

**1. Introduction.** The subroutines below are used to simulate 64-bit MMIX arithmetic on an old-fashioned 32-bit computer—like the one the author had when he wrote MMIXAL and the first MMIX simulators in 1998 and 1999. All operations are fabricated from 32-bit arithmetic, including a full implementation of the IEEE floating point standard, assuming only that the C compiler has a 32-bit unsigned integer type.

Some day 64-bit machines will be commonplace and the awkward manipulations of the present program will look quite archaic. Interested readers who have such computers will be able to convert the code to a pure 64-bit form without difficulty, thereby obtaining much faster and simpler routines. Meanwhile, however, we can simulate the future and hope for continued progress.

This program module has a simple structure, intended to make it suitable for loading with MMIX simulators and assemblers.

```
#include <stdio.h>
#include <string.h>
#include <ctype.h>
  ⟨ Stuff for C preprocessor 2 ⟩
typedef enum { false, true } bool;
  ⟨ Tetrabyte and octabyte type definitions 3 ⟩
  ⟨ Other type definitions 36 ⟩
  ⟨ Global variables 4 ⟩
  ⟨ Subroutines 5 ⟩
```

**2.** Subroutines of this program are declared first with a prototype, as in ANSI C, then with an old-style C function definition. Here are some preprocessor commands that make this work correctly with both new-style and old-style compilers.

```
⟨ Stuff for C preprocessor 2 ⟩ ≡
#define __STDC__
#define ARGS(list) list
#else
#define ARGS(list) ()
#endif
```

This code is used in section 1.

**3.** The definition of type **tetra** should be changed, if necessary, so that it represents an unsigned 32-bit integer.

```
⟨ Tetrabyte and octabyte type definitions 3 ⟩ ≡
typedef unsigned int tetra;    /* for systems conforming to the LP-64 data model */
typedef struct {
    tetra h,l;
} octa;    /* two tetrabytes make one octabyte */
```

This code is used in section 1.

**4.** `#define sign_bit ((unsigned)#80000000)`

```
⟨ Global variables 4 ⟩ ≡
octa zero_octa;    /* zero_octa.h = zero_octa.l = 0 */
octa neg_one = {-1,-1};    /* neg_one.h = neg_one.l = -1 */
octa inf_octa = {#7ff00000,0};    /* floating point +∞ */
octa standard_NaN = {#7ff80000,0};    /* floating point NaN(.5) */
```

See also sections 9, 30, 32, 69, and 75.

This code is used in section 1.

5. It's easy to add and subtract octabytes, if we aren't terribly worried about speed.

⟨Subroutines 5⟩ ≡

```

octa oplus ARGS((octa,octa));
    octa oplus(y,z)    /* compute  $y + z$  */
    octa y,z;
{ octa x;
   $x.h = y.h + z.h$ ;  $x.l = y.l + z.l$ ;
  if ( $x.l < y.l$ )  $x.h++$ ;
  return x;
}

octa ominus ARGS((octa,octa));
    octa ominus(y,z)    /* compute  $y - z$  */
    octa y,z;
{ octa x;
   $x.h = y.h - z.h$ ;  $x.l = y.l - z.l$ ;
  if ( $x.l > y.l$ )  $x.h--$ ;
  return x;
}

```

See also sections 6, 7, 8, 12, 13, 24, 25, 26, 27, 28, 29, 31, 34, 37, 38, 39, 40, 41, 44, 46, 50, 54, 60, 61, 62, 68, 82, 85, 86, 88, 89, 91, and 93.

This code is used in section 1.

6. In the following subroutine, *delta* is a signed quantity that is assumed to fit in a signed tetrabyte.

⟨Subroutines 5⟩ +≡

```

octa incr ARGS((octa,int));
    octa incr(y,delta)    /* compute  $y + \delta$  */
    octa y;
    int delta;
{ octa x;
   $x.h = y.h$ ;  $x.l = y.l + delta$ ;
  if ( $delta \geq 0$ ) {
    if ( $x.l < y.l$ )  $x.h++$ ;
  } else if ( $x.l > y.l$ )  $x.h--$ ;
  return x;
}

```

7. Left and right shifts are only a bit more difficult.

⟨Subroutines 5⟩ +≡

```

octa shift_left ARGS((octa,int));
    octa shift_left(y, s)      /* shift left by s bits, where  $0 \leq s \leq 64$  */
    octa y;
    int s;
{
    while ( $s \geq 32$ ) y.h = y.l, y.l = 0, s -= 32;
    if (s) { register tetra yhl = y.h << s, ylh = y.l >> (32 - s);
        y.h = yhl + ylh; y.l <<= s;
    }
    return y;
}

octa shift_right ARGS((octa,int,int));
    octa shift_right(y, s, u)    /* shift right, arithmetically if u = 0 */
    octa y;
    int s, u;
{
    while ( $s \geq 32$ ) y.l = y.h, y.h = (u ? 0 : -(y.h >> 31)), s -= 32;
    if (s) { register tetra yhl = y.h << (32 - s), ylh = y.l >> s;
        y.h = (u ? 0 : -(y.h >> 31)) << (32 - s) + (y.h >> s); y.l = yhl + ylh;
    }
    return y;
}

```

**8. Multiplication.** We need to multiply two unsigned 64-bit integers, obtaining an unsigned 128-bit product. It is easy to do this on a 32-bit machine by using Algorithm 4.3.1M of *Seminumerical Algorithms*, with  $b = 2^{16}$ .

The following subroutine returns the lower half of the product, and puts the upper half into a global octabyte called *aux*.

⟨Subroutines 5⟩ +≡

```

octa omult ARGS((octa, octa));
    octa omult(y, z)
    octa y, z;
{
    register int i, j, k;
    tetra u[4], v[4], w[8];
    register tetra t;
    octa acc;
    ⟨Unpack the multiplier and multiplicand to u and v 10⟩;
    for (j = 0; j < 4; j++) w[j] = 0;
    for (j = 0; j < 4; j++)
        if (¬v[j]) w[j + 4] = 0;
        else {
            for (i = k = 0; i < 4; i++) {
                t = u[i] * v[j] + w[i + j] + k;
                w[i + j] = t & #ffff, k = t >> 16;
            }
            w[j + 4] = k;
        }
    ⟨Pack w into the outputs aux and acc 11⟩;
    return acc;
}

```

**9.** ⟨Global variables 4⟩ +≡

```

octa aux;      /* secondary output of subroutines with multiple outputs */
bool overflow; /* set by certain subroutines for signed arithmetic */

```

**10.** ⟨Unpack the multiplier and multiplicand to *u* and *v* 10⟩ ≡

```

u[3] = y.h >> 16, u[2] = y.h & #ffff, u[1] = y.l >> 16, u[0] = y.l & #ffff;
v[3] = z.h >> 16, v[2] = z.h & #ffff, v[1] = z.l >> 16, v[0] = z.l & #ffff;

```

This code is used in section 8.

**11.** ⟨Pack *w* into the outputs *aux* and *acc* 11⟩ ≡

```

aux.h = (w[7] << 16) + w[6], aux.l = (w[5] << 16) + w[4];
acc.h = (w[3] << 16) + w[2], acc.l = (w[1] << 16) + w[0];

```

This code is used in section 8.

**12.** Signed multiplication has the same lower half product as unsigned multiplication. The signed upper half product is obtained with at most two further subtractions, after which the result has overflowed if and only if the upper half is unequal to 64 copies of the sign bit in the lower half.

⟨Subroutines 5⟩ +≡

```

octa signed_omult ARGS((octa,octa));
    octa signed_omult(y,z)
    octa y,z;
{
    octa acc;
    acc = omult(y,z);
    if (y.h & sign_bit) aux = ominus(aux,z);
    if (z.h & sign_bit) aux = ominus(aux,y);
    overflow = (aux.h ≠ aux.l ∨ (aux.h ⊕ (aux.h ≫ 1) ⊕ (acc.h & sign_bit)));
    return acc;
}
```

**13. Division.** Long division of an unsigned 128-bit integer by an unsigned 64-bit integer is, of course, one of the most challenging routines needed for MMIX arithmetic. The following program, based on Algorithm 4.3.1D of *Seminumerical Algorithms*, computes octabytes  $q$  and  $r$  such that  $(2^{64}x+y) = qz+r$  and  $0 \leq r < z$ , given octabytes  $x$ ,  $y$ , and  $z$ , assuming that  $x < z$ . (If  $x \geq z$ , it simply sets  $q = x$  and  $r = y$ .) The quotient  $q$  is returned by the subroutine; the remainder  $r$  is stored in  $aux$ .

```

⟨Subroutines 5⟩ +=
  octa odiv ARGS((octa, octa, octa));
    octa odiv( $x, y, z$ )
    octa  $x, y, z$ ;
  {
    register int  $i, j, k, n, d$ ;
    tetra  $u[8], v[4], q[4], mask, qhat, rhat, vh, vmh$ ;
    register tetra  $t$ ;
    octa  $acc$ ;
    ⟨Check that  $x < z$ ; otherwise give trivial answer 14⟩;
    ⟨Unpack the dividend and divisor to  $u$  and  $v$  15⟩;
    ⟨Determine the number of significant places  $n$  in the divisor  $v$  16⟩;
    ⟨Normalize the divisor 17⟩;
    for ( $j = 3$ ;  $j \geq 0$ ;  $j--$ ) ⟨Determine the quotient digit  $q[j]$  20⟩;
    ⟨Unnormalize the remainder 18⟩;
    ⟨Pack  $q$  and  $u$  to  $acc$  and  $aux$  19⟩;
    return  $acc$ ;
  }

```

```

14. ⟨Check that  $x < z$ ; otherwise give trivial answer 14⟩ =
  if ( $x.h > z.h \vee (x.h \equiv z.h \wedge x.l \geq z.l)$ ) {
     $aux = y$ ; return  $x$ ;
  }

```

This code is used in section 13.

```

15. ⟨Unpack the dividend and divisor to  $u$  and  $v$  15⟩ =
   $u[7] = x.h \gg 16, u[6] = x.h \& \#ffff, u[5] = x.l \gg 16, u[4] = x.l \& \#ffff$ ;
   $u[3] = y.h \gg 16, u[2] = y.h \& \#ffff, u[1] = y.l \gg 16, u[0] = y.l \& \#ffff$ ;
   $v[3] = z.h \gg 16, v[2] = z.h \& \#ffff, v[1] = z.l \gg 16, v[0] = z.l \& \#ffff$ ;

```

This code is used in section 13.

```

16. ⟨Determine the number of significant places  $n$  in the divisor  $v$  16⟩ =
  for ( $n = 4$ ;  $v[n-1] \equiv 0$ ;  $n--$ ) ;

```

This code is used in section 13.

**17.** We shift  $u$  and  $v$  left by  $d$  places, where  $d$  is chosen to make  $2^{15} \leq v_{n-1} < 2^{16}$ .

```

⟨ Normalize the divisor 17 ⟩ ≡
    vh = v[n - 1];
    for (d = 0; vh < #8000; d++, vh <= 1) ;
    for (j = k = 0; j < n + 4; j++) {
        t = (u[j] << d) + k;
        u[j] = t & #ffff, k = t >> 16;
    }
    for (j = k = 0; j < n; j++) {
        t = (v[j] << d) + k;
        v[j] = t & #ffff, k = t >> 16;
    }
    vh = v[n - 1];
    vmh = (n > 1 ? v[n - 2] : 0);

```

This code is used in section 13.

```

18. ⟨ Unnormalize the remainder 18 ⟩ ≡
    mask = (1 << d) - 1;
    for (j = 3; j ≥ n; j--) u[j] = 0;
    for (k = 0; j ≥ 0; j--) {
        t = (k << 16) + u[j];
        u[j] = t >> d, k = t & mask;
    }

```

This code is used in section 13.

```

19. ⟨ Pack  $q$  and  $u$  to  $acc$  and  $aux$  19 ⟩ ≡
    acc.h = (q[3] << 16) + q[2], acc.l = (q[1] << 16) + q[0];
    aux.h = (u[3] << 16) + u[2], aux.l = (u[1] << 16) + u[0];

```

This code is used in section 13.

```

20. ⟨ Determine the quotient digit  $q[j]$  20 ⟩ ≡
    {
        ⟨ Find the trial quotient,  $\hat{q}$  21 ⟩;
        ⟨ Subtract  $b^j \hat{q} v$  from  $u$  22 ⟩;
        ⟨ If the result was negative, decrease  $\hat{q}$  by 1 23 ⟩;
        q[j] = qhat;
    }

```

This code is used in section 13.

```

21. ⟨ Find the trial quotient,  $\hat{q}$  21 ⟩ ≡
    t = (u[j + n] << 16) + u[j + n - 1];
    qhat = t / vh, rhat = t - vh * qhat;
    if (n > 1)
        while (qhat ≡ #10000 ∨ qhat * vmh > (rhat << 16) + u[j + n - 2]) {
            qhat--, rhat += vh;
            if (rhat ≥ #10000) break;
        }

```

This code is used in section 20.

**22.** After this step,  $u[j+n]$  will either equal  $k$  or  $k-1$ . The true value of  $u$  would be obtained by subtracting  $k$  from  $u[j+n]$ ; but we don't have to fuss over  $u[j+n]$ , because it won't be examined later.

```

⟨ Subtract  $b^j \hat{q} v$  from  $u$  22 ⟩ ≡
  for ( $i = k = 0; i < n; i++$ ) {
     $t = u[i+j] + \#ffff0000 - k - qhat * v[i];$ 
     $u[i+j] = t \& \#ffff, k = \#ffff - (t \gg 16);$ 
  }

```

This code is used in section 20.

**23.** The correction here occurs only rarely, but it can be necessary—for example, when dividing the number  $\#7fff800100000000$  by  $\#800080020005$ .

```

⟨ If the result was negative, decrease  $\hat{q}$  by 1 23 ⟩ ≡
  if ( $u[j+n] \neq (\text{tetra})k$ ) {
     $qhat--;$ 
    for ( $i = k = 0; i < n; i++$ ) {
       $t = u[i+j] + v[i] + k;$ 
       $u[i+j] = t \& \#ffff, k = t \gg 16;$ 
    }
  }

```

This code is used in section 20.

**24.** Signed division can be reduced to unsigned division in a tedious but straightforward manner. We assume that the divisor isn't zero.

```

⟨ Subroutines 5 ⟩ +=
  octa signed_odiv ARGS((octa, octa));
  octa signed_odiv(y, z)
  octa y, z;
  {
    octa yy, zz, q;
    register int sy, sz;
    if ( $y.h \& sign\_bit$ )  $sy = 2, yy = ominus(zero\_octa, y);$ 
    else  $sy = 0, yy = y;$ 
    if ( $z.h \& sign\_bit$ )  $sz = 1, zz = ominus(zero\_octa, z);$ 
    else  $sz = 0, zz = z;$ 
     $q = odiv(zero\_octa, yy, zz);$ 
     $overflow = false;$ 
    switch ( $sy + sz$ ) {
    case 2 + 1:  $aux = ominus(zero\_octa, aux);$ 
      if ( $q.h \equiv sign\_bit$ )  $overflow = true;$ 
    case 0 + 0: default: return q;
    case 2 + 0: if ( $aux.h \vee aux.l$ )  $aux = ominus(zz, aux);$ 
      goto negate_q;
    case 0 + 1: if ( $aux.h \vee aux.l$ )  $aux = ominus(aux, zz);$ 
      negate_q: if ( $aux.h \vee aux.l$ ) return ominus(neg_one, q);
      else return ominus(zero\_octa, q);
    }
  }

```



**25. Bit fiddling.** The bitwise operators of MMIX are fairly easy to implement directly, but three of them occur often enough to deserve packaging as subroutines.

⟨Subroutines 5⟩ +≡

```

octa oand ARGS((octa,octa));
    octa oand(y,z)      /* compute  $y \wedge z$  */
    octa y,z;
{ octa x;
  x.h = y.h & z.h; x.l = y.l & z.l;
  return x;
}

octa oandn ARGS((octa,octa));
    octa oandn(y,z)     /* compute  $y \wedge \bar{z}$  */
    octa y,z;
{ octa x;
  x.h = y.h & ~z.h; x.l = y.l & ~z.l;
  return x;
}

octa oxor ARGS((octa,octa));
    octa oxor(y,z)      /* compute  $y \oplus z$  */
    octa y,z;
{ octa x;
  x.h = y.h ⊕ z.h; x.l = y.l ⊕ z.l;
  return x;
}

```

**26.** Here’s a fun way to count the number of bits in a tetrabyte. [This classical trick is called the “Gillies–Miller method for sideways addition” in *The Preparation of Programs for an Electronic Digital Computer* by Wilkes, Wheeler, and Gill, second edition (Reading, Mass.: Addison–Wesley, 1957), 191–193. Some of the tricks used here were suggested by Balbir Singh, Peter Rossmanith, and Stefan Schwoon.]

⟨Subroutines 5⟩ +≡

```

int count_bits ARGS((tetra));
    int count_bits(x)
    tetra x;
{
  register int xx = x;
  xx = xx - ((xx >> 1) & #55555555);
  xx = (xx & #33333333) + ((xx >> 2) & #33333333);
  xx = (xx + (xx >> 4)) & #0f0f0f0f;
  xx = xx + (xx >> 8);
  return (xx + (xx >> 16)) & #ff;
}

```

**27.** To compute the nonnegative byte differences of two given tetrabytes, we can carry out the following 20-step branchless computation:

```

⟨Subroutines 5⟩ +=
  tetra byte_diff ARGS((tetra, tetra));
    tetra byte_diff(y, z)
    tetra y, z;
  {
    register tetra d = (y & #00ff00ff) + #01000100 - (z & #00ff00ff);
    register tetra m = d & #01000100;
    register tetra x = d & (m - (m >> 8));
    d = ((y >> 8) & #00ff00ff) + #01000100 - ((z >> 8) & #00ff00ff);
    m = d & #01000100;
    return x + ((d & (m - (m >> 8))) << 8);
  }

```

**28.** To compute the nonnegative wyde differences of two tetrabytes, another trick leads to a 15-step branchless computation. (Research problem: Can *count\_bits*, *byte\_diff*, or *wyde\_diff* be done with fewer operations?)

```

⟨Subroutines 5⟩ +=
  tetra wyde_diff ARGS((tetra, tetra));
    tetra wyde_diff(y, z)
    tetra y, z;
  {
    register tetra a = ((y >> 16) - (z >> 16)) & #10000;
    register tetra b = ((y & #ffff) - (z & #ffff)) & #10000;
    return y - (z ⊕ ((y ⊕ z) & (b - a - (b >> 16))));
  }

```

**29.** The last bitwise subroutine we need is the most interesting: It implements MMIX's MOR and MXOR operations.

```

⟨Subroutines 5⟩ +=
  octa bool_mult ARGS((octa, octa, bool));
    octa bool_mult(y, z, xor)
    octa y, z; /* the operands */
    bool xor; /* do we do xor instead of or? */
  {
    octa o, x;
    register tetra a, b, c;
    register int k;
    for (k = 0, o = y, x = zero_octa; o.h ∨ o.l; k++, o = shift_right(o, 8, 1))
      if (o.l & #ff) {
        a = ((z.h >> k) & #01010101) * #ff;
        b = ((z.l >> k) & #01010101) * #ff;
        c = (o.l & #ff) * #01010101;
        if (xor) x.h ⊕= a & c, x.l ⊕= b & c;
        else x.h |= a & c, x.l |= b & c;
      }
    return x;
  }

```

**30. Floating point packing and unpacking.** Standard IEEE floating binary numbers pack a sign, exponent, and fraction into a tetrabyte or octabyte. In this section we consider basic subroutines that convert between IEEE format and the separate unpacked components.

```
#define ROUND_OFF 1
#define ROUND_UP 2
#define ROUND_DOWN 3
#define ROUND_NEAR 4
⟨ Global variables 4 ⟩ +≡
    int cur_round; /* the current rounding mode */
```

**31.** The *fpack* routine takes an octabyte  $f$ , a raw exponent  $e$ , and a sign  $s$ , and packs them into the floating binary number that corresponds to  $\pm 2^{e-1076} f$ , using a given rounding mode. The value of  $f$  should satisfy  $2^{54} \leq f \leq 2^{55}$ .

Thus, for example, the floating binary number  $+1.0 = \#3ff0000000000000$  is obtained when  $f = 2^{54}$ ,  $e = \#3fe$ , and  $s = '+'$ . The raw exponent  $e$  is usually one less than the final exponent value; the leading bit of  $f$  is essentially added to the exponent. (This trick works nicely for subnormal numbers, when  $e < 0$ , or in cases where the value of  $f$  is rounded upwards to  $2^{55}$ .)

Exceptional events are noted by oring appropriate bits into the global variable *exceptions*. Special considerations apply to underflow, which is not fully specified by Section 7.4 of the IEEE standard: Implementations of the standard are free to choose between two definitions of “tininess” and two definitions of “accuracy loss.” MMIX determines tininess *after* rounding, hence a result with  $e < 0$  is not necessarily tiny; MMIX treats accuracy loss as equivalent to inexactness. Thus, a result underflows if and only if it is tiny and either (i) it is inexact or (ii) the underflow trap is enabled. The *fpack* routine sets U\_BIT in *exceptions* if and only if the result is tiny, X\_BIT if and only if the result is inexact.

```
#define X_BIT (1 << 8)    /* floating inexact */
#define Z_BIT (1 << 9)    /* floating division by zero */
#define U_BIT (1 << 10)   /* floating underflow */
#define O_BIT (1 << 11)   /* floating overflow */
#define I_BIT (1 << 12)   /* floating invalid operation */
#define W_BIT (1 << 13)   /* float-to-fix overflow */
#define V_BIT (1 << 14)   /* integer overflow */
#define D_BIT (1 << 15)   /* integer divide check */
#define E_BIT (1 << 18)   /* external (dynamic) trap bit */

⟨Subroutines 5⟩ +=
  octa fpack ARGS((octa, int, int, int));
    octa fpack(f, e, s, r)
    octa f;    /* the normalized fraction part */
    int e;     /* the raw exponent */
    int s;     /* the sign */
    int r;     /* the rounding mode */
    {
      octa o;
      if (e > #7fd) e = #7ff, o = zero_octa;
      else {
        if (e < 0) {
          if (e < -54) o.h = 0, o.l = 1;
          else { octa oo;
            o = shift_right(f, -e, 1);
            oo = shift_left(o, -e);
            if (oo.l ≠ f.l ∨ oo.h ≠ f.h) o.l |= 1;    /* sticky bit */
          }
          e = 0;
        } else o = f;
      }
    }
  ⟨Round and return the result 33⟩;
}
```

**32.** ⟨Global variables 4⟩ +=

```
int exceptions;    /* bits possibly destined for rA */
```

**33.** Everything falls together so nicely here, it's almost too good to be true! The conditional expression in the case for `ROUND_NEAR` rounds towards an even number in case of a tie.

```

⟨Round and return the result 33⟩ ≡
  if (o.l & 3) exceptions |= X_BIT;
  switch (r) {
  case ROUND_DOWN: if (s == '-') o = incr(o, 3); break;
  case ROUND_UP:   if (s != '-') o = incr(o, 3);
  case ROUND_OFF:  break;
  case ROUND_NEAR: o = incr(o, o.l & 4 ? 2 : 1); break;
  }
  o = shift_right(o, 2, 1);
  o.h += e << 20;
  if (o.h ≥ #7ff00000) exceptions |= O_BIT + X_BIT; /* overflow */
  else if (o.h < #100000) exceptions |= U_BIT; /* tininess */
  if (s == '-') o.h |= sign_bit;
  return o;

```

This code is used in section 31.

**34.** Similarly, *sfpack* packs a short float, from inputs having the same conventions as *fpack*.

```

⟨Subroutines 5⟩ +=
  tetra sfpack ARGS((octa, int, int, int));
  tetra sfpack(f, e, s, r)
    octa f; /* the fraction part */
    int e; /* the raw exponent */
    int s; /* the sign */
    int r; /* the rounding mode */
  {
    register tetra o;
    if (e > #47d) e = #47f, o = 0;
    else {
      o = shift_left(f, 3).h;
      if (f.l & #1fffffff) o |= 1;
      if (e < #380) {
        if (e < #380 - 25) o = 1;
        else { register tetra o0, oo;
          o0 = o;
          o = o >> (#380 - e);
          oo = o << (#380 - e);
          if (oo ≠ o0) o |= 1; /* sticky bit */
        }
        e = #380;
      }
    }
  }
  ⟨Round and return the short result 35⟩;
}

```

**35.**  $\langle$  Round and return the short result [35](#)  $\rangle \equiv$

```

if ( $o \& 3$ ) exceptions |= X_BIT;
switch ( $r$ ) {
  case ROUND_DOWN: if ( $s \equiv \text{'-'}\text{'}$ )  $o += 3$ ; break;
  case ROUND_UP: if ( $s \neq \text{'-'}\text{'}$ )  $o += 3$ ;
  case ROUND_OFF: break;
  case ROUND_NEAR:  $o += (o \& 4 ? 2 : 1)$ ; break;
}
 $o = o \gg 2$ ;
 $o += (e - \#380) \ll 23$ ;
if ( $o \geq \#7f800000$ ) exceptions |= O_BIT + X_BIT;    /* overflow */
else if ( $o < \#800000$ ) exceptions |= U_BIT;          /* tininess */
if ( $s \equiv \text{'-'}\text{'}$ )  $o$  |= sign_bit;
return  $o$ ;

```

This code is used in section [34](#).

**36.** The *funpack* routine is, roughly speaking, the opposite of *fpack*. It takes a given floating point number  $x$  and separates out its fraction part  $f$ , exponent  $e$ , and sign  $s$ . It clears *exceptions* to zero. It returns the type of value found: *zro*, *num*, *inf*, or *nan*. When it returns *num*, it will have set  $f$ ,  $e$ , and  $s$  to the values from which *fpack* would produce the original number  $x$  without exceptions.

```

#define zero_exponent (-1000)    /* zero is assumed to have this exponent */

```

$\langle$  Other type definitions [36](#)  $\rangle \equiv$

```

typedef enum {
    zro, num, inf, nan
} ftype;

```

See also section [59](#).

This code is used in section [1](#).

37. ⟨Subroutines 5⟩ +≡

```

ftype funpack ARGs((octa, octa *, int *, char *));
ftype funpack(x, f, e, s)
  octa x;      /* the given floating point value */
  octa *f;     /* address where the fraction part should be stored */
  int *e;      /* address where the exponent part should be stored */
  char *s;     /* address where the sign should be stored */
{
  register int ee;
  exceptions = 0;
  *s = (x.h & sign_bit ? '-' : '+');
  *f = shift_left(x, 2);
  f ← f & #3fffff;
  ee = (x.h >> 20) & #7ff;
  if (ee) {
    *e = ee - 1;
    f ← f | #400000;
    return (ee < #7ff ? num : f ← f & #400000 ∧ ¬f ← f ? inf : nan);
  }
  if (¬x.l ∧ ¬f ← f) {
    *e = zero_exponent; return zro;
  }
  do { ee --; *f = shift_left(*f, 1); } while (¬(f ← f & #400000));
  *e = ee; return num;
}

```

38. ⟨Subroutines 5⟩ +≡

```

ftype sfunpack ARGs((tetra, octa *, int *, char *));
ftype sfunpack(x, f, e, s)
  tetra x;     /* the given floating point value */
  octa *f;     /* address where the fraction part should be stored */
  int *e;      /* address where the exponent part should be stored */
  char *s;     /* address where the sign should be stored */
{
  register int ee;
  exceptions = 0;
  *s = (x & sign_bit ? '-' : '+');
  f ← f = (x >> 1) & #3fffff, f ← f << 31;
  ee = (x >> 23) & #ff;
  if (ee) {
    *e = ee + #380 - 1;
    f ← f | #400000;
    return (ee < #ff ? num : (x & #7fffffff) ≡ #7f800000 ? inf : nan);
  }
  if (¬(x & #7fffffff)) {
    *e = zero_exponent; return zro;
  }
  do { ee --; *f = shift_left(*f, 1); } while (¬(f ← f & #400000));
  *e = ee + #380; return num;
}

```

**39.** Since MMIX downplays 32-bit operations, it uses *sfpack* and *sfunpack* only when loading and storing short floats, or when converting from fixed point to floating point.

⟨Subroutines 5⟩ +≡

```

octa load_sf ARGS((tetra));
    octa load_sf(z)
    tetra z; /* 32 bits to be loaded into a 64-bit register */
{
    octa f,x; int e; char s; fctype t;
    t = sfunpack(z,&f,&e,&s);
    switch (t) {
    case zro: x = zero_octa; break;
    case num: return fpack(f,e,s,ROUND_OFF);
    case inf: x = inf_octa; break;
    case nan: x = shift_right(f,2,1); x.h |= #7ff00000; break;
    }
    if (s ≡ '-') x.h |= sign_bit;
    return x;
}

```

**40.** ⟨Subroutines 5⟩ +≡

```

tetra store_sf ARGS((octa));
    tetra store_sf(x)
    octa x; /* 64 bits to be loaded into a 32-bit word */
{
    octa f; tetra z; int e; char s; fctype t;
    t = funpack(x,&f,&e,&s);
    switch (t) {
    case zro: z = 0; break;
    case num: return sfpack(f,e,s,cur_round);
    case inf: z = #7f800000; break;
    case nan: if (¬(f.h & #200000)) {
        f.h |= #200000; exceptions |= I_BIT; /* NaN was signaling */
    }
    z = #7f800000 | (f.h << 1) | (f.l >> 31); break;
    }
    if (s ≡ '-') z |= sign_bit;
    return z;
}

```



**41. Floating multiplication and division.** The hardest fixed point operations were multiplication and division; but these two operations are the *easiest* to implement in floating point arithmetic, once their fixed point counterparts are available.

⟨Subroutines 5⟩ +=

```

octa fmult ARGS((octa, octa));
    octa fmult(y, z)
    octa y, z;
{
    ftype yt, zt;
    int ye, ze;
    char ys, zs;
    octa x, xf, yf, zf;
    register int xe;
    register char xs;
    yt = funpack(y, &yf, &ye, &ys);
    zt = funpack(z, &zf, &ze, &zs);
    xs = ys + zs - '+'; /* will be '-' when the result is negative */
    switch (4 * yt + zt) {
    ⟨The usual NaN cases 42⟩;
    default: case 4 * zro + zro: case 4 * zro + num: case 4 * num + zro: x = zero_octa; break;
    case 4 * num + inf: case 4 * inf + num: case 4 * inf + inf: x = inf_octa; break;
    case 4 * zro + inf: case 4 * inf + zro: x = standard_NaN;
        exceptions |= I_BIT; break;
    case 4 * num + num: ⟨Multiply nonzero numbers and return 43⟩;
    }
    if (xs == '-') x.h |= sign_bit;
    return x;
}

```

**42.** ⟨The usual NaN cases 42⟩ =

```

case 4 * nan + nan: if (¬(y.h & #80000)) exceptions |= I_BIT; /* y is signaling */
case 4 * zro + nan: case 4 * num + nan: case 4 * inf + nan:
    if (¬(z.h & #80000)) exceptions |= I_BIT, z.h |= #80000;
    return z;
case 4 * nan + zro: case 4 * nan + num: case 4 * nan + inf:
    if (¬(y.h & #80000)) exceptions |= I_BIT, y.h |= #80000;
    return y;

```

This code is used in sections 41, 44, 46, and 93.

**43.** ⟨Multiply nonzero numbers and **return** 43⟩ =

```

xe = ye + ze - #3fd; /* the raw exponent */
x = omult(yf, shift_left(zf, 9));
if (aux.h ≥ #400000) xf = aux;
else xf = shift_left(aux, 1), xe--;
if (x.h ∨ x.l) xf.l |= 1; /* adjust the sticky bit */
return fpack(xf, xe, xs, cur_round);

```

This code is used in section 41.

44.  $\langle$  Subroutines 5  $\rangle + \equiv$

```

octa fdivide ARGS((octa, octa));
    octa fdivide(y, z)
    octa y, z;
{
    ftype yt, zt;
    int ye, ze;
    char ys, zs;
    octa x, xf, yf, zf;
    register int xe;
    register char xs;

    yt = funpack(y, &yf, &ye, &ys);
    zt = funpack(z, &zf, &ze, &zs);
    xs = ys + zs - '+'; /* will be '-' when the result is negative */
    switch (4 * yt + zt) {
     $\langle$  The usual NaN cases 42  $\rangle$ ;
    case 4 * zro + inf: case 4 * zro + num: case 4 * num + inf: x = zero_octa; break;
    case 4 * num + zro: exceptions |= Z_BIT;
    case 4 * inf + num: case 4 * inf + zro: x = inf_octa; break;
    default: case 4 * zro + zro: case 4 * inf + inf: x = standard_NaN;
        exceptions |= I_BIT; break;
    case 4 * num + num:  $\langle$  Divide nonzero numbers and return 45  $\rangle$ ;
    }
    if (xs  $\equiv$  '-') x.h |= sign_bit;
    return x;
}

```

45.  $\langle$  Divide nonzero numbers and **return** 45  $\rangle \equiv$

```

xe = ye - ze + #3fd; /* the raw exponent */
xf = odiv(yf, zero_octa, shift_left(zf, 9));
if (xf.h  $\geq$  #800000) {
    aux.l |= xf.l & 1;
    xf = shift_right(xf, 1, 1);
    xe++;
}
if (aux.h  $\vee$  aux.l) xf.l |= 1; /* adjust the sticky bit */
return fpack(xf, xe, xs, cur_round);

```

This code is used in section 44.

**46. Floating addition and subtraction.** Now for the bread-and-butter operation, the sum of two floating point numbers. It is not terribly difficult, but many cases need to be handled carefully.

⟨Subroutines 5⟩ +≡

```

octa fplus ARGS((octa,octa));
    octa fplus(y,z)
    octa y,z;
{
    ftype yt,zt;
    int ye,ze;
    char ys,zs;
    octa x,xf,yf,zf;
    register int xe,d;
    register char xs;
    yt = funpack(y, &yf, &ye, &ys);
    zt = funpack(z, &zf, &ze, &zs);
    switch (4 * yt + zt) {
    ⟨The usual NaN cases 42⟩;
    case 4 * zro + num: return fpack(zf, ze, zs, ROUND_OFF); break;    /* may underflow */
    case 4 * num + zro: return fpack(yf, ye, ys, ROUND_OFF); break;    /* may underflow */
    case 4 * inf + inf: if (ys ≠ zs) {
        exceptions |= I_BIT; x = standard_NaN; xs = zs; break;
    }
    case 4 * num + inf: case 4 * zro + inf: x = inf_octa; xs = zs; break;
    case 4 * inf + num: case 4 * inf + zro: x = inf_octa; xs = ys; break;
    case 4 * num + num: if (y.h ≠ (z.h ⊕ #80000000) ∨ y.l ≠ z.l) ⟨Add nonzero numbers and return 47⟩;
    default: case 4 * zro + zro: x = zero_octa;
        xs = (ys ≡ zs ? ys : cur_round ≡ ROUND_DOWN ? '-' : '+'); break;
    }
    if (xs ≡ '-') x.h |= sign_bit;
    return x;
}

```

**47.** ⟨Add nonzero numbers and **return** 47⟩ ≡

```

{ octa o,oo;
  if (ye < ze ∨ (ye ≡ ze ∧ (yf.h < zf.h ∨ (yf.h ≡ zf.h ∧ yf.l < zf.l)))) ⟨Exchange y with z 48⟩;
  d = ye − ze;
  xs = ys, xe = ye;
  if (d) ⟨Adjust for difference in exponents 49⟩;
  if (ys ≡ zs) {
    xf = oplus(yf, zf);
    if (xf.h ≥ #800000) xe ++, d = xf.l & 1, xf = shift_right(xf, 1, 1), xf.l |= d;
  } else {
    xf = ominus(yf, zf);
    if (xf.h ≥ #800000) xe ++, d = xf.l & 1, xf = shift_right(xf, 1, 1), xf.l |= d;
    else while (xf.h < #400000) xe −, xf = shift_left(xf, 1);
  }
  return fpack(xf, xe, xs, cur_round);
}

```

This code is used in section 46.

48.  $\langle \text{Exchange } y \text{ with } z \text{ 48} \rangle \equiv$   
 $\{$   
 $o = yf, yf = zf, zf = o;$   
 $d = ye, ye = ze, ze = d;$   
 $d = ys, ys = zs, zs = d;$   
 $\}$

This code is used in sections 47 and 51.

49. Proper rounding requires two bits to the right of the fraction delivered to *fpack*. The first is the true next bit of the result; the other is a “sticky” bit, which is nonzero if any further bits of the true result are nonzero. Sticky rounding to an integer takes  $x$  into the number  $\lfloor x/2 \rfloor + \lceil x/2 \rceil$ .

Some subtleties need to be observed here, in order to prevent the sticky bit from being shifted left. If we did not shift  $yf$  left 1 before shifting  $zf$  to the right, an incorrect answer would be obtained in certain cases—for example, if  $yf = 2^{54}$ ,  $zf = 2^{54} + 2^{53} - 4$ ,  $d = 52$ .

$\langle \text{Adjust for difference in exponents 49} \rangle \equiv$   
 $\{$   
 $\text{if } (d \leq 2) \text{ } zf = \text{shift\_right}(zf, d, 1); \quad /* \text{ exact result } */$   
 $\text{else if } (d > 54) \text{ } zf.h = 0, zf.l = 1; \quad /* \text{ tricky but OK } */$   
 $\text{else } \{$   
 $\text{if } (ys \neq zs) \text{ } d--, xe--, yf = \text{shift\_left}(yf, 1);$   
 $o = zf;$   
 $zf = \text{shift\_right}(o, d, 1);$   
 $oo = \text{shift\_left}(zf, d);$   
 $\text{if } (oo.l \neq o.l \vee oo.h \neq o.h) \text{ } zf.l |= 1;$   
 $\}$   
 $\}$

This code is used in section 47.

**50.** The comparison of floating point numbers with respect to  $\epsilon$  shares some of the characteristics of floating point addition/subtraction. In some ways it is simpler, and in other ways it is more difficult; we might as well deal with it now.

Subroutine *fepscomp*( $y, z, e, s$ ) returns 2 if  $y$ ,  $z$ , or  $e$  is a NaN or  $e$  is negative. It returns 1 if  $s = 0$  and  $y \approx z$  ( $e$ ) or if  $s \neq 0$  and  $y \sim z$  ( $e$ ), as defined in Section 4.2.2 of *Seminumerical Algorithms*; otherwise it returns 0.

⟨Subroutines 5⟩  $\equiv$

```

int fepscomp ARGs((octa, octa, octa, int));
    int fepscomp( $y, z, e, s$ )
        octa  $y, z, e$ ;    /* the operands */
        int  $s$ ;          /* test similarity? */
    {
        octa  $yf, zf, ef, o, oo$ ;
        int  $ye, ze, ee$ ;
        char  $ys, zs, es$ ;
        register int  $yt, zt, et, d$ ;
         $et = \text{funpack}(e, \&ef, \&ee, \&es)$ ;
        if ( $es \equiv \text{'-'}$ ) return 2;
        switch ( $et$ ) {
        case nan: return 2;
        case inf:  $ee = 10000$ ;
        case num: case zro: break;
        }
         $yt = \text{funpack}(y, \&yf, \&ye, \&ys)$ ;
         $zt = \text{funpack}(z, \&zf, \&ze, \&zs)$ ;
        switch ( $4 * yt + zt$ ) {
        case  $4 * \text{nan} + \text{nan}$ : case  $4 * \text{nan} + \text{inf}$ : case  $4 * \text{nan} + \text{num}$ : case  $4 * \text{nan} + \text{zro}$ : case  $4 * \text{inf} + \text{nan}$ :
            case  $4 * \text{num} + \text{nan}$ : case  $4 * \text{zro} + \text{nan}$ : return 2;
        case  $4 * \text{inf} + \text{inf}$ : return ( $ys \equiv zs \vee ee \geq 1023$ );
        case  $4 * \text{inf} + \text{num}$ : case  $4 * \text{inf} + \text{zro}$ : case  $4 * \text{num} + \text{inf}$ : case  $4 * \text{zro} + \text{inf}$ : return ( $s \wedge ee \geq 1022$ );
        case  $4 * \text{zro} + \text{zro}$ : return 1;
        case  $4 * \text{zro} + \text{num}$ : case  $4 * \text{num} + \text{zro}$ : if ( $\neg s$ ) return 0;
        case  $4 * \text{num} + \text{num}$ : break;
        }
        ⟨Compare two numbers with respect to epsilon and return 51⟩;
    }

```

**51.** The relation  $y \approx z (\epsilon)$  reduces to  $y \sim z (\epsilon/2^d)$ , if  $d$  is the difference between the larger and smaller exponents of  $y$  and  $z$ .

```

⟨ Compare two numbers with respect to epsilon and return 51 ⟩ ≡
  ⟨ Unsubnormalize  $y$  and  $z$ , if they are subnormal 52 ⟩;
  if ( $ye < ze \vee (ye \equiv ze \wedge (yf.h < zf.h \vee (yf.h \equiv zf.h \wedge yf.l < zf.l)))$ ) ⟨ Exchange  $y$  with  $z$  48 ⟩;
  if ( $ze \equiv zero\_exponent$ )  $ze = ye$ ;
   $d = ye - ze$ ;
  if ( $\neg s$ )  $ee \leftarrow d$ ;
  if ( $ee \geq 1023$ ) return 1; /* if  $\epsilon \geq 2$ ,  $z \in N_\epsilon(y)$  */
  ⟨ Compute the difference of fraction parts, o 53 ⟩;
  if ( $\neg o.h \wedge \neg o.l$ ) return 1;
  if ( $ee < 968$ ) return 0; /* if  $y \neq z$  and  $\epsilon < 2^{-54}$ ,  $y \not\sim z$  */
  if ( $ee \geq 1021$ )  $ef = shift\_left(ef, ee - 1021)$ ;
  else  $ef = shift\_right(ef, 1021 - ee, 1)$ ;
  return  $o.h < ef.h \vee (o.h \equiv ef.h \wedge o.l \leq ef.l)$ ;

```

This code is used in section 50.

**52.** ⟨ Unsubnormalize  $y$  and  $z$ , if they are subnormal 52 ⟩ ≡  
**if** ( $ye < 0 \wedge yt \neq zro$ )  $yf = shift\_left(y, 2)$ ,  $ye = 0$ ;  
**if** ( $ze < 0 \wedge zt \neq zro$ )  $zf = shift\_left(z, 2)$ ,  $ze = 0$ ;

This code is used in section 51.

**53.** At this point  $y \sim z$  if and only if

$$yf + (-1)^{[ys=zs]}zf/2^d \leq 2^{ee-1021}ef = 2^{55}\epsilon.$$

We need to evaluate this relation without overstepping the bounds of our simulated 64-bit registers.

When  $d > 2$ , the difference of fraction parts might not fit exactly in an octabyte; in that case the numbers are not similar unless  $\epsilon > 3/8$ , and we replace the difference by the ceiling of the true result. When  $\epsilon < 1/8$ , our program essentially replaces  $2^{55}\epsilon$  by  $\lceil 2^{55}\epsilon \rceil$ . These truncations are not needed simultaneously. Therefore the logic is justified by the facts that, if  $n$  is an integer, we have  $x \leq n$  if and only if  $\lceil x \rceil \leq n$ ;  $n \leq x$  if and only if  $n \leq \lfloor x \rfloor$ . (Notice that the concept of “sticky bit” is *not* appropriate here.)

```

⟨ Compute the difference of fraction parts, o 53 ⟩ ≡
  if ( $d > 54$ )  $o = zero\_octa$ ,  $oo = zf$ ;
  else  $o = shift\_right(zf, d, 1)$ ,  $oo = shift\_left(o, d)$ ;
  if ( $oo.h \neq zf.h \vee oo.l \neq zf.l$ ) { /* truncated result, hence  $d > 2$  */
    if ( $ee < 1020$ ) return 0; /* difference is too large for similarity */
     $o = incr(o, ys \equiv zs ? 0 : 1)$ ; /* adjust for ceiling */
  }
   $o = (ys \equiv zs ? ominus(yf, o) : oplus(yf, o))$ ;

```

This code is used in section 51.

**54. Floating point output conversion.** The *print\_float* routine converts an octabyte to a floating decimal representation that will be input as precisely the same value.

⟨ Subroutines 5 ⟩ +≡

```

static void bignum_times_ten ARGS((bignum *));
static void bignum_dec ARGS((bignum *, bignum *, tetra));
static int bignum_compare ARGS((bignum *, bignum *));
void print_float ARGS((octa));
    void print_float(x)
    octa x;
{
    ⟨ Local variables for print_float 56 ⟩;
    if (x.h & sign_bit) printf("-");
    ⟨ Extract the exponent e and determine the fraction interval [f .. g] or (f .. g) 55 ⟩;
    ⟨ Store f and g as multiprecise integers 63 ⟩;
    ⟨ Compute the significant digits s and decimal exponent e 64 ⟩;
    ⟨ Print the significant digits with proper context 67 ⟩;
}
```

**55.** One way to visualize the problem being solved here is to consider the vastly simpler case in which there are only 2-bit exponents and 2-bit fractions. Then the sixteen possible 4-bit combinations have the following interpretations:

0000	[0 .. 0.125]
0001	(0.125 .. 0.375)
0010	[0.375 .. 0.625]
0011	(0.625 .. 0.875)
0100	[0.875 .. 1.125]
0101	(1.125 .. 1.375)
0110	[1.375 .. 1.625]
0111	(1.625 .. 1.875)
1000	[1.875 .. 2.25]
1001	(2.25 .. 2.75)
1010	[2.75 .. 3.25]
1011	(3.25 .. 3.75)
1100	[3.75 .. $\infty$ ]
1101	NaN(0 .. 0.375)
1110	NaN[0.375 .. 0.625]
1111	NaN(0.625 .. 1)

Notice that the interval is closed,  $[f..g]$ , when the fraction part is even; it is open,  $(f..g)$ , when the fraction part is odd. The printed outputs for these sixteen values, if we actually were dealing with such short exponents and fractions, would be 0., .2, .5, .7, 1., 1.2, 1.5, 1.7, 2., 2.5, 3., 3.5, Inf, NaN.2, NaN, NaN.8, respectively.

⟨Extract the exponent  $e$  and determine the fraction interval  $[f..g]$  or  $(f..g)$  55⟩  $\equiv$

```

f = shift_left(x, 1);
e = f.h >> 21;
f.h &= #1fffff;
if (¬f.h ∧ ¬f.l) ⟨Handle the special case when the fraction part is zero 57⟩
else {
    g = incr(f, 1);
    f = incr(f, -1);
    if (¬e) e = 1; /* subnormal */
    else if (e ≡ #7ff) {
        printf("NaN");
        if (g.h ≡ #100000 ∧ g.l ≡ 1) return; /* the "standard" NaN */
        e = #3ff; /* extreme NaNs come out OK even without adjusting f or g */
    } else f.h |= #200000, g.h |= #200000;
}

```

This code is used in section 54.

**56.** ⟨Local variables for *print\_float* 56⟩  $\equiv$

```

octa f, g; /* lower and upper bounds on the fraction part */
register int e; /* exponent part */
register int j, k; /* all purpose indices */

```

See also section 66.

This code is used in section 54.



**57.** The transition points between exponents correspond to powers of 2. At such points the interval extends only half as far to the left of that power of 2 as it does to the right. For example, in the 4-bit minifloat numbers considered above, case 1000 corresponds to the interval  $[1.875 \dots 2.25]$ .

⟨Handle the special case when the fraction part is zero 57⟩ =

```
{
  if (¬e) {
    printf("0."); return;
  }
  if (e ≡ #7ff) {
    printf("Inf"); return;
  }
  e--;
  f.h = #3fffff, f.l = #ffffff;
  g.h = #400000, g.l = 2;
}
```

This code is used in section 55.

**58.** We want to find the “simplest” value in the interval corresponding to the given number, in the sense that it has fewest significant digits when expressed in decimal notation. Thus, for example, if the floating point number can be described by a relatively short string such as ‘.1’ or ‘37e100’, we want to discover that representation.

The basic idea is to generate the decimal representations of the two endpoints of the interval, outputting the leading digits where both endpoints agree, then making a final decision at the first place where they disagree.

The “simplest” value is not always unique. For example, in the case of 4-bit minifloat numbers we could represent the bit pattern 0001 as either .2 or .3, and we could represent 1001 in five equally short ways: 2.3 or 2.4 or 2.5 or 2.6 or 2.7. The algorithm below tries to choose the middle possibility in such cases.

[A solution to the analogous problem for fixed-point representations, without the additional complication of round-to-even, was used by the author in the program for T<sub>E</sub>X; see *Beauty is Our Business* (Springer, 1990), 233–242.]

Suppose we are given two fractions  $f$  and  $g$ , where  $0 \leq f < g < 1$ , and we want to compute the shortest decimal in the closed interval  $[f..g]$ . If  $f = 0$ , we are done. Otherwise let  $10f = d + f'$  and  $10g = e + g'$ , where  $0 \leq f' < 1$  and  $0 \leq g' < 1$ . If  $d < e$ , we can terminate by outputting any of the digits  $d + 1, \dots, e$ ; otherwise we output the common digit  $d = e$ , and repeat the process on the fractions  $0 \leq f' < g' < 1$ . A similar procedure works with respect to the open interval  $(f..g)$ .

**59.** The program below carries out the stated algorithm by using multiprecision arithmetic on 77-place integers with 28 bits each. This choice facilitates multiplication by 10, and allows us to deal with the whole range of floating binary numbers using fixed point arithmetic. We keep track of the leading and trailing digit positions so that trivial operations on zeros are avoided.

If  $f$  points to a **bignum**, its radix- $2^{28}$  digits are  $f\text{-}dat[0]$  through  $f\text{-}dat[76]$ , from most significant to least significant. We assume that all digit positions are zero unless they lie in the subarray between indices  $f\text{-}a$  and  $f\text{-}b$ , inclusive. Furthermore, both  $f\text{-}dat[f\text{-}a]$  and  $f\text{-}dat[f\text{-}b]$  are nonzero, unless  $f\text{-}a = f\text{-}b = \text{bignum\_prec} - 1$ .

The **bignum** data type can be used with any radix less than  $2^{32}$ ; we will use it later with radix  $10^9$ . The *dat* array is made large enough to accommodate both applications.

```
#define bignum_prec 157 /* would be 77 if we cared only about print_float */
<Other type definitions 36> +≡
typedef struct {
    int a; /* index of the most significant digit */
    int b; /* index of the least significant digit; must be ≥ a */
    tetra dat[bignum_prec]; /* the digits; undefined except between a and b */
} bignum;
```

**60.** Here, for example, is how we go from  $f$  to  $10f$ , assuming that overflow will not occur and that the radix is  $2^{28}$ :

```
<Subroutines 5> +≡
static void bignum_times_ten(f)
    bignum *f;
{
    register tetra *p,*q;
    register tetra x, carry;
    for (p = &f-dat[f-b], q = &f-dat[f-a], carry = 0; p ≥ q; p--) {
        x = *p * 10 + carry;
        *p = x & #ffffff;
        carry = x >> 28;
    }
    *p = carry;
    if (carry) f-a--;
    if (f-dat[f-b] ≡ 0 ∧ f-b > f-a) f-b--;
}
```

**61.** And here is how we test whether  $f < g$ ,  $f = g$ , or  $f > g$ , using any radix whatever:

```
<Subroutines 5> +≡
static int bignum_compare(f,g)
    bignum *f,*g;
{
    register tetra *p,*pp,*q,*qq;
    if (f-a ≠ g-a) return f-a > g-a ? -1 : 1;
    pp = &f-dat[f-b], qq = &g-dat[g-b];
    for (p = &f-dat[f-a], q = &g-dat[g-a]; p ≤ pp; p++, q++) {
        if (*p ≠ *q) return *p < *q ? -1 : 1;
        if (q ≡ qq) return p < pp;
    }
    return -1;
}
```

**62.** The following subroutine subtracts  $g$  from  $f$ , assuming that  $f \geq g > 0$  and using a given radix.

⟨Subroutines 5⟩ +≡

```
static void bignum_dec(f, g, r)
    bignum *f, *g;
    tetra r;      /* the radix */
{
    register tetra *p, *q, *qq;
    register int x, borrow;
    while (g->b > f->b) f->dat[++f->b] = 0;
    qq = &g->dat[g->a];
    for (p = &f->dat[g->b], q = &g->dat[g->b], borrow = 0; q ≥ qq; p--, q--) {
        x = *p - *q - borrow;
        if (x ≥ 0) borrow = 0, *p = x;
        else borrow = 1, *p = x + r;
    }
    for (; borrow; p--)
        if (*p) borrow = 0, *p = *p - 1;
        else *p = r - 1;
    while (f->dat[f->a] ≡ 0) {
        if (f->a ≡ f->b) { /* the result is zero */
            f->a = f->b = bignum_prec - 1, f->dat[bignum_prec - 1] = 0;
            return;
        }
        f->a++;
    }
    while (f->dat[f->b] ≡ 0) f->b--;
}
```

**63.** Armed with these subroutines, we are ready to solve the problem. The first task is to put the numbers into **bignum** form. If the exponent is  $e$ , the number destined for digit  $dat[k]$  will consist of the rightmost 28 bits of the given fraction after it has been shifted right  $c - e - 28k$  bits, for some constant  $c$ . We choose  $c$  so that, when  $e$  has its maximum value `#7fff`, the leading digit will go into position  $dat[1]$ , and so that when the number to be printed is exactly 1 the integer part of  $g$  will also be exactly 1.

```
#define magic_offset 2112 /* the constant c that makes it work */
#define origin 37 /* the radix point follows dat[37] */
⟨Store f and g as multiprecise integers 63⟩ ≡
    k = (magic_offset - e)/28;
    ff.dat[k - 1] = shift_right(f, magic_offset + 28 - e - 28 * k, 1).l & #ffffff;
    gg.dat[k - 1] = shift_right(g, magic_offset + 28 - e - 28 * k, 1).l & #ffffff;
    ff.dat[k] = shift_right(f, magic_offset - e - 28 * k, 1).l & #ffffff;
    gg.dat[k] = shift_right(g, magic_offset - e - 28 * k, 1).l & #ffffff;
    ff.dat[k + 1] = shift_left(f, e + 28 * k - (magic_offset - 28)).l & #ffffff;
    gg.dat[k + 1] = shift_left(g, e + 28 * k - (magic_offset - 28)).l & #ffffff;
    ff.a = (ff.dat[k - 1] ? k - 1 : k);
    ff.b = (ff.dat[k + 1] ? k + 1 : k);
    gg.a = (gg.dat[k - 1] ? k - 1 : k);
    gg.b = (gg.dat[k + 1] ? k + 1 : k);
```

This code is used in section 54.

**64.** If  $e$  is sufficiently small, the fractions  $f$  and  $g$  will be less than 1, and we can use the stated algorithm directly. Of course, if  $e$  is extremely small, a lot of leading zeros need to be lopped off; in the worst case, we may have to multiply  $f$  and  $g$  by 10 more than 300 times. But hey, we don't need to do that extremely often, and computers are pretty fast nowadays.

In the small-exponent case, the computation always terminates before  $f$  becomes zero, because the interval endpoints are fractions with denominator  $2^t$  for some  $t > 50$ .

The invariant relations  $ff.dat[ff.a] \neq 0$  and  $gg.dat[gg.a] \neq 0$  are not maintained by the computation here, when  $ff.a = origin$  or  $gg.a = origin$ . But no harm is done, because *bignum\_compare* is not used.

```

⟨ Compute the significant digits  $s$  and decimal exponent  $e$  64 ⟩ ≡
  if ( $e > \#401$ ) ⟨ Compute the significant digits in the large-exponent case 65 ⟩
  else { /* if  $e \leq \#401$  we have  $gg.a \geq origin$  and  $gg.dat[origin] \leq 8$  */
    if ( $ff.a > origin$ )  $ff.dat[origin] = 0$ ;
    for ( $e = 1, p = s$ ;  $gg.a > origin \vee ff.dat[origin] \equiv gg.dat[origin]$ ; ) {
      if ( $gg.a > origin$ )  $e--$ ;
      else  $*p++ = ff.dat[origin] + '0'$ ,  $ff.dat[origin] = 0$ ,  $gg.dat[origin] = 0$ ;
      bignum_times_ten(&ff);
      bignum_times_ten(&gg);
    }
     $*p++ = ((ff.dat[origin] + 1 + gg.dat[origin]) \gg 1) + '0'$ ; /* the middle digit */
  }
   $*p = '\0'$ ; /* terminate the string  $s$  */

```

This code is used in section 54.

An interesting case arises when the number to be converted is `#44ada56a4b0835bf`, since the interval turns out to be

(69999999999999991611392 .. 70000000000000000000000).

If this were a closed interval, we could simply give the answer `7e22`; but the number `7e22` actually corresponds to `#44ada56a4b0835c0` because of the round-to-even rule. Therefore the correct answer is, say, `6.9999999999999995e22`. This example shows that we need a slightly different strategy in the case of open intervals; we cannot simply look at the first position in which the endpoints have different decimal digits. Therefore we change the invariant relation to  $0 \leq f < g \leq 1$ , when open intervals are involved, and we do not terminate the process when  $f = 0$  or  $g = 1$ .

```

{ register int open = x.l & 1;
  tt.dat[origin] = 10;
  tt.a = tt.b = origin;
  for (e = 1; bignum_compare(&gg, &tt) ≥ open; e++) bignum_times_ten(&tt);
  p = s;
  while (1) {
    bignum_times_ten(&ff);
    bignum_times_ten(&gg);
    for (j = '0'; bignum_compare(&ff, &tt) ≥ 0; j++)
      bignum_dec(&ff, &tt, #10000000), bignum_dec(&gg, &tt, #10000000);
    if (bignum_compare(&gg, &tt) ≥ open) break;
    *p++ = j;
    if (ff.a ≡ bignum_prec - 1 ∧ ¬open) goto done;    /* f = 0 in a closed interval */
  }
  for (k = j; bignum_compare(&gg, &tt) ≥ open; k++) bignum_dec(&gg, &tt, #10000000);
  *p++ = (j + 1 + k) ≫ 1;    /* the middle digit */
done: ;
}

```

This code is used in section 64.

**66.** The length of string  $s$  will be at most 17. For if  $f$  and  $g$  agree to 17 places, we have  $g/f < 1 + 10^{-16}$ ; but the ratio  $g/f$  is always  $\geq (1 + 2^{-52} + 2^{-53})/(1 + 2^{-52} - 2^{-53}) > 1 + 2 \times 10^{-16}$ .

```

bignum ff, gg; /* fractions or numerators of fractions */
bignum tt; /* power of ten (used as the denominator) */
char s[18];
register char *p;

```

**67.** At this point the significant digits are in string  $s$ , and  $s[0] \neq '0'$ . If we put a decimal point at the left of  $s$ , the result should be multiplied by  $10^e$ .

We prefer the output ‘300.’ to the form ‘3e2’, and we prefer ‘.03’ to ‘3e-2’. In general, the output will use an explicit exponent only if the alternative would take more than 18 characters.

```

if (e > 17 ∨ e < (int) strlen(s) - 17) printf("%c%s%se%d", s[0], (s[1] ? ". " : ""), s + 1, e - 1);
else if (e < 0) printf("%.0*d%s", -e, 0, s);
else if ((int) strlen(s) ≥ e) printf("%. *s.%s", e, s, s + e);
else printf("%s%.0*d.", s, e - (int) strlen(s), 0);

```

This code is used in section 54.

**68. Floating point input conversion.** Going the other way, we want to be able to convert a given decimal number into its floating binary equivalent. The following syntax is supported:

```

⟨digit⟩ → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
⟨digit string⟩ → ⟨digit⟩ | ⟨digit string⟩⟨digit⟩
⟨decimal string⟩ → ⟨digit string⟩. | .⟨digit string⟩ | ⟨digit string⟩.⟨digit string⟩
⟨optional sign⟩ → ⟨empty⟩ | + | -
⟨exponent⟩ → e⟨optional sign⟩⟨digit string⟩
⟨optional exponent⟩ → ⟨empty⟩ | ⟨exponent⟩
⟨floating magnitude⟩ → ⟨digit string⟩⟨exponent⟩ | ⟨decimal string⟩⟨optional exponent⟩ |
                        Inf | NaN | NaN.⟨digit string⟩
⟨floating constant⟩ → ⟨optional sign⟩⟨floating magnitude⟩
⟨decimal constant⟩ → ⟨optional sign⟩⟨digit string⟩

```

For example, ‘-3.’ is the floating constant #c008000000000000; ‘1e3’ and ‘1000’ are both equivalent to #408f400000000000; ‘NaN’ and ‘+NaN.5’ are both equivalent to #7ff8000000000000.

The *scan\_const* routine looks at a given string and finds the longest initial substring that matches the syntax of either ⟨decimal constant⟩ or ⟨floating constant⟩. It puts the corresponding value into the global octabyte variable *val*; it also puts the position of the first unscanned character in the global pointer variable *next\_char*. It returns 1 if a floating constant was found, 0 if a decimal constant was found, -1 if nothing was found. A decimal constant that doesn’t fit in an octabyte is computed modulo 2<sup>64</sup>.

The value of *exceptions* set by *scan\_const* is not necessarily correct.

⟨Subroutines 5⟩ +≡

```

static void bignum_double ARGS((bignum *));
int scan_const ARGS((char *));
    int scan_const(s)
    char *s;
{
    ⟨Local variables for scan_const 70⟩;
    val.h = val.l = 0;
    p = s;
    if (*p ≡ '+' ∨ *p ≡ '-') sign = *p++; else sign = '+';
    if (strncmp(p, "NaN", 3) ≡ 0) NaN = true, p += 3;
    else NaN = false;
    if ((isdigit(*p) ∧ ¬NaN) ∨ (*p ≡ '.' ∧ isdigit(*(p+1)))) ⟨Scan a number and return 73⟩;
    if (NaN) ⟨Return the standard NaN 71⟩;
    if (strncmp(p, "Inf", 3) ≡ 0) ⟨Return infinity 72⟩;
    next_char = s; return -1;
}

```

**69.** ⟨Global variables 4⟩ +≡

```

octa val; /* value returned by scan_const */
char *next_char; /* pointer returned by scan_const */

```

**70.** ⟨Local variables for scan\_const 70⟩ ≡

```

register char *p, *q; /* for string manipulations */
register bool NaN; /* are we processing a NaN? */
int sign; /* '+' or '-' */

```

See also sections 76 and 81.

This code is used in section 68.

**71.**  $\langle \text{Return the standard NaN } 71 \rangle \equiv$   

```
{
  next_char = p;
  val.h = #600000, exp = #3fe;
  goto packit;
}
```

This code is used in section 68.

**72.**  $\langle \text{Return infinity } 72 \rangle \equiv$   

```
{
  next_char = p + 3;
  goto make_it_infinite;
}
```

This code is used in section 68.

**73.** We saw above that a string of at most 17 digits is enough to characterize a floating point number, for purposes of output. But a much longer buffer for digits is needed when we're doing input. For example, consider the borderline quantity  $(1 + 2^{-53})/2^{1022}$ ; its decimal expansion, when written out exactly, is a number with more than 750 significant digits: 2.2250738585...8125e-308. If *any one* of those digits is increased, or if additional nonzero digits are added as in 2.2250738585...81250000001e-308, the rounded value is supposed to change from #0010000000000000 to #0010000000000001.

We assume here that the user prefers a perfectly correct answer to a speedy almost-correct one, so we implement the most general case.

$\langle \text{Scan a number and return } 73 \rangle \equiv$   

```
{
  for (q = buf0, dec_pt = (char *)0; isdigit(*p); p++) {
    val = oplus(val, shift_left(val, 2)); /* multiply by 5 */
    val = incr(shift_left(val, 1), *p - '0');
    if (q > buf0  $\vee$  *p  $\neq$  '0') {
      if (q < buf_max) *q++ = *p;
      else if (*(q - 1)  $\equiv$  '0') *(q - 1) = *p; }
  }
  if (NaN) *q++ = '1';
  if (*p  $\equiv$  '.')  $\langle \text{Scan a fraction part } 74 \rangle$ ;
  next_char = p;
  exp = 0;
  if (*p  $\equiv$  'e'  $\wedge$   $\neg$  NaN)  $\langle \text{Scan an exponent } 77 \rangle$ ;
  if (dec_pt)  $\langle \text{Return a floating point constant } 78 \rangle$ ;
  if (sign  $\equiv$  '-') val = ominus(zero_octa, val);
  return 0;
}
```

This code is used in section 68.

**74.**  $\langle \text{Scan a fraction part } 74 \rangle \equiv$

```

{
  dec_pt = q;
  p++;
  for (zeros = 0; isdigit(*p); p++)
    if (*p == '0' & q == buf0) zeros++;
    else if (q < buf_max) *q++ = *p;
    else if (*(q - 1) == '0') *(q - 1) = *p;
}

```

This code is used in section 73.

**75.** The buffer needs room for eight digits of padding at the left, followed by up to  $1022+53-307$  significant digits, followed by a “sticky” digit at position  $\text{buf\_max} - 1$ , and eight more digits of padding.

```

#define buf0 (buf + 8)
#define buf_max (buf + 777)
 $\langle \text{Global variables } 4 \rangle + \equiv$ 
static char buf[785] = "00000000"; /* where we put significant input digits */

```

**76.**  $\langle \text{Local variables for } \text{scan\_const } 70 \rangle + \equiv$

```

register char *dec_pt; /* position of decimal point in buf */
register int exp; /* scanned exponent; later used for raw binary exponent */
register int zeros = 0; /* leading zeros removed after decimal point */

```

**77.** Here we don’t advance  $\text{next\_char}$  and force a decimal point until we know that a syntactically correct exponent exists.

The code here will convert extra-large inputs like ‘9e+999999999999999’ into  $\infty$  and extra-small inputs into zero. Strange inputs like ‘-00.0e9999999’ must also be accommodated. (But we *don’t* try to deliver precise answers when there are a billion or more leading zeros.)

$\langle \text{Scan an exponent } 77 \rangle \equiv$

```

{ register char exp_sign;
  p++;
  if (*p == '+' || *p == '-') exp_sign = *p++; else exp_sign = '+';
  if (isdigit(*p)) {
    for (exp = *p++ - '0'; isdigit(*p); p++)
      if (exp < 1000000000) exp = 10 * exp + *p - '0';
    if (!dec_pt) dec_pt = q, zeros = 0;
    if (exp_sign == '-') exp = -exp;
    next_char = p;
  }
}

```

This code is used in section 73.

**78.**  $\langle \text{Return a floating point constant } 78 \rangle \equiv$

```

{
   $\langle \text{Move the digits from } \text{buf} \text{ to } \text{ff } 79 \rangle$ ;
   $\langle \text{Determine the binary fraction and binary exponent } 83 \rangle$ ;
  packit:  $\langle \text{Pack and round the answer } 84 \rangle$ ;
  return 1;
}

```

This code is used in section 73.



**79.** Now we get ready to compute the binary fraction bits, by putting the scanned input digits into a multiprecision fixed-point accumulator *ff* that spans the full necessary range. After this step, the number that we want to convert to floating binary will appear in *ff.dat*[*ff.a*], *ff.dat*[*ff.a* + 1], ..., *ff.dat*[*ff.b*]. The radix-10<sup>9</sup> digit in *ff*[36 − *k*] is understood to be multiplied by 10<sup>9<sup>k</sup></sup>, for 36 ≥ *k* ≥ −120.

```

⟨ Move the digits from buf to ff 79 ⟩ ≡
    x = buf + 341 + zeros − dec_pt − exp;
    if (q ≡ buf0 ∨ x ≥ 1413) {
        exp = −99999; goto packit;
    }
    if (x < 0) {
        make_it_infinite: exp = 99999; goto packit;
    }
    ff.a = x/9;
    for (p = q; p < q + 8; p++) *p = '0'; /* pad with trailing zeros */
    q = q − 1 − (q + 341 + zeros − dec_pt − exp) % 9; /* compute stopping place in buf */
    for (p = buf0 − x % 9, k = ff.a; p ≤ q ∧ k ≤ 156; p += 9, k++)
        ⟨ Put the 9-digit number *p...*(p + 8) into ff.dat[k] 80 ⟩;
    ff.b = k − 1;
    for (x = 0; p ≤ q; p += 9)
        if (strncmp(p, "000000000", 9) ≠ 0) x = 1;
    ff.dat[156] += x; /* nonzero digits that fall off the right are sticky */
    while (ff.dat[ff.b] ≡ 0) ff.b--;

```

This code is used in section 78.

```

80. ⟨ Put the 9-digit number *p...*(p + 8) into ff.dat[k] 80 ⟩ ≡
    {
        for (x = *p − '0', pp = p + 1; pp < p + 9; pp++) x = 10 * x + *pp − '0';
        ff.dat[k] = x;
    }

```

This code is used in section 79.

```

81. ⟨ Local variables for scan_const 70 ⟩ +=
    register int k, x;
    register char *pp;
    bignum ff, tt;

```

**82.** Here's a subroutine that is dual to *bignum\_times\_ten*. It changes *f* to *2f*, assuming that overflow will not occur and that the radix is  $10^9$ .

⟨Subroutines 5⟩ +≡

```
static void bignum_double(f)
    bignum *f;
{
    register tetra *p,*q;
    register int x, carry;
    for (p = &f->dat[f->b], q = &f->dat[f->a], carry = 0; p ≥ q; p--) {
        x = *p + *p + carry;
        if (x ≥ 1000000000) carry = 1, *p = x - 1000000000;
        else carry = 0, *p = x;
    }
    *p = carry;
    if (carry) f->a--;
    if (f->dat[f->b] ≡ 0 ∧ f->b > f->a) f->b--;
}
```

**83.** ⟨Determine the binary fraction and binary exponent 83⟩ ≡

```
val = zero_octa;
if (ff.a > 36) {
    for (exp = #3fe; ff.a > 36; exp--) bignum_double(&ff);
    for (k = 54; k; k--) {
        if (ff.dat[36]) {
            if (k ≥ 32) val.h |= 1 << (k - 32); else val.l |= 1 << k;
            ff.dat[36] = 0;
            if (ff.b ≡ 36) break; /* break if ff now zero */
        }
        bignum_double(&ff);
    }
} else {
    tt.a = tt.b = 36, tt.dat[36] = 2;
    for (exp = #3fe; bignum_compare(&ff, &tt) ≥ 0; exp++) bignum_double(&tt);
    for (k = 54; k; k--) {
        bignum_double(&ff);
        if (bignum_compare(&ff, &tt) ≥ 0) {
            if (k ≥ 32) val.h |= 1 << (k - 32); else val.l |= 1 << k;
            bignum_dec(&ff, &tt, 1000000000);
            if (ff.a ≡ bignum_prec - 1) break; /* break if ff now zero */
        }
    }
}
if (k ≡ 0) val.l |= 1; /* add sticky bit if ff nonzero */
```

This code is used in section 78.

**84.** We need to be careful that the input `'NaN.9999999999999999999'` doesn't get rounded up; it is supposed to yield `#7fffffffffffffff`.

Although the input `'NaN.0'` is illegal, strictly speaking, we silently convert it to `#7ff0000000000001`—a number that would be output as `'NaN.0000000000000002'`.

$$\langle \text{Pack and round the answer } 84 \rangle \equiv$$

```
val = fpack(val, exp, sign, ROUND_NEAR);
```

$$\mathbf{if} \ (NaN) \ \{$$

```
if ((val.h & #7fffffff) == #40000000) val.h |= #7fffffff, val.l = #ffffff;
```

```

else if ((val.h & #7fffffff) ≡ #3ff00000 ∧ ¬val.l) val.h |= #40000000, val.l = 1;

```

```
else val.h |= #40000000;
```

$$\}$$

This code is used in section 78.

**85. Floating point remainders.** In this section we implement the remainder of the floating point operations—one of which happens to be the operation of taking the remainder.

The easiest task remaining is to compare two floating point quantities. Routine *fcomp* returns  $-1$  if  $y < z$ ,  $0$  if  $y = z$ ,  $+1$  if  $y > z$ , and  $+2$  if  $y$  and  $z$  are unordered.

⟨Subroutines 5⟩  $\equiv$

```

int fcomp ARGS((octa, octa));
    int fcomp(y, z)
    octa y, z;
{
    ftype yt, zt;
    int ye, ze;
    char ys, zs;
    octa yf, zf;
    register int x;
    yt = funpack(y, &yf, &ye, &ys);
    zt = funpack(z, &zf, &ze, &zs);
    switch (4 * yt + zt) {
    case 4 * nan + nan: case 4 * zro + nan: case 4 * num + nan: case 4 * inf + nan: case 4 * nan + zro:
        case 4 * nan + num: case 4 * nan + inf: return 2;
    default: case 4 * zro + zro: return 0;
    case 4 * zro + num: case 4 * num + zro: case 4 * zro + inf: case 4 * inf + zro: case 4 * num + num:
        case 4 * num + inf: case 4 * inf + num: case 4 * inf + inf:
        if (ys ≠ zs) x = 1;
        else if (y.h > z.h) x = 1;
        else if (y.h < z.h) x = -1;
        else if (y.l > z.l) x = 1;
        else if (y.l < z.l) x = -1;
        else return 0;
        break;
    }
    return (ys ≡ '-' ? -x : x);
}

```

**86.** Several MMIX operations act on a single floating point number and accept an arbitrary rounding mode. For example, consider the operation of rounding to the nearest floating point integer:

⟨Subroutines 5⟩ +≡

```

octa fintegerize ARGS((octa,int));
    octa fintegerize(z,r)
    octa z;      /* the operand */
    int r;      /* the rounding mode */
{
    ftype zt;
    int ze;
    char zs;
    octa xf, zf;
    zt = funpack(z, &zf, &ze, &zs);
    if (¬r) r = cur_round;
    switch (zt) {
    case nan: if (¬(z.h & #80000)) { exceptions |= I_BIT; z.h |= #80000; }
    case inf: case zro: default: return z;
    case num: ⟨Integerize and return 87⟩;
    }
}

```

**87.** ⟨Integerize and **return** 87⟩ ≡

```

if (ze ≥ 1074) return fpack(zf, ze, zs, ROUND_OFF);    /* already an integer */
if (ze ≤ 1020) xf.h = 0, xf.l = 1;
else { octa oo;
    xf = shift_right(zf, 1074 − ze, 1);
    oo = shift_left(xf, 1074 − ze);
    if (oo.l ≠ zf.l ∨ oo.h ≠ zf.h) xf.l |= 1;    /* sticky bit */
}
switch (r) {
case ROUND_DOWN: if (zs ≡ '−') xf = incr(xf, 3); break;
case ROUND_UP: if (zs ≠ '−') xf = incr(xf, 3);
case ROUND_OFF: break;
case ROUND_NEAR: xf = incr(xf, xf.l & 4 ? 2 : 1); break;
}
xf.l &= #ffffffc;
if (ze ≥ 1022) return fpack(shift_left(xf, 1074 − ze), ze, zs, ROUND_OFF);
if (xf.l) xf.h = #3ff00000, xf.l = 0;
if (zs ≡ '−') xf.h |= sign_bit;
return xf;

```

This code is used in section 86.

**88.** To convert floating point to fixed point, we use *fixit*.

⟨Subroutines 5⟩ +≡

```

octa fixit ARGS((octa,int));
    octa fixit(z,r)
    octa z;      /* the operand */
    int r;      /* the rounding mode */
{
    ftype zt;
    int ze;
    char zs;
    octa zf,o;
    zt = funpack(z, &zf, &ze, &zs);
    if (¬r) r = cur_round;
    switch (zt) {
    case nan: case inf: exceptions |= I_BIT; return z;
    case zro: default: return zero_octa;
    case num: if (funpack(fintegerize(z,r), &zf, &ze, &zs) ≡ zro) return zero_octa;
        if (ze ≤ 1076) o = shift_right(zf, 1076 − ze, 1);
        else {
            if (ze > 1085 ∨ (ze ≡ 1085 ∧ (zf.h > #400000 ∨
                (zf.h ≡ #400000 ∧ (zf.l ∨ zs ≠ '−'))))) exceptions |= W_BIT;
            if (ze ≥ 1140) return zero_octa;
            o = shift_left(zf, ze − 1076);
        }
    }
    return (zs ≡ '−' ? ominus(zero_octa,o) : o);
}

```

**89.** Going the other way, we can specify not only a rounding mode but whether the given fixed point octabyte is signed or unsigned, and whether the result should be rounded to short precision.

⟨Subroutines 5⟩ +≡

```

octa floatit ARGS((octa,int,int,int));
    octa floatit(z,r,u,p)
    octa z;    /* octabyte to float */
    int r;      /* rounding mode */
    int u;      /* unsigned? */
    int p;      /* short precision? */
{
    int e; char s;
    register int t;
    exceptions = 0;
    if ( $\neg z.h \wedge \neg z.l$ ) return zero_octa;
    if ( $\neg r$ ) r = cur_round;
    if ( $\neg u \wedge (z.h \ \& \ sign\_bit)$ ) s = '-', z = ominus(zero_octa, z); else s = '+';
    e = 1076;
    while (z.h < #400000) e--, z = shift_left(z, 1);
    while (z.h ≥ #800000) {
        e++;
        t = z.l & 1;
        z = shift_right(z, 1, 1);
        z.l |= t;
    }
    if (p) ⟨Convert to short float 90⟩;
    return fpack(z, e, s, r);
}

```

**90.** ⟨Convert to short float 90⟩ ≡

```

{
    register int ex; register tetra t;
    t = sfpack(z, e, s, r);
    ex = exceptions;
    sfunpack(t, &z, &e, &s);
    exceptions = ex;
}

```

This code is used in section 89.

**91.** The square root operation is more interesting.

⟨Subroutines 5⟩ +≡

```

octa froot ARGS((octa,int));
    octa froot(z,r)
    octa z;      /* the operand */
    int r;      /* the rounding mode */
{
    ftype zt;
    int ze;
    char zs;
    octa x,xf,rf,zf;
    register int xe,k;
    if ( $\neg r$ ) r = cur_round;
    zt = funpack(z, &zf, &ze, &zs);
    if (zs ≡ '−' ∧ zt ≠ zro) exceptions |= I_BIT, x = standardNaN;
    else switch (zt) {
        case nan: if ( $\neg(z.h \& \#80000)$ ) exceptions |= I_BIT, z.h |=  $\#80000$ ;
        return z;
        default: case inf: case zro: x = z; break;
        case num: ⟨Take the square root and return 92⟩;
    }
    if (zs ≡ '−') x.h |= sign_bit;
    return x;
}

```

**92.** The square root can be found by an adaptation of the old pencil-and-paper method. If  $n = \lfloor \sqrt{s} \rfloor$ , where  $s$  is an integer, we have  $s = n^2 + r$  where  $0 \leq r \leq 2n$ ; this invariant can be maintained if we replace  $s$  by  $4s + (0, 1, 2, 3)$  and  $n$  by  $2n + (0, 1)$ . The following code implements this idea with  $2n$  in *xf* and  $r$  in *rf*. (It could easily be made to run about twice as fast.)

⟨Take the square root and **return** 92⟩ ≡

```

xf.h = 0, xf.l = 2;
xe = (ze +  $\#3fe$ )  $\gg$  1;
if (ze & 1) zf = shift_left(zf, 1);
rf.h = 0, rf.l = (zf.h  $\gg$  22) − 1;
for (k = 53; k; k−−) {
    rf = shift_left(rf, 2); xf = shift_left(xf, 1);
    if (k ≥ 43) rf = incr(rf, (zf.h  $\gg$  (2 * (k − 43))) & 3);
    else if (k ≥ 27) rf = incr(rf, (zf.l  $\gg$  (2 * (k − 27))) & 3);
    if ((rf.l > xf.l ∧ rf.h ≥ xf.h) ∨ rf.h > xf.h) {
        xf.l++; rf = ominus(rf, xf); xf.l++;
    }
}
if (rf.h ∨ rf.l) xf.l++; /* sticky bit */
return fpack(xf, xe, '−', r);

```

This code is used in section 91.



**93.** And finally, the genuine floating point remainder. Subroutine *fremstep* either calculates  $y \bmod z$  or reduces  $y$  to a smaller number having the same remainder with respect to  $z$ . In the latter case the `E_BIT` is set in *exceptions*. A third parameter, *delta*, gives a decrease in exponent that is acceptable for incomplete results; if *delta* is sufficiently large, say 2500, the correct result will always be obtained in one step of *fremstep*.

⟨Subroutines 5⟩ +=

```

octa fremstep ARGS((octa,octa,int));
    octa fremstep(y, z, delta)
    octa y, z;
    int delta;
{
    ftype yt, zt;
    int ye, ze;
    char xs, ys, zs;
    octa x, xf, yf, zf;
    register int xe, thresh, odd;
    yt = funpack(y, &yf, &ye, &ys);
    zt = funpack(z, &zf, &ze, &zs);
    switch (4 * yt + zt) {
    ⟨The usual NaN cases 42⟩;
    case 4 * zro + zro: case 4 * num + zro: case 4 * inf + zro: case 4 * inf + num: case 4 * inf + inf:
        x = standard_NaN;
        exceptions |= I_BIT; break;
    case 4 * zro + num: case 4 * zro + inf: case 4 * num + inf: return y;
    case 4 * num + num: ⟨Remainderize nonzero numbers and return 94⟩;
    default: zero_out: x = zero_octa;
    }
    if (ys ≡ '−') x.h |= sign_bit;
    return x;
}

```

**94.** If there's a huge difference in exponents and the remainder is nonzero, this computation will take a long time. One could compute  $(2^n y) \bmod z$  much more quickly for large  $n$  by using  $O(\log n)$  multiplications modulo  $z$ , but the floating remainder operation isn't important enough to justify such expensive hardware.

Results of floating remainder are always exact, so the rounding mode is immaterial.

⟨Remainderize nonzero numbers and **return** 94⟩ ≡

```

    odd = 0; /* will be 1 if we've subtracted an odd multiple of z from y */
    thresh = ye − delta;
    if (thresh < ze) thresh = ze;
    while (ye ≥ thresh) ⟨Reduce (ye, yf) by a multiple of zf; goto zero_out if the remainder is zero, goto
        try_complement if appropriate 95⟩;
    if (ye ≥ ze) {
        exceptions |= E_BIT; return fpack(yf, ye, ys, ROUND_OFF);
    }
    if (ye < ze − 1) return fpack(yf, ye, ys, ROUND_OFF);
    yf = shift_right(yf, 1, 1);
try_complement: xf = ominus(zf, yf), xe = ze, xs = '+' + '−' − ys;
    if (xf.h > yf.h ∨ (xf.h ≡ yf.h ∧ (xf.l > yf.l ∨ (xf.l ≡ yf.l ∧ ¬odd)))) xf = yf, xs = ys;
    while (xf.h < #400000) xe −, xf = shift_left(xf, 1);
    return fpack(xf, xe, xs, ROUND_OFF);

```

This code is used in section 93.

**95.** Here we are careful not to change the sign of  $y$ , because a remainder of 0 is supposed to inherit the original sign of  $y$ .

⟨ Reduce  $(ye, yf)$  by a multiple of  $zf$ ; **goto** *zero\_out* if the remainder is zero, **goto** *try\_complement* if appropriate 95 ⟩  $\equiv$

```
{
  if ( $yf.h \equiv zf.h \wedge yf.l \equiv zf.l$ ) goto zero_out;
  if ( $yf.h < zf.h \vee (yf.h \equiv zf.h \wedge yf.l < zf.l)$ ) {
    if ( $ye \equiv ze$ ) goto try_complement;
     $ye \text{---}, yf = \text{shift\_left}(yf, 1)$ ;
  }
   $yf = \text{ominus}(yf, zf)$ ;
  if ( $ye \equiv ze$ )  $odd = 1$ ;
  while ( $yf.h < \#400000$ )  $ye \text{---}, yf = \text{shift\_left}(yf, 1)$ ;
}
```

This code is used in section 94.

## 96. Index.

- `--STDC--`: [2](#).
- `a`: [28](#), [29](#), [59](#).
- `acc`: [8](#), [11](#), [12](#), [13](#), [19](#).
- accuracy loss: [31](#).
- `ARGS`: [2](#), [5](#), [6](#), [7](#), [8](#), [12](#), [13](#), [24](#), [25](#), [26](#), [27](#), [28](#),  
[29](#), [31](#), [34](#), [37](#), [38](#), [39](#), [40](#), [41](#), [44](#), [46](#), [50](#), [54](#),  
[68](#), [85](#), [86](#), [88](#), [89](#), [91](#), [93](#).
- `aux`: [8](#), [9](#), [11](#), [12](#), [13](#), [14](#), [19](#), [24](#), [43](#), [45](#).
- `b`: [28](#), [29](#), [59](#).
- bignum**: [54](#), [59](#), [60](#), [61](#), [62](#), [66](#), [68](#), [81](#), [82](#).
- `bignum_compare`: [54](#), [61](#), [64](#), [65](#), [83](#).
- `bignum_dec`: [54](#), [62](#), [65](#), [83](#).
- `bignum_double`: [68](#), [82](#), [83](#).
- `bignum_prec`: [59](#), [62](#), [65](#), [83](#).
- `bignum_times_ten`: [54](#), [60](#), [64](#), [65](#), [82](#).
- binary-to-decimal conversion: [54](#).
- bool**: [1](#), [9](#), [29](#), [70](#).
- `bool_mult`: [29](#).
- `borrow`: [62](#).
- `buf`: [75](#), [76](#), [79](#).
- `buf_max`: [73](#), [74](#), [75](#).
- `buf0`: [73](#), [74](#), [75](#), [79](#).
- `byte_diff`: [27](#), [28](#).
- `c`: [29](#).
- `carry`: [60](#), [82](#).
- `count_bits`: [26](#), [28](#).
- `cur_round`: [30](#), [40](#), [43](#), [45](#), [46](#), [47](#), [86](#), [88](#), [89](#), [91](#).
- `d`: [13](#), [27](#), [46](#), [50](#).
- `D_BIT`: [31](#).
- `dat`: [59](#), [60](#), [61](#), [62](#), [63](#), [64](#), [65](#), [79](#), [80](#), [82](#), [83](#).
- `dec_pt`: [73](#), [74](#), [76](#), [77](#), [79](#).
- decimal-to-binary conversion: [68](#).
- `delta`: [6](#), [93](#), [94](#).
- `done`: [65](#).
- `e`: [31](#), [34](#), [37](#), [38](#), [39](#), [40](#), [50](#), [56](#), [89](#).
- `E_BIT`: [31](#), [93](#), [94](#).
- `ee`: [37](#), [38](#), [50](#), [51](#), [53](#).
- `ef`: [50](#), [51](#), [53](#).
- `es`: [50](#).
- `et`: [50](#).
- `ex`: [90](#).
- `exceptions`: [31](#), [32](#), [33](#), [35](#), [36](#), [37](#), [38](#), [40](#), [41](#), [42](#),  
[44](#), [46](#), [68](#), [86](#), [88](#), [89](#), [90](#), [91](#), [93](#), [94](#).
- `exp`: [71](#), [73](#), [76](#), [77](#), [79](#), [83](#), [84](#).
- `exp_sign`: [77](#).
- `f`: [31](#), [34](#), [37](#), [38](#), [39](#), [40](#), [56](#), [60](#), [61](#), [62](#), [82](#).
- `false`: [1](#), [24](#), [68](#).
- `fcomp`: [85](#).
- `fdivide`: [44](#).
- `fepscomp`: [50](#).
- `ff`: [63](#), [64](#), [65](#), [66](#), [79](#), [80](#), [81](#), [83](#).
- `fintegerize`: [86](#), [88](#).
- `fixit`: [88](#).
- `floatit`: [89](#).
- `fmult`: [41](#).
- `fpack`: [31](#), [34](#), [36](#), [39](#), [43](#), [45](#), [46](#), [47](#), [49](#), [84](#),  
[87](#), [89](#), [92](#), [94](#).
- `fplus`: [46](#).
- `fremstep`: [93](#).
- `froot`: [91](#).
- ftype**: [36](#), [37](#), [38](#), [39](#), [40](#), [41](#), [44](#), [46](#), [85](#), [86](#),  
[88](#), [91](#), [93](#).
- `funpack`: [36](#), [37](#), [40](#), [41](#), [44](#), [46](#), [50](#), [85](#), [86](#),  
[88](#), [91](#), [93](#).
- `g`: [56](#), [61](#), [62](#).
- `gg`: [63](#), [64](#), [65](#), [66](#).
- Gill, Stanley: [26](#).
- Gillies, Donald Bruce: [26](#).
- `h`: [3](#).
- `i`: [8](#), [13](#).
- `I_BIT`: [31](#), [40](#), [41](#), [42](#), [44](#), [46](#), [86](#), [88](#), [91](#), [93](#).
- `incr`: [6](#), [33](#), [53](#), [55](#), [73](#), [87](#), [92](#).
- `inf`: [36](#), [37](#), [38](#), [39](#), [40](#), [41](#), [42](#), [44](#), [46](#), [50](#), [85](#),  
[86](#), [88](#), [91](#), [93](#).
- `inf_octa`: [4](#), [39](#), [41](#), [44](#), [46](#).
- `isdigit`: [68](#), [73](#), [74](#), [77](#).
- `j`: [8](#), [13](#), [56](#).
- `k`: [8](#), [13](#), [29](#), [56](#), [81](#), [91](#).
- Knuth, Donald Ervin: [58](#).
- `l`: [3](#).
- `list`: [2](#).
- `load_sf`: [39](#).
- `m`: [27](#).
- `magic_offset`: [63](#).
- `make_it_infinite`: [72](#), [79](#).
- `mask`: [13](#), [18](#).
- Miller, Jeffrey Charles Percy: [26](#).
- multiprecision conversion: [54](#), [68](#).
- multiprecision division: [13](#).
- multiprecision multiplication: [8](#).
- `n`: [13](#).
- `NaN`: [68](#), [70](#), [73](#), [84](#).
- `nan`: [36](#), [37](#), [38](#), [39](#), [40](#), [42](#), [50](#), [85](#), [86](#), [88](#), [91](#).
- `neg_one`: [4](#), [24](#).
- `negate_q`: [24](#).
- `next_char`: [68](#), [69](#), [71](#), [72](#), [73](#), [77](#).
- `num`: [36](#), [37](#), [38](#), [39](#), [40](#), [41](#), [42](#), [44](#), [46](#), [50](#), [85](#),  
[86](#), [88](#), [91](#), [93](#).
- `o`: [29](#), [31](#), [34](#), [47](#), [50](#), [88](#).
- `O_BIT`: [31](#), [33](#), [35](#).
- `oand`: [25](#).
- `oandn`: [25](#).

**octa:** [3](#), [4](#), [5](#), [6](#), [7](#), [8](#), [9](#), [12](#), [13](#), [24](#), [25](#), [29](#), [31](#),  
[34](#), [37](#), [38](#), [39](#), [40](#), [41](#), [44](#), [46](#), [47](#), [50](#), [54](#), [56](#),  
[69](#), [85](#), [86](#), [87](#), [88](#), [89](#), [91](#), [93](#).

**odd:** [93](#), [94](#), [95](#).

**odiv:** [13](#), [24](#), [45](#).

**ominus:** [5](#), [12](#), [24](#), [47](#), [53](#), [73](#), [88](#), [89](#), [92](#), [94](#), [95](#).

**omult:** [8](#), [12](#), [43](#).

**oo:** [31](#), [34](#), [47](#), [49](#), [50](#), [53](#), [87](#).

**open:** [65](#).

**oplus:** [5](#), [47](#), [53](#), [73](#).

**origin:** [63](#), [64](#), [65](#).

**overflow:** [9](#), [12](#), [24](#).

**oxor:** [25](#).

**o0:** [34](#).

**p:** [60](#), [61](#), [62](#), [66](#), [70](#), [82](#), [89](#).

**packit:** [71](#), [78](#), [79](#).

**pp:** [61](#), [80](#), [81](#).

**print\_float:** [54](#), [59](#).

**printf:** [54](#), [55](#), [57](#), [67](#).

prototypes for functions: [2](#).

**q:** [13](#), [24](#), [60](#), [61](#), [62](#), [70](#), [82](#).

**qhat:** [13](#), [20](#), [21](#), [22](#), [23](#).

**qq:** [61](#), [62](#).

**r:** [31](#), [34](#), [62](#), [86](#), [88](#), [89](#), [91](#).

radix conversion: [54](#), [68](#).

**rf:** [91](#), [92](#).

**rhat:** [13](#), [21](#).

Rossmannith, Peter: [26](#).

ROUND\_DOWN: [30](#), [33](#), [35](#), [46](#), [87](#).

ROUND\_NEAR: [30](#), [33](#), [35](#), [84](#), [87](#).

ROUND\_OFF: [30](#), [33](#), [35](#), [39](#), [46](#), [87](#), [94](#).

ROUND\_UP: [30](#), [33](#), [35](#), [87](#).

**s:** [7](#), [31](#), [34](#), [37](#), [38](#), [39](#), [40](#), [50](#), [66](#), [68](#), [89](#).

**scan\_const:** [68](#), [69](#).

Schwoon, Stefan: [26](#).

**sfpack:** [34](#), [39](#), [40](#), [90](#).

**sfunpack:** [38](#), [39](#), [90](#).

**shift\_left:** [7](#), [31](#), [34](#), [37](#), [38](#), [43](#), [45](#), [47](#), [49](#), [51](#), [52](#),  
[53](#), [55](#), [63](#), [73](#), [87](#), [88](#), [89](#), [92](#), [94](#), [95](#).

**shift\_right:** [7](#), [29](#), [31](#), [33](#), [39](#), [45](#), [47](#), [49](#), [51](#), [53](#),  
[63](#), [87](#), [88](#), [89](#), [94](#).

**sign:** [68](#), [70](#), [73](#), [84](#).

**sign\_bit:** [4](#), [12](#), [24](#), [33](#), [35](#), [37](#), [38](#), [39](#), [40](#), [41](#), [44](#),  
[46](#), [54](#), [87](#), [89](#), [91](#), [93](#).

**signed\_odiv:** [24](#).

**signed\_omult:** [12](#).

Singh, Balbir: [26](#).

standard\_NaN: [4](#), [41](#), [44](#), [46](#), [91](#), [93](#).

sticky bit: [31](#), [34](#), [49](#), [53](#), [79](#), [87](#).

**store\_sf:** [40](#).

**strlen:** [67](#).

**strncmp:** [68](#), [79](#).

**sy:** [24](#).

syntax of floating point constants: [68](#).

system dependencies: [3](#).

**sz:** [24](#).

**t:** [8](#), [13](#), [39](#), [40](#), [89](#), [90](#).

**tetra:** [3](#), [7](#), [8](#), [13](#), [23](#), [26](#), [27](#), [28](#), [29](#), [34](#), [38](#), [39](#),  
[40](#), [54](#), [59](#), [60](#), [61](#), [62](#), [82](#), [90](#).

**thresh:** [93](#), [94](#).

tininess: [31](#).

**true:** [1](#), [24](#), [68](#).

**try\_complement:** [94](#), [95](#).

**tt:** [65](#), [66](#), [81](#), [83](#).

**u:** [7](#), [8](#), [13](#), [89](#).

U\_BIT: [31](#), [33](#), [35](#).

underflow: [31](#).

**v:** [8](#), [13](#).

V\_BIT: [31](#).

**val:** [68](#), [69](#), [71](#), [73](#), [83](#), [84](#).

**vh:** [13](#), [17](#), [21](#).

**vmh:** [13](#), [17](#), [21](#).

**w:** [8](#).

W\_BIT: [31](#), [88](#).

Wheeler, David John: [26](#).

Wilkes, Maurice Vincent: [26](#).

**wyde\_diff:** [28](#).

**x:** [5](#), [6](#), [13](#), [25](#), [26](#), [27](#), [29](#), [37](#), [38](#), [39](#), [40](#), [41](#), [44](#),  
[46](#), [54](#), [60](#), [62](#), [81](#), [82](#), [85](#), [91](#), [93](#).

X\_BIT: [31](#), [33](#), [35](#).

**xe:** [41](#), [43](#), [44](#), [45](#), [46](#), [47](#), [49](#), [91](#), [92](#), [93](#), [94](#).

**xf:** [41](#), [43](#), [44](#), [45](#), [46](#), [47](#), [86](#), [87](#), [91](#), [92](#), [93](#), [94](#).

**xor:** [29](#).

**xs:** [41](#), [43](#), [44](#), [45](#), [46](#), [47](#), [93](#), [94](#).

**xx:** [26](#).

**y:** [5](#), [6](#), [7](#), [8](#), [12](#), [13](#), [24](#), [25](#), [27](#), [28](#), [29](#), [41](#), [44](#),  
[46](#), [50](#), [85](#), [93](#).

**ye:** [41](#), [43](#), [44](#), [45](#), [46](#), [47](#), [48](#), [50](#), [51](#), [52](#), [85](#),  
[93](#), [94](#), [95](#).

**yf:** [41](#), [43](#), [44](#), [45](#), [46](#), [47](#), [48](#), [49](#), [50](#), [51](#), [52](#),  
[53](#), [85](#), [93](#), [94](#), [95](#).

**yhl:** [7](#).

**ylh:** [7](#).

**ys:** [41](#), [44](#), [46](#), [47](#), [48](#), [49](#), [50](#), [53](#), [85](#), [93](#), [94](#).

**yt:** [41](#), [44](#), [46](#), [50](#), [52](#), [85](#), [93](#).

**yy:** [24](#).

**z:** [5](#), [8](#), [12](#), [13](#), [24](#), [25](#), [27](#), [28](#), [29](#), [39](#), [40](#), [41](#), [44](#),  
[46](#), [50](#), [85](#), [86](#), [88](#), [89](#), [91](#), [93](#).

Z\_BIT: [31](#), [44](#).

**ze:** [41](#), [43](#), [44](#), [45](#), [46](#), [47](#), [48](#), [50](#), [51](#), [52](#), [85](#), [86](#),  
[87](#), [88](#), [91](#), [92](#), [93](#), [94](#), [95](#).

**zero\_exponent:** [36](#), [37](#), [38](#), [51](#).

**zero\_octa:** [4](#), [24](#), [29](#), [31](#), [39](#), [41](#), [44](#), [45](#), [46](#), [53](#),  
[73](#), [83](#), [88](#), [89](#), [93](#).

*zero\_out*: [93](#), [95](#).  
*zeros*: [74](#), [76](#), [77](#), [79](#).  
*zf*: [41](#), [43](#), [44](#), [45](#), [46](#), [47](#), [48](#), [49](#), [50](#), [51](#), [52](#), [53](#),  
[85](#), [86](#), [87](#), [88](#), [91](#), [92](#), [93](#), [94](#), [95](#).  
*zro*: [36](#), [37](#), [38](#), [39](#), [40](#), [41](#), [42](#), [44](#), [46](#), [50](#), [52](#),  
[85](#), [86](#), [88](#), [91](#), [93](#).  
*zs*: [41](#), [44](#), [46](#), [47](#), [48](#), [49](#), [50](#), [53](#), [85](#), [86](#), [87](#),  
[88](#), [91](#), [93](#).  
*zt*: [41](#), [44](#), [46](#), [50](#), [52](#), [85](#), [86](#), [88](#), [91](#), [93](#).  
*zz*: [24](#).

- ⟨ Add nonzero numbers and **return** 47 ⟩ Used in section 46.
- ⟨ Adjust for difference in exponents 49 ⟩ Used in section 47.
- ⟨ Check that  $x < z$ ; otherwise give trivial answer 14 ⟩ Used in section 13.
- ⟨ Compare two numbers with respect to epsilon and **return** 51 ⟩ Used in section 50.
- ⟨ Compute the difference of fraction parts,  $o$  53 ⟩ Used in section 51.
- ⟨ Compute the significant digits in the large-exponent case 65 ⟩ Used in section 64.
- ⟨ Compute the significant digits  $s$  and decimal exponent  $e$  64 ⟩ Used in section 54.
- ⟨ Convert to short float 90 ⟩ Used in section 89.
- ⟨ Determine the binary fraction and binary exponent 83 ⟩ Used in section 78.
- ⟨ Determine the number of significant places  $n$  in the divisor  $v$  16 ⟩ Used in section 13.
- ⟨ Determine the quotient digit  $q[j]$  20 ⟩ Used in section 13.
- ⟨ Divide nonzero numbers and **return** 45 ⟩ Used in section 44.
- ⟨ Exchange  $y$  with  $z$  48 ⟩ Used in sections 47 and 51.
- ⟨ Extract the exponent  $e$  and determine the fraction interval  $[f..g]$  or  $(f..g)$  55 ⟩ Used in section 54.
- ⟨ Find the trial quotient,  $\hat{q}$  21 ⟩ Used in section 20.
- ⟨ Global variables 4, 9, 30, 32, 69, 75 ⟩ Used in section 1.
- ⟨ Handle the special case when the fraction part is zero 57 ⟩ Used in section 55.
- ⟨ If the result was negative, decrease  $\hat{q}$  by 1 23 ⟩ Used in section 20.
- ⟨ Integerize and **return** 87 ⟩ Used in section 86.
- ⟨ Local variables for *print\_float* 56, 66 ⟩ Used in section 54.
- ⟨ Local variables for *scan\_const* 70, 76, 81 ⟩ Used in section 68.
- ⟨ Move the digits from *buf* to *ff* 79 ⟩ Used in section 78.
- ⟨ Multiply nonzero numbers and **return** 43 ⟩ Used in section 41.
- ⟨ Normalize the divisor 17 ⟩ Used in section 13.
- ⟨ Other type definitions 36, 59 ⟩ Used in section 1.
- ⟨ Pack and round the answer 84 ⟩ Used in section 78.
- ⟨ Pack  $q$  and  $u$  to *acc* and *aux* 19 ⟩ Used in section 13.
- ⟨ Pack  $w$  into the outputs *aux* and *acc* 11 ⟩ Used in section 8.
- ⟨ Print the significant digits with proper context 67 ⟩ Used in section 54.
- ⟨ Put the 9-digit number  $*p \dots *(p+8)$  into *ff.dat*[ $k$ ] 80 ⟩ Used in section 79.
- ⟨ Reduce  $(ye, yf)$  by a multiple of  $zf$ ; **goto** *zero\_out* if the remainder is zero, **goto** *try\_complement* if appropriate 95 ⟩ Used in section 94.
- ⟨ Remainderize nonzero numbers and **return** 94 ⟩ Used in section 93.
- ⟨ Return a floating point constant 78 ⟩ Used in section 73.
- ⟨ Return infinity 72 ⟩ Used in section 68.
- ⟨ Return the standard NaN 71 ⟩ Used in section 68.
- ⟨ Round and return the result 33 ⟩ Used in section 31.
- ⟨ Round and return the short result 35 ⟩ Used in section 34.
- ⟨ Scan a fraction part 74 ⟩ Used in section 73.
- ⟨ Scan a number and **return** 73 ⟩ Used in section 68.
- ⟨ Scan an exponent 77 ⟩ Used in section 73.
- ⟨ Store  $f$  and  $g$  as multiprecise integers 63 ⟩ Used in section 54.
- ⟨ Stuff for C preprocessor 2 ⟩ Used in section 1.
- ⟨ Subroutines 5, 6, 7, 8, 12, 13, 24, 25, 26, 27, 28, 29, 31, 34, 37, 38, 39, 40, 41, 44, 46, 50, 54, 60, 61, 62, 68, 82, 85, 86, 88, 89, 91, 93 ⟩ Used in section 1.
- ⟨ Subtract  $b^j \hat{q} v$  from  $u$  22 ⟩ Used in section 20.
- ⟨ Take the square root and **return** 92 ⟩ Used in section 91.
- ⟨ Tetrabyte and octabyte type definitions 3 ⟩ Used in section 1.
- ⟨ The usual NaN cases 42 ⟩ Used in sections 41, 44, 46, and 93.
- ⟨ Unnormalize the remainder 18 ⟩ Used in section 13.
- ⟨ Unpack the dividend and divisor to  $u$  and  $v$  15 ⟩ Used in section 13.
- ⟨ Unpack the multiplier and multiplicand to  $u$  and  $v$  10 ⟩ Used in section 8.

⟨ Unsubnormalize  $y$  and  $z$ , if they are subnormal [52](#) ⟩    Used in section [51](#).

## MMIX-ARITH

	Section	Page
Introduction .....	<a href="#">1</a>	1
Multiplication .....	<a href="#">8</a>	4
Division .....	<a href="#">13</a>	6
Bit fiddling .....	<a href="#">25</a>	9
Floating point packing and unpacking .....	<a href="#">30</a>	11
Floating multiplication and division .....	<a href="#">41</a>	17
Floating addition and subtraction .....	<a href="#">46</a>	19
Floating point output conversion .....	<a href="#">54</a>	23
Floating point input conversion .....	<a href="#">68</a>	30
Floating point remainders .....	<a href="#">85</a>	36
Index .....	<a href="#">96</a>	43