_zero(float_bits f) {
    /* Decompose bit representation into parts */
    unsigned sign = f>>31;
    unsigned exp =  f>>23 & 0xFF;
    unsigned frac = f     & 0x7FFFFF;
    if (exp == 0) {
        /* Denormalized.  Set fraction to 0 */
        frac = 0;
    }
    /* Reassemble bits */
    return (sign << 31) | (exp << 23) | frac;
}
```

### 2.92 ◆◆

Following the bit-level floating-point coding rules, implement the function with the following prototype:

```
/* Compute -f.  If f is NaN, then return f. */
float_bits float_negate(float_bits f);
```

For floating-point number $f$, this function computes $-f$. If $f$ is *NaN*, your function should simply return $f$.

Test your function by evaluating it for all $2^{32}$ values of argument f and comparing the result to what would be obtained using your machine's floating-point operations.

**2.93** ◆◆

Following the bit-level floating-point coding rules, implement the function with the following prototype:

```
/* Compute |f|.  If f is NaN, then return f. */
float_bits float_absval(float_bits f);
```

For floating-point number $f$, this function computes $|f|$. If $f$ is *NaN*, your function should simply return $f$.

Test your function by evaluating it for all $2^{32}$ values of argument f and comparing the result to what would be obtained using your machine's floating-point operations.

**2.94** ◆◆◆

Following the bit-level floating-point coding rules, implement the function with the following prototype:

```
/* Compute 2*f.  If f is NaN, then return f. */
float_bits float_twice(float_bits f);
```

For floating-point number $f$, this function computes $2.0 \cdot f$. If $f$ is *NaN*, your function should simply return $f$.

Test your function by evaluating it for all $2^{32}$ values of argument f and comparing the result to what would be obtained using your machine's floating-point operations.

**2.95** ◆◆◆

Following the bit-level floating-point coding rules, implement the function with the following prototype:

```
/* Compute 0.5*f.  If f is NaN, then return f. */
float_bits float_half(float_bits f);
```

For floating-point number $f$, this function computes $0.5 \cdot f$. If $f$ is *NaN*, your function should simply return $f$.

Test your function by evaluating it for all $2^{32}$ values of argument f and comparing the result to what would be obtained using your machine's floating-point operations.

**2.96** ◆◆◆◆

Following the bit-level floating-point coding rules, implement the function with the following prototype:

```
/*
 * Compute (int) f.
 * If conversion causes overflow or f is NaN, return 0x80000000
 */
int float_f2i(float_bits f);
```

For floating-point number $f$, this function computes (int) $f$. Your function should round toward zero. If $f$ cannot be represented as an integer (e.g., it is out of range, or it is *NaN*), then the function should return 0x80000000.

Test your function by evaluating it for all $2^{32}$ values of argument f and comparing the result to what would be obtained using your machine's floating-point operations.

### 2.97 ◆◆◆◆

Following the bit-level floating-point coding rules, implement the function with the following prototype:

```
/* Compute (float) i */
float_bits float_i2f(int i);
```

For argument i, this function computes the bit-level representation of (float) i.

Test your function by evaluating it for all $2^{32}$ values of argument f and comparing the result to what would be obtained using your machine's floating-point operations.

## Solutions to Practice Problems

### Solution to Problem 2.1  (page 73)

Understanding the relation between hexadecimal and binary formats will be important once we start looking at machine-level programs. The method for doing these conversions is in the text, but it takes a little practice to become familiar.

A. 0x25B9D2 to binary:

| Hexadecimal | 2 | 5 | B | 9 | D | 2 |
|---|---|---|---|---|---|---|
| Binary | 0010 | 0101 | 1101 | 1001 | 1101 | 0010 |

B. Binary 1100 1001 0111 1011 to hexadecimal:

| Binary | 1100 | 1001 | 0111 | 1011 |
|---|---|---|---|---|
| Hexadecimal | C | 9 | 7 | B |

C. 0xA8B3D to binary:

| Hexadecimal | A | 8 | B | 3 | D |
|---|---|---|---|---|---|
| Binary | 1010 | 1000 | 1011 | 0011 | 1101 |

D. Binary 11 0010 0010 1101 1001 0110 to hexadecimal:

| Binary | 11 | 0010 | 0010 | 1101 | 1001 | 0110 |
|---|---|---|---|---|---|---|
| Hexadecimal | 3 | 2 | 2 | D | 9 | 6 |

### Solution to Problem 2.2  (page 73)

This problem gives you a chance to think about powers of 2 and their hexadecimal representations.