

# The LEDA class **real** number\*

Christoph Burnikel      Kurt Mehlhorn      Stefan Schirra

January 16, 1996

## Abstract

We describe the implementation of the LEDA [MN95, Näh95] data type **real**. Every integer is a real and reals are closed under the operations addition, subtraction, multiplication, division and squareroot. The main features of the data type real are

- The user-interface is similar to that of the built-in data type double.
- All comparison operators  $\{>, \geq, <, \leq, =\}$  are *exact*. In order to determine the sign of a real number  $x$  the data type first computes a rational number  $q$  such that  $|x| \leq q$  implies  $x = 0$  and then computes an approximation of  $x$  of sufficient precision to decide the sign of  $x$ . The user may assist the data type by providing a separation bound  $q$ .
- The data type also allows to evaluate real expressions with arbitrary precision. One may either set the mantissae length of the underlying floating point system and then evaluate the expression with that mantissa length or one may specify an error bound  $q$ . The data type then computes an approximation with absolute error at most  $q$ .

The implementation of the data type real is based on the LEDA data types **integer** and **bigfloat** which are the types of arbitrary precision integers and floating point numbers, respectively. The implementation takes various shortcuts for increased efficiency, e.g., a **double** approximation of any real number together with an error bound is maintained and tests are first performed on these approximations. A high precision computation is only started when the test on the **double** approximation is inconclusive.

---

\*This work was supported in part by the ESPRIT Basic Research Actions Program of the EC under contract No. 7141(ALCOM II) and the BMFT(Förderungskennzeichen ITS 9103).

## Contents

1. Introduction . . . . .	2
2. The manual pages of data type real . . . . .	3
3. A Short Review of Floating Point Arithmetic . . . . .	7
7. The implementation . . . . .	11
18. Constructors and destructors . . . . .	17
22. Basic functions . . . . .	20
31. The arithmetic operations . . . . .	24
50. Sign determination . . . . .	33
65. The functions improve and related functions . . . . .	39

## 1. Introduction

The need for exact real arithmetic or at least higher precision arithmetic than standard **double** precision arithmetic arises frequently. One prime example is computational geometry. Other examples are computational metrology and constrained logic programming; we refer the reader to [YD95] for an extensive discussion. Almost all algorithms in computational geometry are formulated under the assumption of exact real arithmetic. When implemented naively using floating point arithmetic these algorithms tend to be unreliable, i.e., to either crash or produce wrong results at least on some inputs. Two approaches have been taken to remedy this situation. The first approach sticks with floating point arithmetic and reformulates the algorithms. This approach has been successful for a small number of problems only, for example, two-dimensional convex hulls [LM92], line arrangements [FM91] and Delaunay triangulations [For92]. The other approach uses exact real arithmetic. Although it has been known for a long time [Mig92, Yap94] that exact computation with algebraic numbers<sup>1</sup> is possible, it has also been felt that the overhead would be prohibitive. Serious experimental studies appeared only recently. In [KLN91, Fv93, BJMM93, MN94b] it is shown that the overhead for exact rational arithmetic is small in the context of geometric computations. In particular, floating point filters as introduced in [Fv93] are very effective. A floating point filter first approximates a computation using **doubles**. Only if a test cannot be performed conclusively on the basis of the **double** approximation, exact arithmetic is used as a back-up. In [BMS94] the same idea is generalized to a computation involving non-rational algebraic numbers, and in [Yap93, YD95] it is argued that it can be applied to all numerical computations arising in computational geometry. In this paper we define and implement the LEDA data type **real**. Any integer is a **real** and the class of **reals** is closed under addition, subtraction, multiplication, division and squareroot. The full repertoire of relational operators ( $=, \neq, <, \leq, >, \geq,$ ) is available on **reals**. The outcome is *exact*. In addition, arbitrary floating point approximations of **real** numbers can be computed. We define the class **real** in section 2 and also give two examples of the use of **real** numbers: an incircle test arising in the computation of Voronoi diagrams of line segments and the computation of line segment intersection. Starting in section 7 we give the implementation of **reals**. Our implementation is similar to the proposal by Yap and Dube [YD95]. In section 3 we review some facts about floating point computation.

---

<sup>1</sup>A real number is algebraic if it is the root of a polynomial with integer coefficients.

## 2. The manual pages of data type real

### 1. Definition

An instance  $x$  of the data type real is an algebraic real number. There are many ways to construct a real: either by conversion from *double*, *bigfloat*, *integer* or *rational* or by applying one of the arithmetic operators  $+$ ,  $-$ ,  $*$ ,  $/$  or  $\sqrt{\quad}$  to real numbers. One may test the sign of a real number or compare two real numbers by any of the comparison relations  $=$ ,  $\neq$ ,  $<$ ,  $\leq$ ,  $>$  and  $\geq$ . The outcome of such a test is *exact*. There is also a non-standard version of the sign function: the call  $x.sign(integer\ q)$  computes the sign of  $x$  under the precondition that  $|x| \leq 2^{-q}$  implies  $x = 0$ . This version of the sign function allows the user to assist the data type in the computation of the sign of  $x$ , see the example below.

There are several functions to compute approximations of reals. The calls  $tobigfloat(x)$  and  $x.get\_bigfloat\_error()$  return *bigfloats*  $xnum$  and  $xerr$  such that  $|xnum - x| \leq xerr$ . The user may set a bound on  $xerr$ . More precisely, after the call  $x.improve\_approximation\_to(integer\ q)$  the data type guarantees  $xerr \leq 2^{-q}$ . One can also ask for *double* approximations of a real number  $x$ . The calls  $todouble(x)$  and  $x.get\_double\_error()$  return *doubles*  $xnum$  and  $xerr$  such that  $|xnum - x| \leq |xnum| * xerr$ . Note that  $xerr = \infty$  is possible.

### 2. Creation

reals may be constructed from data types *double*, *bigfloat*, *long*, *int* and *integer*. The default constructor  $real()$  initializes the real to zero.

### 3. Operations

*double*  $todouble(real)$  returns the current double approximation of  $x$ .

*bigfloat*  $tobigfloat(real)$  returns the current bigfloat approximation of  $x$ .

*double*  $x.get\_double\_error()$   
returns the relative error of the current double approximation of  $x$ , i.e.,  $|x - todouble(x)| \leq x.get\_double\_error() \cdot |todouble(x)|$ .

*bigfloat*  $x.get\_bigfloat\_error()$   
returns the absolute error of the current bigfloat approximation of  $x$ , i.e.,  $|x - tobigfloat(x)| \leq x.get\_bigfloat\_error()$ .

*int*  $sign(real\ x)$  returns the sign of (the exact value of)  $x$ .

*int*  $x.sign(integer\ q)$  as above. *Precondition:* if  $|x| \leq 2^{-q}$  then  $x = 0$ .

<i>void</i>	<i>x.improve_approximation_to(integer q)</i>	(re-)computes the approximation of <i>x</i> such that <i>x.get_bigfloat_error()</i> $\leq 2^{-q}$ after the call <i>x.improve_approximation_to(q)</i> .
<i>void</i>	<i>x.compute_with_precision(long k)</i>	(re-)computes the approximation of <i>x</i> ; each numerical operation is carried out with <i>k</i> binary places.
<i>void</i>	<i>x.guarantee_relative_error(long k)</i>	(re-)computes an approximation of <i>x</i> such that the relative error of the approximation is less than $2^{-k}$ .
<i>ostream&amp; ostream&amp; O</i>	$\ll$ <i>x</i>	writes the closest interval that is known to contain <i>x</i> to the output stream <i>O</i>
<i>istream&amp; istream&amp; I</i>	$\gg$ <i>real&amp; x</i>	reads <i>real</i> number <i>x</i> from the output stream <i>I</i> (in <i>double</i> format)
<i>real</i>	<i>abs(real x)</i>	absolute value of <i>x</i>
<i>real</i>	<i>sq(real x)</i>	square of <i>x</i>
<i>real</i>	<i>dist(real x, real y)</i>	euclidean distance of vector (x,y) to the origin
<i>real</i>	<i>powi(real x, int n)</i>	<i>n</i> .th power of <i>x</i>

#### 4. Implementation

A real is represented by the expression which defines it and a *double* approximation  $\hat{x}$  together with a relative error bound  $\epsilon_x$ . The arithmetic operators  $+$ ,  $-$ ,  $*$ ,  $\sqrt{\quad}$  take constant time. When the sign of a real number needs to be determined, the data type first computes a number  $q$ , if not already given as an argument to *sign*, such that  $|x| \leq 2^{-q}$  implies  $x = 0$ . The bound  $q$  is computed as described in [Mig92]. Using *bigfloat* arithmetic, the data type then computes an interval  $I$  of maximal length  $2^{-q}$  that contains  $x$ . If  $I$  contains zero, then  $x$  itself is equal to zero. Otherwise, the sign of any point in  $I$  is returned as the sign of  $x$ .

Two shortcuts are used to speed up the computation of the sign. Firstly, if the *double* approximation already suffices to determine the sign, then no *bigfloat* approximation is computed at all. Secondly, the *bigfloat* approximation is first computed only with small

precision. The precision is then doubled until either the sign can be decided (i.e., if the current approximation interval does not contain zero) or the full precision  $2^{-q}$  is reached.

## 5. Example

We give two examples of the use of the data type real. The examples deal with the Voronoi diagram of line segments and the intersection of line segments, respectively.

The following incircle test arises in the computation of Voronoi diagrams of line segments [BMS94]. For  $i$ ,  $1 \leq i \leq 3$ , let  $l_i : a_i x + b_i y + c_i = 0$  be a line in two-dimensional space and let  $p = (0, 0)$  be the origin. In general, there are two circles passing through  $p$  and touching  $l_1$  and  $l_2$ . The centers of these circles have homogeneous coordinates  $(x_v, y_v, z_v)$ , where

$$\begin{aligned} x_v &= a_1 c_2 + a_2 c_1 \pm \text{sign}(s) \sqrt{2c_1 c_2 (\sqrt{N} + D)} \\ y_v &= b_1 c_2 + b_2 c_1 \pm \text{sign}(r) \sqrt{2c_1 c_2 (\sqrt{N} - D)} \\ z_v &= \sqrt{N} - a_1 a_2 - b_1 b_2 \end{aligned}$$

and

$$\begin{aligned} s &= b_1 |a_2| - b_2 |a_1|, & N &= (a_1^2 + b_1^2)(a_2^2 + b_2^2) \\ r &= a_1 |b_2| - a_2 |b_1|, & D &= a_1 a_2 - b_1 b_2. \end{aligned}$$

Let us concentrate on one of these (say, we take the plus sign in both cases). The test whether  $l_3$  intersects, touches or misses the circle amounts to determining the sign of

$$E := \text{dist}^2(v, l_3) - \text{dist}^2(v, p) = \frac{(a_3 x_v + b_3 y_v + c_3)^2}{a_3^2 + b_3^2} - (x_v^2 + y_v^2)$$

We show how to compute the sign of  $\tilde{E} := (a_3^2 + b_3^2) \cdot E$  using our data type real. We assume that all coefficients  $a_i, b_i$  and  $c_i$ ,  $1 \leq i \leq 3$ , are  $k$  bit integers, i.e., integers whose absolute value is bounded by  $2^k - 1$ . In [BMS94] we showed that for  $\tilde{E} \neq 0$  we have  $|\tilde{E}| \geq 2^{-24k-26}$ . These observations suggest the following program.

```

int INCIRCLE(int k, real a1, real b1, real c1, real a2, real b2, real c2, real a3, real
b3, real c3 )
{
  real RN = sqrt((a1 * a1 + b1 * b1) * (a2 * a2 + b2 * b2));
  real A = a1 * c2 + a2 * c1;
  real B = b1 * c2 + b2 * c1;
  real C = 2 * c1 * c2;
  real D = a1 * a2 - b1 * b2;
  real s = b1 * fabs(a2) - b2 * fabs(a1);
  real r = a1 * fabs(b2) - a2 * fabs(b1);
  int signx = sign(s);
  int signy = sign(r);

```

```

real  $x_v = A + \text{sign}_x * \text{sqrt}(C * (RN + D));$ 
real  $y_v = B - \text{sign}_y * \text{sqrt}(C * (RN - D));$ 
real  $z_v = RN - (a_1 * a_2 + b_1 * b_2);$ 
real  $P = a_3 * x_v + b_3 * y_v + c_3 * z_v;$ 
real  $D_3^2 = a_3 * a_3 + b_3 * b_3;$ 
real  $R^2 = x_v * x_v + y_v * y_v;$ 
real  $E = P * P - D_3^2 * R^2;$ 
return  $E.\text{sign}(24 * k + 26);$ 
}

```

We turn to the line segment intersection problem next. Assume that all endpoints have  $k$ -bit integer homogeneous coordinates. This implies that the intersection points have homogeneous coordinates  $(X, Y, W)$  where  $X, Y$  and  $W$  are  $(4k + 3)$ -bit integers. The Bentley–Ottmann plane sweep algorithm for segment intersection [MN94a] needs to sort points by their  $x$ -coordinates, i.e., to compare fractions  $X_1/W_1$  and  $X_2/W_2$  where  $X_1, X_2, W_1, W_2$  are as above. This is tantamount to determining the sign of the  $8k + 7$  bit integer  $X_1 * W_2 - X_2 * W_1$ . If all variables  $X_i, W_i$  are declared **real** then their sign test will be performed quite efficiently. First, a *double* approximation is computed and then, if necessary, *bigfloat* approximations of increasing precision. In many cases, the *double* approximation already determines the sign. In this way, the user of the data type *real* gets the efficiency of a floating point filter [Fv93, MN94b] without any work on his side. This is in marked contrast to [Fv93, MN94b] and will be incorporated into [MN94a].

### 3. A Short Review of Floating Point Arithmetic

We review some facts about double precision floating point arithmetic according to the IEEE standard [IEE87]. A **double** is specified by a sign  $s \in \{0, 1\}$ , an exponent  $e \in [0 \dots 2047]$  and a binary fraction  $f = f_1 f_2 \dots f_{52}$  with  $f_i \in \{0, 1\}$  for all  $i, 1 \leq i \leq 52$ . We also use  $f$  to denote the number  $\sum_i f_i 2^{-i}$ . The number  $v$  represented by a triple  $(s, e, f)$  is defined as follows:

- if  $e = 2047$  and  $f \neq 0$  then no number is represented ( $v$  is a *NaN*)
- if  $e = 2047$  and  $f = 0$  then  $v = (-1)^s \cdot \infty$
- if  $0 < e < 2047$  then  $v = (-1)^s \cdot (1 + f) \cdot 2^{e-1023}$
- if  $e = 0$  and  $f \neq 0$  then  $v = (-1)^s \cdot f \cdot 2^{-1022}$  and  $v$  is called a denormalized number
- if  $e = 0$  and  $f = 0$  then  $v = (-1)^s \cdot 0$  (= signed zero)

Let  $MinDbl = 2^{-1022}$  and  $MaxDbl = (2 - 2^{-52}) \cdot 2^{1023}$ . We call a real number *approximable* if either  $a = 0$  or  $MinDbl \leq |a| < MaxDbl + 2^{-53} \cdot 2^{1023}$ . For any approximable number  $a$  let  $round(a)$  be the floating point number nearest to  $a$ . If there are two floating point numbers equally near to  $a$  then  $round(a)$  is the one where the least significant bit is zero. For any of the arithmetic operations  $+, -, *, /, \sqrt{\phantom{x}}$ , the IEEE standard guarantees the following. Let  $z$  be the *exact* result of an operation applied to floating point operands.

- if  $z$  is approximable then  $round(z)$  is delivered,
- if  $|z| \geq MaxDbl + 2^{-53} \cdot 2^{1023}$  then an overflow exception is signaled and  $\pm\infty$  is delivered where the sign is determined by  $z$ ,
- if  $0 < |z| < MinDbl$  then the result may be zero, a denormalized number or  $\pm MinDbl$  with the restriction that  $z$  and the delivered result do not have opposite signs<sup>2</sup>; in some cases the underflow is signaled,
- if any argument of an operation is infinity or not a number, then the result is guaranteed to be infinity or not a number.

We need a relation between an approximable number  $a$  and its floating point approximation  $round(a)$ . Assume  $a \neq 0$  and let  $e$  be such that  $2^e \leq |a| \leq 2^{e+1}$ . Then  $round(a) = a + \rho(a)$  where  $|\rho(a)| \leq 2^{-53} \cdot 2^e$  and  $2^e \leq |round(a)| < 2^{e+1}$ . Thus

$$\frac{|a - round(a)|}{|a|} \leq \frac{2^{-53} \cdot 2^e}{2^e} = 2^{-53}$$

and

$$\frac{|a - round(a)|}{|round(a)|} \leq \frac{2^{-53} \cdot 2^e}{2^e} = 2^{-53}.$$

The quantity  $2^{-53}$  is so important that it deserves a special name. We call it the *precision* of floating point arithmetic and use  $eps$  to denote it. By the two bounds derived above we have

---

<sup>2</sup>The standard is actually more specific but this is of no importance for this paper



$$\text{round}(a) \in [(1 - \text{eps}) \cdot a, (1 + \text{eps}) \cdot a] \subset [a/(1 + \text{eps}), a \cdot (1 + \text{eps})]$$

and

$$a \in [(1 - \text{eps}) \cdot \text{round}(a), (1 + \text{eps}) \cdot \text{round}(a)]$$

for any non-zero approximable number  $a$ .

We will frequently derive expressions for error bounds. These expressions involve the operators  $+$ ,  $*$ ,  $/$  and have non-negative operands. When such an expression is evaluated using floating point arithmetic, further error incurs which must be compensated by post-multiplying with a *correction factor*. We now show how to compute this factor. We first determine for any expression  $E$  a number  $\text{ind}(E)$  such that

$$\hat{E} \in [E/(1 + \text{eps})^{\text{ind}(E)}, E \cdot (1 + \text{eps})^{\text{ind}(E)}]$$

where  $\hat{E}$  is the value of  $E$  computed in floating point arithmetic. We assume first that no underflow or overflow occurs in the computation of  $\hat{E}$ . This implies that all intermediate results are approximable. If  $E$  is an operand then  $\text{ind}(E) = 0$ , if  $E = E_1 + E_2$  then ( $\oplus$  denotes floating point addition)

$$\hat{E} = \hat{E}_1 \oplus \hat{E}_2 \geq \frac{\hat{E}_1 + \hat{E}_2}{1 + \text{eps}} \geq \frac{E_1/(1 + \text{eps})^{\text{ind}(E_1)} + E_2/(1 + \text{eps})^{\text{ind}(E_2)}}{1 + \text{eps}}$$

and hence  $\text{ind}(E_1 + E_2) = 1 + \max(\text{ind}(E_1), \text{ind}(E_2))$ , and if  $E = E_1 \cdot E_2$  or  $E = E_1/E_2$  then

$$\hat{E} = \hat{E}_1 \odot \hat{E}_2 \geq \frac{\hat{E}_1 \cdot \hat{E}_2}{1 + \text{eps}} \geq \frac{E_1 \cdot E_2}{(1 + \text{eps})^{1 + \text{ind}(E_1) + \text{ind}(E_2)}}$$

and hence  $\text{ind}(E_1 \cdot E_2) = 1 + \text{ind}(E_1) + \text{ind}(E_2)$  and analogously  $\text{ind}(E_1/E_2) = 1 + \text{ind}(E_1) + \text{ind}(E_2)$ . The upper bound for  $\hat{E}$  can be shown by similar reasoning. We now choose our correction factor as  $1 + k \cdot \text{eps}$  where  $k$  is a power of 2 and  $(1 + \text{eps})^{\text{ind}(E)+1} \leq 1 + k \cdot \text{eps}$ . Then

$$E \leq (1 + \text{eps})^{\text{ind}(E)} \cdot \hat{E} \leq \frac{(1 + k \cdot \text{eps}) \cdot \hat{E}}{1 + \text{eps}} \leq (1 \oplus k \odot \text{eps}) \odot \hat{E},$$

where the last inequality follows from the fact that  $1 \oplus k \odot \text{eps} = 1 + k \cdot \text{eps}$  since  $k$  is a power of two. For small values of  $\text{ind}(E)$  one can choose  $k$  to be the smallest power of two with  $\text{ind}(E) + 2 \leq k$ .

What happens when there is overflow or underflow in the evaluation of  $E$ ? We detect overflow by inspecting the final result; it must be infinity or not a number. So assume that the evaluation of some subexpression  $E_s$  underflows. Let us also assume that there is no division operator on the path from the subexpression  $E_s$  to the root of the expression  $E$  (in our error bounds this will be the case for any subexpression that can underflow). Then we may conceptually write  $E$  as  $E = E_{a(s)} \cdot E_s + E_{b(s)}$  since any expression over  $+$  and  $*$  is a linear function of each of its subexpressions. Let  $M_{a(s)}$  be a bound on  $E_{a(s)}$ . The value of  $\hat{E}_s$  is less than  $\text{MinDbl}$  and due to underflow it may be set to zero. The underflow of  $E_s$  can therefore decrease the value of  $E$  by no more than  $\text{MinDbl} \cdot M_{a(s)}$ . Summing over all subexpressions we conclude that underflow can decrease the value of  $E$

by no more than  $MinDbl \cdot \sum_s M_{a(s)}$  where  $s$  ranges over all subexpressions of  $E$ . This implies that we have  $E \leq (1 + eps)^{ind(E)} \cdot \hat{E} + MinDbl \cdot \sum_s M_{a(s)}$  where  $\hat{E}$  is the computed value (now with the possibility of underflow). In our applications, we will always have  $\hat{E} \geq eps$ . In that case, if we guarantee that  $MinDbl \cdot (\sum_s M_{a(s)}/eps) \leq eps/2$  and choose  $k$  such that  $(1 + eps)^{ind(E)+1} + (1 + eps) \cdot eps/2 \leq 1 + k \cdot eps$  then

$$\begin{aligned}
E &\leq (1 + eps)^{ind(E)} \cdot \hat{E} + MinDbl \cdot \sum_s M_{a(s)} \\
&\leq ((1 + eps)^{ind(E)} + (MinDbl \cdot \sum_s M_{a(s)})/eps) \cdot \hat{E} \\
&\leq ((1 + eps)^{ind(E)} + eps/2) \cdot \hat{E} \\
&\leq (1 + k \cdot eps) \cdot \hat{E} / (1 + eps) \leq (1 \oplus k \odot eps) \odot \hat{E}
\end{aligned}$$

Let us consider the example  $E = eps + |x/z| \cdot \epsilon_x + |y/z| \cdot \epsilon_y$ . The expression  $E$  has index 4 and the two multiplications and the two divisions may underflow. The sum of the multipliers of the “dangerous” subexpressions is  $M := \epsilon_x + \epsilon_y + 1 + 1$ . If  $\epsilon_x$  and  $\epsilon_y$  are both bounded by  $2^{915}$  then clearly  $MinDbl \cdot M \leq eps^2/2$  and hence we may use the correction factor  $1 + 8 \cdot eps$ . In fact, *the factor  $1 + 8 \cdot eps$  suffices for all bounds in this paper provided that incoming error bounds are bounded by  $2^{915}$* . We define both quantities as global constants.

`<Global functions and constants 3 > ≡`

```

double pow2(long);
const double eps = pow2(-53);
const double correction = 1 + 8 * eps;
const double MaxError = pow2(915);
const double MinDbl = pow2(-1022);
const double twoMinDbl = pow2(-1021);
const double MaxDbl = (2 - pow2(-52)) * pow2(1023);
const double Infinity = pow2(1024);

```

See also chunks 5, 6, and 29.

This code is used in chunk 13.

#### 4. Bigfloats.

LEDA’s **bigfloats** are a generalization of **doubles**. The differences are:

- The exponent range is the set of all integers (type **integer**) and hence there is no underflow and overflow. This implies that all numbers are approximable and that  $round(a) = 0$  implies  $a = 0$ .
- The mantissa length can be prescribed by the user. There are two ways to do this. One way is to call the function **bigfloat**::*set\_glob\_prec*(**int**  $p$ ) to set the mantissa length to  $p$ . The other way is to call *add*( $x, y, p$ ) for the addition of **bigfloats**  $x$  and  $y$  with mantissa length  $p$ , and analogously *sub*, *mul*, *div* for subtraction, multiplication and division, respectively<sup>3</sup>.

---

<sup>3</sup>The motivation for this syntax is to be able to evaluate different operations in a large expression with different precisions.

5. We need functions to decide whether **integers** and **bigfloats** are exactly representable by **doubles**.

```
< Global functions and constants 3 > +≡  
  int isdouble(const integer &x)  
  { return (x ≡ integer(todouble(bigfloat(x)))); }  
  int isdouble(const bigfloat &x)  
  { return (x ≡ bigfloat(todouble(x))); }
```

6. Also we use a constant for zero as an **integer**.

```
< Global functions and constants 3 > +≡  
  const integer zero_integer(0);
```

## 7. The implementation

The implementation defines the files `real.h` and `real.c`. The first file contains the interface of class `real` and must be included into any program using `reals`. The latter file can be compiled into `real.o`. Any application of `reals` must specify `real.o` in the linking phase.

The public functions of class `real` are exactly the ones defined in the manual.

```
<real.h 7> ≡
#ifndef REAL_H
#define REAL_H
#include "/KM/usr/burnikel/cweb/bigfloat/new/bigfloat.h"
#include <LEDA/rational.h>
class real_rep;
class real
{
  <private data and functions of type real 15>
public:
  real();
  real(double);
  real(int);
  real(const integer &);
  real(const bigfloat &);
  real(const rational &);
  real(const real &);
  real(real_rep &);
  ~real();

  real &operator=(const real &);
  friend double todouble(const real &);
  inline double todouble() const { return ::todouble(*this); }
  friend bigfloat tobigfloat(const real &);
  inline bigfloat tobigfloat() const { return ::tobigfloat(*this); }
  double get_double_error() const;
  bigfloat get_bigfloat_error() const;
  inline friend int sign(const real &x) { return x.sign(); }
  int sign() const;
  int sign(const integer &q) const;
  int sign(long) const;
  void improve_approximation_to(const integer &q) const;
  void compute_with_precision(long k) const;
  void guarantee_relative_error(long k) const;
  friend ostream &operator<<(ostream &O, const real &x);
  friend istream &operator>>(istream &I, real &x);
  <Comparison operators 8>
  <Arithmetic operators 9>
};
```

```

    < LEDA functions 11 >
    < Assignment operators 10 >
    < Declaration of mathematical functions in real.h 12 >
#endif

```

8. We have the full list of comparison operators.

```

< Comparison operators 8 > ≡
friend int operator < (const real &, const real &);
friend int operator ≤ (const real &, const real &);
friend int operator > (const real &, const real &);
friend int operator ≥ (const real &, const real &);
friend int operator ≡ (const real &, const real &);
friend int operator ≠ (const real &, const real &);

```

This code is used in chunks 7 and 17.

9. All arithmetic operators are friends of classes **real** and **real\_rep**.

```

< Arithmetic operators 9 > ≡
friend real operator + (const real &, const real &);
friend real operator - (const real &, const real &);
friend real operator * (const real &, const real &);
friend real operator / (const real &, const real &);
friend real operator - (const real &);
friend real sqrt (const real &);

```

This code is used in chunks 7 and 17.

10. We also have some fancy assignment operators.

```

< Assignment operators 10 > ≡
inline real operator += (real &x, const real &y)
{ x = x + y; return x; }
inline real operator -= (real &x, const real &y)
{ x = x - y; return x; }
inline real operator *= (real &x, const real &y)
{ x = x * y; return x; }

```

This code is used in chunk 7.

11. The following functions have to be defined for every LEDA type.

```

< LEDA functions 11 > ≡
inline void Print (const real &x, ostream &out) { out << x; }
inline void Read (real &x, istream &in) { in >> x; }
inline int compare (const real &x, const real &y) { return (x - y).sign(); }
inline char *Type_Name (const real *) { return "real"; }

```

This code is used in chunk 7.

12. Mathematical functions. This list is still small.

```
<Declaration of mathematical functions in real.h 12> ≡  
    real abs(const real &x);  
    real sq(const real &x);  
    real dist(const real &x, const real &y);  
    real powi(const real &x, int n);
```

This code is used in chunk 7.

13. The file `real.c` has a similar structure as the file `real.h`. For each group of functions there is a section containing the implementation. In addition, we define some constants and the class `real_rep`.

```
<real.c 13> ≡  
#include "real.h"  
#include <math.h>  
    <Global functions and constants 3>  
    <Basic declarations 16>  
    <Declaration of class real_rep 17>  
    <Constructors and destructors 18>  
    <Basic functions 23>  
    <Operators 32>  
    <Sign functions 50>  
    <Functions for improving the approximation 65>  
    <Implementation of mathematical functions in real.c 30>
```

14. Realisation of class `real`.

A **real** may be exactly representable by a **double**. In this case, we call it a **basic real**. Otherwise, if a **real** is not (known to be) a **double**, it is defined by an expression tree whose leaves are given by basic **reals**. In that case we maintain the expression tree and the real only needs to know a pointer to the root of the tree. The nodes of expression trees are realized by the class `real_rep` which is the work horse of our implementation. The member functions and all operations, except the arithmetic operators, are first implemented on `real_reps` and then lifted to **reals**. The design to separate an object of type **real** from its tree representation has the following advantages:

- Copying a **real** takes constant time since it basically amounts to copying a pointer to a `real_rep`.
- Since reals may share subtrees, we have to provide our own memory management. The part of the information that can be automatically allocated is stored in class **real** itself, and the dynamically allocated part is stored in class `real_rep`.

15. The private data of class **real** is organized as follows.

- If `is_double` is true then the **real** is basic and the value is stored in the member `value`. If the **real** is part of an expression tree then `PTR` is defined and points to a `real_rep` also containing the value. Storing the value twice simplifies programming; of course, it also takes extra space.

- If *is\_double* is *false*, *PTR* points to the root of an expression tree defining the **real**; the content of *value* is meaningless in this case.

```

⟨private data and functions of type real 15⟩ ≡
  bool is_double;
  double value;
  real_rep *PTR;

```

See also chunk 25.

This code is used in chunk 7.

## 16. The class **real\_rep**.

A **real\_rep** object is either a leaf, i.e., represents a **double** number or a **bigfloat** number, or it is constructed by an arithmetic operation from two other **real\_rep** objects. We define the enumeration type *constructortype* to distinguish these cases.

A **real\_rep** may contain a **bigfloat** approximation to the exact value of the represented number. Since this approximation is space-consuming and not always needed we encapsulate it in a structure called *bf\_app*. The structure is organized as follows.

- *NUM\_BF* is the **bigfloat** approximation,
- *ERROR\_BF* bounds the error of the approximation *NUM\_BF*, i.e.,  $|NUM\_BF - NUM\_EXACT| \leq ERROR\_BF$ .

```

⟨Basic declarations 16⟩ ≡
  enum constructortype {
    DOUBLE = 0, BIGFLOAT = 1, NEGATION = 2, SQUAREROOT = 3,
    ADDITION = 4, SUBTRACTION = 5, MULTIPLICATION = 6,
    DIVISION = 7
  };
  struct app_bf {
    bigfloat NUM_BF;
    bigfloat ERROR_BF;
    app_bf() {}
    app_bf(const bigfloat &x, const bigfloat &error_x = pZero)
    { NUM_BF = x; ERROR_BF = error_x; }
  };

```

See also chunks 52 and 83.

This code is used in chunk 13.

17. We can now define the class **real\_rep**. Any object of class **real\_rep** contains a **double** approximation *NUM*, a **double** *ERROR*, a pointer *APP\_BF\_PTR* to a **bigfloat** approximation (of type **app\_bf**), the **constructortype** *CON*, pointers *OP1* and *OP2* to the two operands, and a *count* of the number of pointers to the object. Furthermore there is a long member *status* that we use for efficient searching through the expression trees and a pointer *sep\_bound\_ptr* to a structure of type *sep\_bound* that we define in the section about sign computation.

We will always use *NUM\_EXACT* to denote the exact value of a **real\_rep**. The quantities *NUM*, *ERROR*, *NUM\_BF* and *ERROR\_BF* satisfy the following inequalities: Either *ERROR* = ∞ or *ERROR* ≤ *MaxError*, *NUM* = 0 or *MinDbl* ≤ *NUM* ≤ *MaxDbl* (i.e., *NUM* is not a denormalized number) and

$$|NUM\_EXACT - NUM| \leq |NUM| \cdot ERROR.$$

If *APP\_BF\_PTR* ≠ *nil* then

$$|NUM\_EXACT - NUM\_BF| \leq ERROR\_BF < \infty.$$

Note that *ERROR* is the *relative* error of the **double** approximation and *ERROR\_BF* is the *absolute* error of the **bigfloat** approximation. Also note that the **bigfloat** approximation may not exist (if *APP\_BF\_PTR* ≡ *nil*) and that some numbers have no **double** approximation (if *ERROR* ≡ ∞). However, if the **bigfloat** approximation of a **real\_rep** is defined, then it is also defined for its children.

If a **real** has a **double** and a **bigfloat** approximation then the **double** approximation is *essentially* the conversion of the **bigfloat** value to **double**. Since our algorithms do not exploit any connection between the **double** and the **real** approximation we formalize no connection.

The member functions of class **real\_rep** do the following:

- *init\_app\_bf*( ) initializes the **bigfloat** approximation,
- *adjust\_dbl*( ) adjusts the **double** approximation to an improved **bigfloat** approximation,
- *compute\_op*(**int** *p*) computes the last operation (given by member *CON*) in precision *p*, using bigfloat arithmetic,
- *compute\_approximation* implements the function **real**::*compute\_with\_precision*,
- *sign\_with\_separation\_bound*(**int** *k*, **bool** *supported*) returns the sign of the represented number *x*; if the boolean parameter *supported* is true, the sign computation is supported by the information that  $|x| \leq 2^k$  implies  $x = 0$ ,
- *guarantee\_bound\_two\_to*(**integer** *e*) improves the **bigfloat** approximation of the number such that *ERROR\_BF* is bounded by  $2^e$ ,
- *num\_bf*( ) and *error\_bf*( ) return references to the values *APP\_BF\_PTR*→*NUM\_BF* and *APP\_BF\_PTR*→*ERROR\_BF*,
- *num\_exp*( ) and *err\_exp*( ) compute the binary logarithms of nonzero *NUM\_BF* and *ERROR\_BF*, rounded up towards the next integer,
- *exact*( ) tells whether *ERROR\_BF* is zero.

⟨ Declaration of class **real\_rep** 17 ⟩ ≡  
**struct sep\_bound;**  
**class real\_rep {**  
**friend class real;**



```

double NUM;
double ERROR;
app_bf *APP_BF_PTR;
constructortype CON;
real_rep *OP1, *OP2;
int count;
long status;
sep_bound *sep_bound_ptr;
real_rep();
real_rep(double);
real_rep(const bigfloat &);
~real_rep();
LEDA_MEMORY(real_rep)
void init_app_bf();
friend double todouble(const real &);
friend bigfloat tobigfloat(const real &);
void adjust_dbl();
void compute_op(long);
void compute_approximation(long);
int sign_with_separation_bound (const integer &, bool = false );
void guarantee_bound_two_to(const integer &);
inline bigfloat &num_bf() { return APP_BF_PTR->NUM_BF; }
inline bigfloat &error_bf() { return APP_BF_PTR->ERROR_BF; }
inline bool exact() { return isZero(error_bf()); }
inline integer num_exp() { return log2(num_bf()); }
inline integer err_exp() { return log2(error_bf()); }
< Local functions of class real_rep 53 >
friend ostream &operator<<(ostream &, const real &);
< Arithmetic operators 9 >
< Comparison operators 8 >
};

```

This code is used in chunk 13.

## 18. Constructors and destructors

We start with the constructors for **class real**. The first three construct basic **reals** from an empty argument list, from a **double**, and from an **int**. They are easy to implement. The next two – for **integer** and **bigfloat** – construct **reals** that in general are not representable by **doubles** but are defined by a single **real\_rep** node. Rationals are constructed differently. We first construct **reals** for numerator and denominator, respectively, and then return their **real** division. Hence the result is a pointer to a tree with three nodes unless the rational is a **double**. The constructor that takes a **real\_rep** argument simply assigns the address of the **real\_rep** to **PTR** and sets *is\_double*  $\equiv$  *false*.

⟨ Constructors and destructors 18 ⟩  $\equiv$

```
real::real() { is_double = true; value = 0; PTR = nil; }
real::real(int x) { is_double = true; value = x; PTR = nil; }
real::real(double x)
{ is_double = true; value = x; PTR = nil; }
real::real(const bigfloat &x)
{
  if (is_double = isdouble(x)) {
    value = ::todouble(x);
    PTR = nil;
  }
  else PTR = new real_rep (x);
}
real::real(const integer &x)
{
  if (is_double = isdouble(x)) {
    value = x.todouble();
    PTR = nil;
  }
  else PTR = new real_rep (bigfloat(x));
}
real::real(const rational &x)
{
  real OP1 = *new real_rep (x.numerator());
  real OP2 = *new real_rep (x.denominator());
  *this = OP1/OP2;
}
real::real(real_rep &x_rep) { is_double = false; PTR = &x_rep; }
```

See also chunks 19, 20, and 21.

This code is used in chunk 13.

**19.** Copy constructor and assignment operator are treated in direct succession because they are closely related. Both functions distinguish the two possible cases of *is\_double*. If the argument is not basic, the copy constructor simply copies the pointer to the **real\_rep** of its arguments and records the fact that the number of **reals** sharing the representation

has increased by one. The assignment operator does basically the same, but is difficult to implement correctly. Notice that we must take care that

1. an existing pointer  $x.PTR \neq nil$  is *not* copied for  $x.is\_double \equiv true$
2. an existing reference  $PTR$  is correctly deleted
3. everything remains correct even for  $x.PTR \equiv PTR$  or even  $*this = x$

In addition, if  $PTR$  pointed to some **real\_rep** node, it decreases the count of this node by one and deletes the node if the count goes to zero.

⟨Constructors and destructors 18⟩ +≡

```

real::real(const real &x)
{
  if (is_double = x.is_double) { value = x.value; PTR = nil; }
  else { x.PTR→count++; PTR = x.PTR;
  }
}

real &real::operator=(const real &x)
{
  /* we avoid problems by ruling out the assignment  $x = x$ ; */
  if (this ≡ &x) return (*this);    /* first the actions that must anyway be
      performed, in the respective cases of is_double */
  if (is_double = x.is_double) value = x.value;
  else x.PTR→count++;    /* second, an existing reference in PTR must always be
      deleted; if  $x.PTR \equiv PTR \neq nil$  and  $is\_double \equiv false$ , this leaves the number
      of references to *PTR unchanged, because we increased x.PTR before, and
      for  $is\_double \equiv true$  we in fact lose one reference because PTR is set to nil
      at the end */
  if (PTR) if (--PTR→count ≡ 0) delete PTR;    /* last, we set this→PTR */
  if (is_double) PTR = nil; else PTR = x.PTR;
  return (*this);
}

```

20. Beside the constructors for class **real** we need three constructors for class **real\_rep**. The first one is the straightforward empty constructor. The second one takes **double**s and is only needed when one of the **real** operators needs to initialize a leaf representing a basic **real**. The third constructor takes **bigfloat**s and is used by the corresponding **real** :: constructor. Note that even in this case we have to initialize the **double** approximation by calling the function *adjust\_dbl()*. It is defined in section 24.

⟨Constructors and destructors 18⟩ +≡

```

real_rep::real_rep()
{ APP_BF_PTR = nil; sep_bound_ptr = nil; count = 1; status = 0; }

real_rep::real_rep(double x)
{ APP_BF_PTR = nil; sep_bound_ptr = nil;
  OP1 = OP2 = nil; CON = DOUBLE;
  NUM = x; ERROR = 0;
  count = 1; status = 0;
}

```

```

real_rep::real_rep(const bigfloat &x)
{
    APP_BF_PTR = new app_bf (x,0); sep_bound_ptr = nil;
    OP1 = OP2 = nil; CON = BIGFLOAT;
    adjust_dbl();
    count = 1; status = 0;
}

```

21. We discuss the destructors of classes **real** and **real\_rep**. When a **real** is destroyed we simply decrease the counter of the corresponding **real\_rep** (if it exists). When the counter reaches zero we also delete the **real\_rep**. In order to delete a **real\_rep** we free the storage allocated for the **real\_rep** (this is done automatically except for the storage allocated for the optional **bigfloat** approximation and the memory addressed by *sep\_bound\_ptr*) and decrease the counters of the operands (if any) defining the **real\_rep**. If any counter reaches zero we destroy the operands recursively.

```

⟨ Constructors and destructors 18 ⟩ +≡
real::~real()
{ if (PTR) if (--PTR→count ≡ 0) delete PTR;
}

real_rep::~real_rep()
{
    if (APP_BF_PTR) delete APP_BF_PTR;
    if (sep_bound_ptr) delete sep_bound_ptr;
    if (OP1 ∧ (--OP1→count ≡ 0)) delete OP1;
    if (OP2 ∧ (--OP2→count ≡ 0)) delete OP2;
#ifdef DEBUG_REAL
    APP_BF_PTR = nil;
    OP1 = OP2 = nil;
#endif
}

```

## 22. Basic functions

Class **real\_rep** has the basic functions *init\_app\_bf*, *adjust\_dbl*, *num\_exp* and *err\_exp*, and the basic functions of class **real** are *todouble*, *get\_dbl\_error*, *tobigfloat*, *get\_precision*, and the stream operators  $\ll$  and  $\gg$ . We start with the basic functions of class **real\_rep**.

23. The function *init\_app\_bf*() initializes the **bigfloat** approximation of a **real\_rep** (if not defined already). If the **double** approximation is meaningful then we simply convert it to **bigfloat** format. Otherwise the functions calls itself recursively for the operands and computes the operation given by *CON* in precision 32.

```

⟨ Basic functions 23 ⟩ ≡
void real_rep::init_app_bf()
{
    if (APP_BF_PTR) return;
    if (OP1) OP1->init_app_bf();
    if (OP2) OP2->init_app_bf();
    if (ERROR < Infinity) {
        APP_BF_PTR = new app_bf (bigfloat (NUM),
            abs(bigfloat (ERROR * NUM * correction)));
        return;
    }
    APP_BF_PTR = new app_bf ();
    compute_op(53);
}

```

See also chunks 24, 26, 27, and 28.

This code is used in chunk 13.

24. The function *adjust\_dbl*() adjusts the **double** approximation (*NUM*, *ERROR*) to the current **bigfloat** approximation (which is guaranteed to exist when *adjust\_dbl* is called). The adjustment is non-trivial since the **bigfloat** error bound is absolute and the **double** error bound is relative and since we have to take care of underflow. We need to distinguish several cases, two of which are simple. If  $|NUM\_BF| > MaxDbl$  or  $ERROR\_BF > MaxDbl$  then we simply set *ERROR* to infinity and if  $NUM\_BF = ERROR\_BF = 0$  then we set  $NUM = ERROR = 0$ . So assume that neither is the case. Let  $n = NUM\_EXACT$ ,  $nbf = NUM\_BF$ ,  $ebf = ERROR\_BF$ ,  $\hat{e} = \mathbf{double}(ERROR\_BF)$ , and  $\hat{n} = \mathbf{double}(NUM\_BF)$ . Assume first that  $\hat{n} \neq 0$ . Then

$$\begin{aligned}
 |\hat{n} - n| &\leq |\hat{n} - nbf| + |nbf - n| \\
 &\leq eps \cdot |\hat{n}| + ebf \\
 &\leq (eps + ebf/|\hat{n}|) \cdot |\hat{n}| \\
 &\leq (eps + Max(MinDbl, (1 + eps) \cdot \hat{e})/|\hat{n}|) \cdot |\hat{n}|,
 \end{aligned}$$

since  $ebf \leq Max(MinDbl, (1 + eps)\hat{e})$ . The formula above has index 4 (cf. section 3) and hence  $1 + 8eps$  is a suitable correction factor for its floating point evaluation (this correction factor even allows to drop the multiplication of  $\hat{e}$  by  $(1 + eps)$ ).

Assume next that  $\hat{n} = 0$ . Then  $|n| \leq ebf + MinDbl \leq (1 + eps) \cdot \hat{e} + 2MinDbl$ . We set  $NUM = \hat{e} \oplus 2MinDbl$ . Then

$$|n - NUM| \leq 2(1 + eps)^2 NUM \leq 2 \odot (1 + 8eps) \odot NUM.$$

In either case, if *ERROR* exceeds *MaxError* (this will also be the case when the division  $|\hat{e}|/|\hat{n}|$  overflows) we invalidate our computation and set *ERROR* to infinity.

```

⟨Basic functions 23⟩ +≡
void real_rep::adjust_dbl()
{
    double n_head = todouble(num_bf());
    double e_head = todouble(error_bf());
    if (n_head ≠ 0) {
        NUM = n_head;
        ERROR = eps + Max(MinDbl, e_head)/fabs(n_head);
    }
    else {
        NUM = e_head + 2 * MinDbl;
        ERROR = 2;
    }
    ERROR = ERROR * correction;
    if ((ERROR ≥ MaxError) ∨ (fabs(n_head) ≡ Infinity)) {
        NUM = 1; ERROR = Infinity;
    }
}

```

**25.** We declare a private function **real**::*adjust* to adjust the state of a real. That is, we check whether a real is exactly a **double** and if so, we turn it into a basic **real**. In any case we adjust the **double** approximation. Not that we have to cast away the constancy of the **real** before changing its private data.

```

⟨private data and functions of type real 15⟩ +≡
void adjust() const;

```

```

26.⟨Basic functions 23⟩ +≡
void real::adjust() const
{
    if (is_double) return;
    if (PTR→exact() ∧ isdouble(PTR→num_bf())) {
        ((real &) *this).is_double = true;
        ((real &) *this).value = ::todouble(PTR→num_bf());
        return;
    }
    PTR→adjust_dbl();
}

```

27. We give out a **real** as an interval containing its exact value. The size of the interval depends on the quality of the **bigfloat** approximation.

We assume that **real** input reads from a **double** variable.

```

⟨Basic functions 23⟩ +≡
ostream &operator<<(ostream &out, const real &x)
{
    bigfloat low_bound, upp_bound;
    long prec;
    if (x.is_double) low_bound = upp_bound = x.value;
    else {
        x.PTR→init_app_bf();
        prec = (log2(x.PTR→num_bf()) - log2(x.PTR→error_bf())).tolong();
        if (prec < 53) prec = 53;
        low_bound = sub(x.tobigfloat(), x.PTR→error_bf(), prec, TO_N_INF);
        upp_bound = add(x.tobigfloat(), x.PTR→error_bf(), prec, TO_P_INF);
    }
    out << "[" << low_bound << ", ";
    out << upp_bound << "]";
    return out;
}

istream &operator>>(istream &in, real &x)
{ double x_num; in >> x_num;
  x = real(x_num); return in;
}

```

28. The other basic functions of class **real** do the obvious things.

```

⟨Basic functions 23⟩ +≡
double todouble(const real &x)
{
    if (x.is_double) return x.value;
    if (x.PTR→APP_BF_PTR) x.PTR→adjust_dbl();
    return x.PTR→NUM;
}

double real::get_double_error() const
{
    if (is_double) return 0;
    if (PTR→APP_BF_PTR) PTR→adjust_dbl();
    return PTR→ERROR;
}

bigfloat tobigfloat(const real &x)
{
    if (x.is_double) return x.value;
    else {
        x.PTR→init_app_bf();
        return x.PTR→num_bf();
    }
}

```

```

    }
}
bigfloat real::get_bigfloat_error() const
{
    if (is_double) return 0;
    else {
        PTR→init_app_bf();
        return PTR→error_bf();
    }
}

```

29. Now we list several auxiliary functions that we use in later sections. We need to compute

- $ldexp\_bf(x, k) = x \cdot 2^k$  for **bigfloats**  $x$  and integers  $k$ ,
- the maximum of two integers,
- the sign of a **double**.

```

⟨ Global functions and constants 3 ⟩ +≡
inline bigfloat ldexp_bf(const bigfloat &x, const integer &k)
{ return bigfloat(x.get_significant(), x.get_exponent() + k); }
integer Maximum(integer x, integer y) { return (x > y ? x : y); }
int sign(double x)
{ if (x ≡ 0) return 0; else return (x > 0 ? 1 : -1); }

```

30. Mathematical functions.

We implement some of the most frequently used functions from “math.h” for type **real**.

```

⟨ Implementation of mathematical functions in real.c 30 ⟩ ≡
real abs(const real &x) { return (x.sign() < 0 ? -x : x); }
real sq(const real &x) { return x * x; }
real dist(const real &x, const real &y) { return sqrt(x * x + y * y); }
real powi(const real &x, int n)
{
    real y = x, z = 1;
    int n_prefix = n;
    while (n_prefix > 0) {
        if (n_prefix % 2) z = z * y;
        n_prefix = n_prefix / 2;
        y = y * y;
    }
    return z;
}

```

This code is used in chunk 13.



### 31. The arithmetic operations

We give the implementation of the arithmetic operators  $+$ ,  $-$ ,  $*$ ,  $/$  and  $\sqrt{\quad}$ . The implementations have a common structure. If both operands are basic, we check if the **double** operation already gives an *exact* result which may then be returned immediately. We believe that the time for checking this is well invested because, if the check succeeds

1. it saves space,
2. it saves the time for computing error bounds and doing dynamic memory allocation,
3. it may detect “degenerate” values without expensive **bigfloat** arithmetic.

If one of the operands is not a **double**, or if the result is not exactly representable in **double** format, we create **real\_rep** nodes for the operands, if necessary, and also we create new node for the result  $z$ . There we record the type of the operation in the *CON* field, store pointers to the operands, increase the counters of the operands, and then compute the **double** approximation and the error bound for  $z$ .

Clearly, if either operand is not approximable (or no approximation is known), we cannot compute an approximation of  $z$ . So in our numerical analysis we will tacitly assume that both operands are approximable.

32. We start with the binary operations addition, subtraction, and multiplication.

⟨ Operators 32 ⟩ ≡

```
real operator+(const real &x,const real &y)
{
  ⟨if addition is exact return double result 36⟩
  ⟨create nodes for two operands 33⟩
  real_rep &z_rep = *new real_rep ();
  ⟨initialize z_rep as result of a binary operation 34⟩
  z_rep.CON = ADDITION;
  z_rep.NUM = x.PTR→NUM + y.PTR→NUM;
  ⟨compute error for addition or subtraction 38⟩
  ⟨if approximation is not valid set error to Infinity 35⟩
  return z_rep;
}

real operator-(const real &x,const real &y)
{
  ⟨if subtraction is exact return double result 37⟩
  ⟨create nodes for two operands 33⟩
  real_rep &z_rep = *new real_rep ();
  ⟨initialize z_rep as result of a binary operation 34⟩
  z_rep.CON = SUBTRACTION;
  z_rep.NUM = x.PTR→NUM - y.PTR→NUM;
  ⟨compute error for addition or subtraction 38⟩
  ⟨if approximation is not valid set error to Infinity 35⟩
  return z_rep;
}
```

```

real operator * (const real &x, const real &y)
{
  ⟨if multiplication is exact return double result 39⟩
  ⟨create nodes for two operands 33⟩
  real_rep &z_rep = *new real_rep ();
  ⟨initialize z_rep as result of a binary operation 34⟩
  z_rep.CON = MULTIPLICATION;
  z_rep.NUM = x.PTR→NUM * y.PTR→NUM;
  ⟨compute error for multiplication 40⟩
  ⟨if approximation is not valid set error to Infinity 35⟩
  return z_rep;
}

```

See also chunks 41, 45, and 49.

This code is used in chunk 13.

**33.** We give the common parts of the implementations first. First we create nodes for the operands  $x$  and  $y$ , if not already existing. We simply have to use the **real\_rep** constructor. Note that we have to cast away the constancy of  $x$  and  $y$  away before manipulating their data.

```

⟨create nodes for two operands 33⟩ ≡
if ( $\neg$ x.PTR) ((real &) x).PTR = new real_rep (x.value);
if ( $\neg$ y.PTR) ((real &) y).PTR = new real_rep (y.value);
double &x_num = x.PTR→NUM;
double &y_num = y.PTR→NUM;

```

This code is used in chunks 32 and 41.

**34.** We initialize the result  $z\_rep$ .

```

⟨initialize z_rep as result of a binary operation 34⟩ ≡
z_rep.OP1 = x.PTR; z_rep.OP2 = y.PTR;
z_rep.OP1→count++; z_rep.OP2→count++;
double &NUM = z_rep.NUM;

```

This code is used in chunks 32 and 41.

**35.** Checking the validity of the **double** approximation is done as follows. If one of the incoming error bounds is *Infinity*, our new error bound cannot be valid. The same holds if the computed approximation  $z\_rep.NUM$  is *NaN* or  $\pm\infty$ , or if the computed error bound  $z\_rep.NUM$  is larger than *MaxError* or *NaN*. Note that we do not need to check for the *NaN* cases separately because a comparison with a *NaN* always gives the answer “no”.

```

⟨if approximation is not valid set error to Infinity 35⟩ ≡
if ( $\neg$ ((fabs(NUM) + x.PTR→ERROR + y.PTR→ERROR < Infinity)  $\wedge$  (z_rep.ERROR <
  MaxError))) {
  z_rep.ERROR = Infinity;
  NUM = 1;
}

```

This code is used in chunks 32 and 41.

**36.** Now we skip to the addition-specific parts. It is possible to test efficiently whether a **double** addition  $s = a + b$  is exact. In fact, one can show that – irrespective of which rounding mode is active or whether the operands are normalized or denormalized –, one of the two subtractions  $(a \oplus b) \ominus b$  or  $(a \oplus b) \ominus a$  is exact. This implies immediately that the addition is exact if and only if both of the following equalities hold<sup>4</sup> (which is easy to test):

$$(a \oplus b) \ominus b = a \wedge (a \oplus b) \ominus a = b.$$

```

⟨ if addition is exact return double result 36 ⟩ ≡
  double NUM_result;
  if (x.is_double ∧ y.is_double) {
    NUM_result = x.value + y.value;
    if ((NUM_result - x.value ≡ y.value) ∧ ((NUM_result - y.value) ≡ x.value))
      return NUM_result;
  }

```

This code is used in chunk 32.

**37.** The same test procedure applies to subtraction.

```

⟨ if subtraction is exact return double result 37 ⟩ ≡
  double NUM_result;
  if (x.is_double ∧ y.is_double) {
    NUM_result = x.value - y.value;
    if ((NUM_result + y.value ≡ x.value) ∧ (x.value - NUM_result ≡ y.value))
      return NUM_result;
  }

```

This code is used in chunk 32.

**38.** Now we give the error analysis for addition and subtraction. Because they are essentially the same, we only deal with addition.

Here and in the analysis of the operations to follow we use  $x$  and  $y$  to denote the operands,  $\hat{x}$  and  $\hat{y}$  to denote the **double** approximations of  $x$  and  $y$  and  $\epsilon_x$  and  $\epsilon_y$  to denote the relative error of  $\hat{x}$  and  $\hat{y}$ . Thus  $|x - \hat{x}| \leq \epsilon_x \cdot |\hat{x}|$  and  $|y - \hat{y}| \leq \epsilon_y \cdot |\hat{y}|$ . Let  $\hat{z}$  denote the floating point approximation of the result of the floating point addition of  $\hat{x}$  and  $\hat{y}$ . Assume first that  $|\hat{z}| \geq \text{MinDbl}$ . Then

$$\begin{aligned}
|\hat{z} - (x + y)| &\leq |\hat{z} - (\hat{x} + \hat{y})| + |\hat{x} - x| + |\hat{y} - y| \\
&\leq |\hat{z}| \cdot eps + |\hat{x}| \cdot \epsilon_x + |\hat{y}| \cdot \epsilon_y \\
&\leq |\hat{z}| \cdot (eps + \left|\frac{\hat{x}}{\hat{z}}\right| \cdot \epsilon_x + \left|\frac{\hat{y}}{\hat{z}}\right| \cdot \epsilon_y)
\end{aligned}$$

Dividing both sides of the inequality by  $|\hat{z}|$  we obtain

$$\epsilon_z \leq eps + \left|\frac{\hat{x}}{\hat{z}}\right| \cdot \epsilon_x + \left|\frac{\hat{y}}{\hat{z}}\right| \cdot \epsilon_y \quad (1)$$

---

<sup>4</sup>The sufficiency of this condition follows immediately from [Knu80] Theorem 4.4.4 B

This formula has index 4. According to section 3 we may therefore use the correction factor  $1 + 8 \cdot eps$  provided that  $\epsilon_x$  and  $\epsilon_y$  are both bounded by  $2^{900}$ .

Assume next that  $|\hat{z}| < MinDbl$ . Then  $|\hat{x} + \hat{y}| < MinDbl$  and hence

$$\begin{aligned} |z| &\leq |x - \hat{x}| + |y - \hat{y}| + |\hat{x} + \hat{y}| \\ &\leq |\hat{x}| \cdot \epsilon_x + |\hat{y}| \cdot \epsilon_y + MinDbl \\ &\leq (1 + eps)^2 \cdot (|\hat{x}| \odot \epsilon_x \oplus |\hat{y}| \odot \epsilon_y) + 3MinDbl \\ &\leq (1 + eps)^3 \cdot (|\hat{x}| \odot \epsilon_x \oplus |\hat{y}| \odot \epsilon_y \oplus 4MinDbl) \end{aligned}$$

where the term  $2MinDbl$  accounts for the possible underflow in the multiplications. With  $NUM = (|\hat{x}| \odot \epsilon_x \oplus |\hat{y}| \odot \epsilon_y \oplus 4MinDbl)$  we have  $|x + y - NUM| \leq 2 \cdot (1 + eps)^3 |NUM|$ . We may therefore set  $ERROR = 2 * correction$ .

```

⟨ compute error for addition or subtraction 38 ⟩ ≡
  if (fabs(NUM) ≥ MinDbl) {
    z_rep.ERROR = fabs(x_num/NUM) * x.PTR→ERROR + fabs(y_num/NUM) *
      y.PTR→ERROR + eps;
  }
  else {
    NUM = fabs(x_num)*x.PTR→ERROR + fabs(y_num)*y.PTR→ERROR + 4*MinDbl;
    z_rep.ERROR = 2;
  }
  z_rep.ERROR *= correction;

```

This code is used in chunk 32.

**39.** We turn to multiplication-specific parts. Checking exactness of **double** multiplication seems more difficult than for addition. However, if both operands are “single-precision”, that is, if the last 26 bits are zero, we know that the exact value of the multiplication is representable in **double** format. Hence in this case the **double** multiplication is guaranteed to compute exactly.

Note that the position of the **double**’s low word is machine-dependent. For example, for SUN SPARC’s the *BIG\_ENDIAN* byte ordering is active, an INTEL PC’s use *LITTLE\_ENDIAN* byte ordering.

```

⟨ if multiplication is exact return double result 39 ⟩ ≡
  double NUM_result;
  long *x_ptr, *y_ptr;
  if (x.is_double ∧ y.is_double) {
    NUM_result = x.value * y.value;
    x_ptr = (long *) &(x.value);
    y_ptr = (long *) &(y.value);
  #ifndef LITTLE_ENDIAN
    x_ptr++; y_ptr++;
  #endif
  if (¬(*x_ptr & #07ffffff) ∧ ¬(*y_ptr & #07ffffff)) {
    NUM_result = x.value * y.value;
    return NUM_result;
  }

```

}

This code is used in chunk 32.

**40.** Next we analyze multiplication. If either  $\hat{x} = 0$  or  $\hat{y} = 0$  then  $z = 0$ . If  $|\hat{z}| < MinDbl$  then  $|\hat{x} \cdot \hat{y}| \leq MinDbl$  and hence

$$\begin{aligned} |x \cdot y| &\leq |\hat{x}| \cdot (1 + \epsilon_x) \cdot |\hat{y}| \cdot (1 + \epsilon_y) \\ &\leq MinDbl \cdot (1 + \epsilon_x) \cdot (1 + \epsilon_y) \end{aligned}$$

We set  $NUM$  to  $MinDbl \cdot (1 + \epsilon_x) \cdot (1 + \epsilon_y)$  and  $ERROR$  to  $2 * correction$ . If  $|\hat{z}| \geq MinDbl$  and  $|\hat{z}| \neq \infty$  then

$$\begin{aligned} |\hat{z} - (x \cdot y)| &\leq |\hat{z} - \hat{x} \cdot \hat{y}| + |\hat{x} \cdot (\hat{y} - y)| + |y \cdot (\hat{x} - x)| \\ &\leq |\hat{z}| \cdot eps + |\hat{x} \cdot \hat{y}| \cdot \epsilon_y + (|\hat{y}| + |y - \hat{y}|) \cdot |\hat{x}| \cdot \epsilon_x \\ &\leq |\hat{z}| \cdot eps + |\hat{x} \cdot \hat{y}| \cdot \epsilon_y + |\hat{x} \cdot \hat{y}| \cdot (1 + \epsilon_y) \cdot \epsilon_x \\ &= |\hat{z}| \cdot eps + |\hat{x} \cdot \hat{y}| \cdot (\epsilon_x + \epsilon_y + \epsilon_x \cdot \epsilon_y) \end{aligned}$$

Because of  $|\hat{x} \cdot \hat{y}| \leq (1 + eps) \cdot |\hat{z}|$ , we finally get after dividing both sides by the term  $|\hat{z}|$  the estimate

$$\epsilon_{x \cdot y} \leq (\epsilon_x + \epsilon_y + \epsilon_x \cdot \epsilon_y) \cdot (1 + eps) + eps \quad (2)$$

Since we post-multiply the error bound by the correction factor  $1 + 8 \cdot eps$  we may drop the factor  $1 + eps$ .

```

⟨ compute error for multiplication 40 ⟩ ≡
  double &qx = x.PTR→ERROR;
  double &qy = y.PTR→ERROR;
  if (fabs(NUM) ≥ MinDbl) {
    z_rep.ERROR = qx + qy + qx * qy + eps;
  }
  else {
    if ((x_num ≠ 0) ∧ (y_num ≠ 0)) {
      NUM = MinDbl * (1 + qx) * (1 + qy);
      z_rep.ERROR = 2;
    }
    else return 0;
  }
  z_rep.ERROR *= correction;

```

This code is used in chunk 32.

**41.** Now we come to divisions. The only part where division globally differs from the previous operations is the treatment of the zero-exception. Also, we must set the error to Infinity if the sign of the divisor is unsure.

```

⟨ Operators 32 ⟩ +≡
  real operator/(const real &x, const real &y)
  {
    ⟨ check for division by zero 42 ⟩
    ⟨ if division is exact return double result 43 ⟩
    ⟨ create nodes for two operands 33 ⟩
  }

```

```

real_rep &z_rep = *new real_rep ();
⟨ initialize z_rep as result of a binary operation 34 ⟩
z_rep.CON = DIVISION;
NUM = x_num/y_num;
⟨ compute error for division 44 ⟩
⟨ if approximation is not valid set error to Infinity 35 ⟩
return z_rep;
}

```

**42.** We have to check for division by zero. This is easy to do if the divisor is an exact **double**. (If the sign of the divisor is unsure, we later set the error to infinity).

```

⟨ check for division by zero 42 ⟩ ≡
if (y.is_double) {
    if (y.value ≡ 0) error_handler(1, "real::operator/:Division by zero");
}
else if ((y.PTR→NUM ≡ 0) ∧ (y.PTR→ERROR < Infinity))
    error_handler(1, "real::operator/:Division by zero");

```

This code is used in chunk 41.

**43.** One obvious sufficient condition for the exactness of **double** division is the following. Suppose that both the computed division result *NUM\_result* and the dividend *y.NUM* are single-precision. Then the multiplication *NUM\_result \* y.NUM* is exact and we only have to compare with the dividend *x.NUM* to see if *NUM\_result* was exact.

```

⟨ if division is exact return double result 43 ⟩ ≡
double NUM_result;
long *num_ptr, *y_ptr;
if (x.is_double ∧ y.is_double) {
    NUM_result = x.value * y.value;
    num_ptr = (long *) &NUM_result;
    y_ptr = (long *) &y.value;
#ifndef LITTLE_ENDIAN
    num_ptr++; y_ptr++;
#endif
    if (¬(*num_ptr & #07fffffff) ∧ ¬(*y_ptr & #07fffffff))
        if (NUM_result * y.value ≡ x.value) return NUM_result;
}

```

This code is used in chunk 41.

**44.** We analyze the error in **double** division. If  $\hat{x} = 0$  then  $x = 0$  and hence  $z = 0$ . If  $\hat{y} = 0$  then  $y = 0$  and the operation is illegal. If  $\epsilon_y \geq 1$  then  $y$  may be zero and we can give no approximation for the result. So let us assume  $\epsilon_y < 1$ ,  $\hat{x} \neq 0$  and  $\hat{y} \neq 0$ . If  $|\hat{x}/\hat{y}| < MinDbl$  then  $|\hat{x}| \circlearrowleft |\hat{y}| \leq MinDbl$  and

$$|x/y| \leq \frac{|\hat{x}|(1 + \epsilon_x)}{|\hat{y}|(1 - \epsilon_y)} \leq \frac{MinDbl(1 + \epsilon_x)}{1 - \epsilon_y}.$$

We use  $NUM = MinDbl(1 + \epsilon_x)/(1 - \epsilon_y)$  and  $ERROR = 2 * correction$ . Finally, if  $|\hat{z}| \geq MinDbl$  then

$$\begin{aligned}
|\widehat{x/y} - x/y| &\leq |\widehat{x/y} - \hat{x}/\hat{y}| + |\hat{x}/\hat{y} - x/y| \\
&\leq |\widehat{x/y}| \cdot eps + \left| \frac{\hat{x} \cdot y - x \cdot \hat{y}}{y \cdot \hat{y}} \right| \\
&\leq |\widehat{x/y}| \cdot eps + \frac{|\hat{x} \cdot (\hat{y} - y)| + |\hat{y} \cdot (\hat{x} - x)|}{|y \cdot \hat{y}|} \\
&\leq |\widehat{x/y}| \cdot eps + \left| \frac{\hat{x} \cdot \hat{y}}{y \cdot \hat{y}} \right| \cdot (\epsilon_x + \epsilon_y)
\end{aligned}$$

Dividing by  $|\widehat{x/y}|$  yields

$$\begin{aligned}
\epsilon_{x/y} &\leq eps + |\hat{y}/y| \cdot (\epsilon_x + \epsilon_y) \cdot (1 + eps) \\
&\leq eps + \frac{|\hat{y}| \cdot (1 + eps)}{|\hat{y}| - |\hat{y} - y|} \cdot (\epsilon_x + \epsilon_y) \\
&\leq eps + \frac{\epsilon_x + \epsilon_y}{1 - \epsilon_y} \cdot (1 + eps)
\end{aligned}$$

Again, we may absorb the factor  $1 + eps$  into the correction factor.

```

⟨ compute error for division 44 ⟩ ≡
  if (y.PTR→ERROR ≥ 1) {
    z_rep.ERROR = Infinity; NUM = 1;
    return z_rep;
  }
  if (fabs(NUM) ≥ MinDbl)
    z_rep.ERROR = (x.PTR→ERROR + y.PTR→ERROR)/(1 - y.PTR→ERROR) + eps;
  else {
    if (x_num ≠ 0) {
      NUM = MinDbl * (1 + x.PTR→ERROR)/(1 - y.PTR→ERROR);
      z_rep.ERROR = 2;
    }
    else return 0;
  }
  z_rep.ERROR *= correction;

```

This code is used in chunk 41.

**45.** We come to the squareroot operation. Here we have to take care if the argument is negative or if the sign of the argument is unknown. Otherwise we proceed exactly as before, only we have just one argument.

```

⟨ Operators 32 ⟩ +≡
  real sqrt(const real &x)
  {
    ⟨ check for negative argument 46 ⟩
    ⟨ if squareroot is exact return double result 47 ⟩
    if (¬x.PTR) ((real &) x).PTR = new real_rep(x.value);
  }

```

```

real_rep &z_rep = *new real_rep ();
z_rep.OP1 = x.PTR; z_rep.OP1→count++;
z_rep.OP2 = nil;
z_rep.CON = SQUAREROOT;
if (x.PTR→ERROR ≥ 1) {
    z_rep.ERROR = Infinity; z_rep.NUM = 1;
    return real(z_rep);
}
if (x.PTR→NUM < 0)
    error_handler(1, "real::operator□sqrt:negative□argument");
z_rep.NUM = sqrt(x.PTR→NUM);
⟨compute error for squareroot 48⟩
return z_rep;
}

```

46. If the argument of squareroot is negative, we exit with an error message.

```

⟨check for negative argument 46⟩ ≡
if ((x.is_double) ∧ (x.value < 0))
    error_handler(1, "real::operator□sqrt:negative□argument");

```

This code is used in chunk 45.

47. We give a simple test whether the squareroot is exact. If the squareroot approximation is only “single-precision” and if squaring gives back the original value, then it must be exact. Vice versa, if the approximation has a one after the 26<sup>th</sup> most significant digit, then it can never give back the original result by squaring.

```

⟨if squareroot is exact return double result 47⟩ ≡
double NUM_result;
long *num_ptr;
if (x.is_double) {
    NUM_result = sqrt(x.value);
    num_ptr = (long *) &NUM_result;
#ifndef LITTLE_ENDIAN
    num_ptr++;
#endif
if (¬(*num_ptr & #07ffffff)) return NUM_result;
}

```

This code is used in chunk 45.

48. It remains to analyze the error in squareroot.

If  $\hat{x} = 0$  then  $x = z = 0$ . If  $\hat{x} \neq 0$  and  $\epsilon_x > 1$  or  $\hat{x} < 0$  then  $x$  may be negative and we can give no meaningful approximation. So let us finally assume that  $\hat{x} > 0$  and  $\epsilon_x \leq 1$ . Then  $\hat{x} \geq MinDbl$  and hence  $\hat{z} \geq MinDbl$ . We have

$$|\widehat{\sqrt{x}} - \sqrt{x}| \leq |\widehat{\sqrt{x}} - \sqrt{\hat{x}}| + |\sqrt{\hat{x}} - \sqrt{x}|$$



$$\begin{aligned}
&\leq \widehat{\sqrt{x}} \cdot eps + \frac{|(\sqrt{\hat{x}} - \sqrt{x}) \cdot (\sqrt{\hat{x}} + \sqrt{x})|}{\sqrt{\hat{x}} + \sqrt{x}} \\
&= \widehat{\sqrt{x}} \cdot eps + \frac{|\hat{x} - x|}{\sqrt{\hat{x}} + \sqrt{x}} \\
&\leq \widehat{\sqrt{x}} \cdot eps + \frac{\hat{x} \cdot \epsilon_x}{\sqrt{\hat{x}}}
\end{aligned}$$

Dividing by  $\widehat{\sqrt{x}}$  yields

$$\epsilon_{\sqrt{x}} \leq eps + \epsilon_x \cdot (1 + eps) \quad (3)$$

```

⟨ compute error for squareroot 48 ⟩ ≡
  z_rep.ERROR = x.PTR→ERROR + eps;
  z_rep.ERROR *= correction;

```

This code is used in chunk 45.

**49.** Finally we implement the unary “−” operator. Because there is no error in this calculation, we only have to copy the operand with its approximation negated.

```

⟨ Operators 32 ⟩ +≡
  real operator−(const real &x)
  {
    if (x.is_double) return −x.value;
    real_rep &z_rep = *new real_rep ();
    z_rep.CON = NEGATION;
    z_rep.OP1 = x.PTR; z_rep.OP1→count++; z_rep.OP2 = nil;
    z_rep.NUM = −x.PTR→NUM; z_rep.ERROR = x.PTR→ERROR;
    return z_rep;
  }

```

## 50. Sign determination

We discuss the different versions of the function `real::sign`. They all compute the sign of the `real` and their implementation rests on the function

`real_rep::sign_with_separation_bound`.

⟨ Sign functions 50 ⟩ ≡

```
int real::sign() const
{
    if (is_double) return ::sign(value);
    else return PTR→sign_with_separation_bound(zero_integer, false);
}
int real::sign(const integer &q) const
{
    if (is_double) return ::sign(value);
    else return PTR→sign_with_separation_bound(q, true);
}
int real::sign(long p) const
{
    if (is_double) return ::sign(value);
    else return PTR→sign_with_separation_bound(p, true);
}
```

See also chunks 54, 56, 57, 64, and 86.

This code is used in chunk 13.

**51.** For the implementation of the function `sign_with_separation_bound` we need the concept of a *separation bound*. A separation bound for a real number  $x$  is a number  $q$  such that  $x \leq q$  implies  $x = 0$ . We allow the user to provide a separation bound. Following Mignotte [Mig92] we compute for each `real`  $x$  quantities  $M$  and  $deg$  such that  $2^{-M}$  is a separation bound for  $x$ . The quantities  $M$  and  $deg$  are defined inductively according to the following rules.

If  $x$  is an integer then

$$M(x) = \lceil \log_2(|x|) \rceil, deg(x) = 1,$$

if  $x$  is a `bigfloat`,  $x = significant \cdot 2^{exponent}$ , then we have

$$M = 1 + Max(significant.length() + exponent, -exponent),$$

if  $x = y + z$  or  $x = y - z$  then

$$\begin{aligned} M(x) &= deg(y) * deg(z) + deg(z) * M(y) + deg(y) * M(z), \\ deg(x) &= deg(y) * deg(z), \end{aligned}$$

if  $x = y * z$  or  $x = y/z$  then

$$\begin{aligned} M(x) &= deg(z) * M(y) + deg(y) * M(z), \\ deg(x) &= deg(y) * deg(z) \end{aligned}$$

and if  $x = \sqrt{y}$  then

$$M(x) = M(y), deg(x) = 2 * deg(y).$$

52. We store the quantities  $M$  and  $deg$  in a structure.

```

⟨ Basic declarations 16 ⟩ +≡
    struct sep_bound {
        integer M, deg;
        sep_bound() {}
    };

```

53. The following function computes the separation bound of the number represented by a node. It assumes that the separation bounds of the children already exist.

```

⟨ Local functions of class real_rep 53 ⟩ ≡
    void compute_separation_bound();

```

See also chunks 55, 67, 69, and 85.

This code is used in chunk 17.

54.⟨ Sign functions 50 ⟩ +≡

```

void real_rep::compute_separation_bound()
{
    sep_bound *ptr = sep_bound_ptr;
    switch (CON) {
    case DOUBLE: case BIGFLOAT:
        {
            bigfloat &num = APP_BF_PTR->NUM_BF;
            ptr->deg = 1;
            if (isZero(num)) ptr->M = zero_integer;
            else {
                if (sign(num.get_exponent()) ≥ 0) ptr->M = log2(num);
                else ptr->M = 1 + Maximum(num.get_precision(), -num.get_exponent());
            }
            break;
        }
    case NEGATION: ptr->deg = OP1->sep_bound_ptr->deg;
        ptr->M = OP1->sep_bound_ptr->M;
        break;
    case MULTIPLICATION: case DIVISION: ptr->M = OP1->sep_bound_ptr->M *
        OP2->sep_bound_ptr->deg + OP2->sep_bound_ptr->M * OP1->sep_bound_ptr->deg;
        ptr->deg = OP1->sep_bound_ptr->deg * OP2->sep_bound_ptr->deg;
        break;
    case ADDITION: case SUBTRACTION:
        ptr->deg = OP1->sep_bound_ptr->deg * OP2->sep_bound_ptr->deg;
        ptr->M = ptr->deg + OP1->sep_bound_ptr->M * OP2->sep_bound_ptr->deg +
            OP2->sep_bound_ptr->M * OP1->sep_bound_ptr->deg;
        break;
    case SQUAREROOT: ptr->M = OP1->sep_bound_ptr->M;
        ptr->deg = OP1->sep_bound_ptr->deg ≪ 1;
        break;
    }
}

```

```

}
}

```

55. We also need a recursive function that initializes the separation bound for the whole expression tree of a **real\_rep**.

```

⟨ Local functions of class real_rep 53 ⟩ +≡
    void initialize_separation_bound( );

```

```

56.⟨ Sign functions 50 ⟩ +≡
    void real_rep :: initialize_separation_bound( )
    {
        if (sep_bound_ptr) return;
        if (OP1) OP1→initialize_separation_bound( );
        if (OP2) OP2→initialize_separation_bound( );
        sep_bound_ptr = new sep_bound ( );
        compute_separation_bound( );
    }

```

57. Now we come to the function *sign\_with\_separation\_bound*, the core of this section. The basic structure of the function is as follows. We first decide whether the current approximations already determine the sign. If so, we return the sign immediately. Otherwise, we first compute a separation bound for the **real\_rep**. In the main part of *sign*, we repeatedly improve the quality of the **bigfloat** approximation. As long as the current precision is less than  $q$  and less than the separation bound and the current approximation does not determine the sign we square the precision and repeat. If after that loop the sign is still not clear from the current approximation, the sign must be zero and we can remove the operand pointers.

```

⟨ Sign functions 50 ⟩ +≡
    int real_rep :: sign_with_separation_bound(const integer &q, bool supported)
    {
        ⟨ if current approximations determine the sign return it 58 ⟩
        initialize_separation_bound( );
        ⟨ compute bigfloat approximations of increasing precision 60 ⟩
        if ⟨ sign is clear from the bigfloat approximation 59 ⟩ return sign(num_bf( ));
        ⟨ set number to zero 61 ⟩
        return 0;
    }

```

58. When do the current approximations determine the sign? The **double** approximations does when either  $NUM \equiv 0$  or  $ERROR < 1$ .

```

⟨ if current approximations determine the sign return it 58 ⟩ ≡
    if ((ERROR ≠ Infinity) ∧ ((ERROR < 1) ∨ (NUM ≡ 0))) return ::sign(NUM);
    init_app_bf( );
    if ⟨ sign is clear from the bigfloat approximation 59 ⟩ return sign(num_bf( ));

```

This code is used in chunk 57.

**59.** The **bigfloat** approximation determines the sign when  $ERROR\_BF \equiv 0$  or  $ERROR\_BF < abs(NUM\_BF)$ .

$\langle$  sign is clear from the bigfloat approximation 59  $\rangle \equiv$   
 $(exact() \vee (\neg isZero(num\_bf()) \wedge (num\_exp() > err\_exp())))$

This code is used in chunks 57, 58, and 60.

**60.** Our current approximations do not determine the sign and we have computed a maximal precision bound  $q$ . Now we repeatedly improve the working precision until either precision  $q$  or the bound  $bound$  is reached or the sign is determined.

$\langle$  compute bigfloat approximations of increasing precision 60  $\rangle \equiv$   
**bool** *user\_bound\_is\_reached* = *false*;  
**integer** *relative\_precision* = 26;  
**integer** *absolute\_precision*;  
**do** {  
    *relative\_precision* = *relative\_precision* \* 2;  
    *absolute\_precision* = *err\_exp()* - *relative\_precision*;  
    *guarantee\_bound\_two\_to*(*absolute\_precision*);  
    *user\_bound\_is\_reached* = *supported*  $\wedge$  (*absolute\_precision*  $\geq$  *q*);  
**}** **while** ( $\neg$ (*user\_bound\_is\_reached*  $\vee$   $\langle$  sign is clear from the bigfloat approximation 59  $\rangle$ ));

This code is used in chunk 57.

**61.** Finally we assume that the number is smaller than the computed separation bound and consequently is zero.

$\langle$  set number to zero 61  $\rangle \equiv$   
*NUM* = *ERROR* = 0;  
*num\_bf()* = *error\_bf()* = *pZero*;  
 $\langle$  remove operand pointers 62  $\rangle$

This code is used in chunks 57 and 66.

**62.**  $\langle$  remove operand pointers 62  $\rangle \equiv$   
{  
    **if** (*OP1*) **if** ( $--OP1 \rightarrow count \equiv 0$ ) **delete** *OP1*;  
    **if** (*OP2*) **if** ( $--OP2 \rightarrow count \equiv 0$ ) **delete** *OP2*;  
    *OP1* = *OP2* = *nil*;  
    *CON* = *BIGFLOAT*;  
}

This code is used in chunks 61, 66, and 68.

**63.** The comparison operators could have been directly derived from the sign function in the obvious way. But then in most cases we would have to allocate an obsolete new node by calling the subtraction operator. This can be avoided if we use a straightforward test that tells us if the comparison is easy and can be done by comparing the **double**s.

```

⟨check if comparison is easy 63⟩ ≡
double x_error, y_error;
bool is_easy;
if (x.is_double ∧ y.is_double) is_easy = true;
else {
  if (x.is_double) x_error = 0;
  else {
    ((real &) x).value = x.PTR→NUM;
    x_error = fabs(x.value) * x.PTR→ERROR;
  }
  if (y.is_double) y_error = 0;
  else {
    ((real &) y).value = y.PTR→NUM;
    y_error = fabs(y.value) * y.PTR→ERROR;
  }
  is_easy = (fabs(x.value - y.value) > (x_error + y_error + twoMinDbl) * correction);
}

```

This code is used in chunk 64.

**64.** Now we can give the comparison operators. If we know that the comparison is easy, then we compare the **double** values. Otherwise we call the sign function.

```

⟨Sign functions 50⟩ +≡
int operator≡(const real &x, const real &y)
{
  ⟨check if comparison is easy 63⟩
  if (is_easy) return (x.value ≡ y.value);
  return ((y - x).sign() ≡ 0);
}
int operator≠(const real &x, const real &y)
{
  ⟨check if comparison is easy 63⟩;
  if (is_easy) return (x.value ≠ y.value);
  return ((y - x).sign() ≠ 0);
}
int operator < (const real &x, const real &y)
{
  ⟨check if comparison is easy 63⟩;
  if (is_easy) return (x.value < y.value);
  return ((y - x).sign() > 0);
}
int operator > (const real &x, const real &y)
{
  ⟨check if comparison is easy 63⟩;
  if (is_easy) return (x.value > y.value);
  return ((y - x).sign() < 0);
}

```

```
int operator≤(const real &x, const real &y)
{
    ⟨check if comparison is easy 63⟩;
    if (is_easy) return (x.value ≤ y.value);
    return ((y - x).sign() ≥ 0);
}

int operator≥(const real &x, const real &y)
{
    ⟨check if comparison is easy 63⟩;
    if (is_easy) return (x.value ≥ y.value);
    return ((y - x).sign() ≤ 0);
}
```

## 65. The functions `improve` and related functions

Function `real::improve_approximation_to` is realized by calling function `real_rep::guarantee_bound_two_to` with the sign of the argument inversed. The reason to invert the sign is that this syntax makes it easier to understand the numerical details of `guarantee_bound_two_to`.

```
⟨Functions for improving the approximation 65⟩ ≡
void real::improve_approximation_to(const integer &p) const
{
    if (is_double) return;
    PTR→init_app_bf();
    PTR→guarantee_bound_two_to(-p);
    adjust();
}
```

See also chunks 66, 68, 72, 76, 77, 78, 79, 80, 81, 82, 84, and 87.

This code is used in chunk 13.

66. The function `guarantee_bound_two_to` is among the core of our data structure. The call `guarantee_bound_two_to(integer e)` improves `error_bf()` to a value  $\leq 2^e$ . For each of the operations addition, subtraction, multiplication, division and squareroot, we have a function `improve_operation` to do this job. In case that the **bigfloat** approximation is exact, we remove the obsolete pointers `OP1` and `OP2`. If the function changes the **bigfloat** approximation of the `real_rep` it also updates its separation bound. After doing so we check whether the new separation bound and the new approximation together imply that the number is zero.

```
⟨Functions for improving the approximation 65⟩ +≡
void real_rep::guarantee_bound_two_to(const integer &e)
{
    bool current_precision_is_sufficient = isZero(error_bf()) ∨ (err_exp() ≤ e);
    if (current_precision_is_sufficient) return;
    switch (CON) {
    case DOUBLE: case BIGFLOAT: return;
    case NEGATION: OP1→guarantee_bound_two_to(e);
        num_bf() = -OP1→num_bf();
        error_bf() = OP1→error_bf();
        break;
    case ADDITION: improve_add(e); break;
    case SUBTRACTION: improve_sub(e); break;
    case MULTIPLICATION: improve_mul(e); break;
    case DIVISION: improve_div(e); break;
    case SQUAREROOT: improve_sqrt(e); break;
    }
    if (exact()) ⟨remove operand pointers 62⟩
    if (sep_bound_ptr) {
        compute_separation_bound();
        if (log2(error_bf()) + abs(num_bf())) ≤ -sep_bound_ptr→M) {
```



```

    <set number to zero 61>
    sep_bound_ptr→M = zero_integer;
    sep_bound_ptr→deg = 1;
  }
}
adjust_dbl();
}

```

67. At this moment we list only the prototypes of the functions *improve\_op*. The function bodies will be given later, when we present the functions *compute\_op*.

```

<Local functions of class real_rep 53> +≡
void improve_add(const integer &);
void improve_sub(const integer &);
void improve_mul(const integer &);
void improve_div(const integer &);
void improve_sqrt(const integer &);

```

68. We implement the function **real\_rep**::*compute\_op* that we used in **real**::*init\_app\_bf*. The basic structure of *compute\_op* is as follows. Suppose that the **real\_rep** object was constructed by the operation *op* and has operands that represent numbers  $x$  and  $y$  with approximations  $\hat{x}$  and  $\hat{y}$ . Again we distinguish the respective cases of *op*. First we compute  $\hat{x}$  *op*  $\hat{y}$  in precision  $p$ . The error in this operation alone is called the *operation-induced* error and is given by  $|\hat{x}$  *op*  $\hat{y}| \cdot 2^{-p}$ . Then we compute the error arising from the inaccuracies of  $\hat{x}$  and  $\hat{y}$ , which we call the *operand-induced* error. For multiplication, division and squareroot we use special functions to compute the operand-induced error. The overall error of the approximation finally is the sum of the operand-induced and the operation-induced error. We remark that for division the sign of the dividend must be known to be nonzero and similarly the argument of squareroot must be known to be nonnegative.

```

<Functions for improving the approximation 65> +≡
void real_rep::compute_op(long p)
{
  bool isexact;
  bigfloat operation_error, operand_error;
  switch (CON) {
  case DOUBLE: case BIGFLOAT: return;
  case NEGATION: num_bf() = -OP1→num_bf();
    error_bf() = OP1→error_bf();
    break;
  case ADDITION:
    num_bf() = add(OP1→num_bf(), OP2→num_bf(), p, TO_NEAREST, isexact);
    <compute operand error for addition and subtraction 71>
    break;
  case SUBTRACTION:
    num_bf() = sub(OP1→num_bf(), OP2→num_bf(), p, TO_NEAREST, isexact);

```

```

    < compute operand error for addition and subtraction 71 >
    break;
case MULTIPLICATION:
    num_bf() = mul(OP1→num_bf(), OP2→num_bf(), p, TO_NEAREST, isexact);
    compute_mul_error(operand_error);
    break;
case DIVISION:
    if (OP2→sign_with_separation_bound(zero_integer, false) ≡ 0)
        error_handler(1, "compute_op:division_by_zero");
    num_bf() = div(OP1→num_bf(), OP2→num_bf(), p, TO_NEAREST, isexact);
    compute_div_error(operand_error);
    break;
case SQUAREROOT:
    if (OP1→sign_with_separation_bound(zero_integer, false) < 0)
        error_handler(1, "compute_op:sqrt_of_negative_number");
    num_bf() = sqrt(OP1→num_bf(), p, TO_NEAREST, isexact);
    compute_sqrt_error(operand_error);
}
if (¬isexact) operation_error = ldexp_bf(num_bf(), -p);
error_bf() = operation_error + operand_error;
if (exact()) < remove operand pointers 62 >
    adjust_dbl();
}

```

69. We list the prototypes of the functions *compute\_op\_error*.

```

< Local functions of class real_rep 53 > +≡
void compute_mul_error(bigfloat &);
void compute_div_error(bigfloat &);
void compute_sqrt_error(bigfloat &);

```

70. We now give the details of *compute\_op\_error* and *improve\_op* for the various cases of numerical operations *op*. In the context of *compute\_op\_error* we want to compute the error  $q$  of an approximation  $\hat{z}$  of some real number  $z$ . For the function *improve*, we face the opposite problem. There we are *given* a bound  $q$  and we want to determine an approximation  $\hat{z}$  with error at most  $q$ . In both cases, the number  $z$  is defined by two (one in the case of squareroot) other numbers  $x$  and  $y$  by means of the operation *op*.

71. Let's deal with addition and subtraction first. We have

$$\begin{aligned}
 |\hat{z} - z| &\leq |\hat{z} - (\hat{x} + \hat{y})| + |\hat{x} - x| + |\hat{y} - y| \\
 &\leq \epsilon \cdot |\hat{z}| + qx + qy
 \end{aligned}
 \tag{4}$$

where  $\epsilon$  is the relative precision of the **bigfloat** addition. We use this formula to compute the error bounds for addition and subtraction directly with **bigfloats**.

```

< compute operand error for addition and subtraction 71 > ≡
    operand_error = add(OP1→error_bf(), OP2→error_bf(), 16, TO_INF);

```

This code is used in chunk 68.

**72.** Let us discuss how to improve the operations addition and subtraction. Here we need to determine two bounds  $qx$  and  $qy$  and a relative precision  $\epsilon$  such that the following is true: If  $\hat{x}$  and  $\hat{y}$  are approximations of  $x$  and  $y$  of error  $qx$  and  $qy$  respectively and  $\hat{z}$  is computed from  $\hat{x}$  and  $\hat{y}$  by performing operation  $op$  with relative precision  $\epsilon$  then  $\hat{z}$  has maximal error  $q$ . Looking again at Inequality 4, we are done if we bound the first term (the operation induced error) by  $q/2$  and the second and third term (the operand errors) by  $q/4$ . That is, we need

$$\epsilon \leq \frac{q}{2|\hat{z}|}, qx \leq q/4, qy \leq q/4.$$

The second two conditions are easy to satisfy. Taking  $q = 2^e$ , we only need to call  $x.guarantee\_bound\_two\_to(e - 2)$  and  $y.guarantee\_bound\_two\_to(e - 2)$ . How can we guarantee the first condition? The problem is that we do not know  $\hat{z}$ . So we replace  $\log(\hat{z})$  by the crude bound  $Max(\log(\hat{x}), \log(\hat{y})) + 1$ .

```

<Functions for improving the approximation 65> +≡
void real_rep::improve_add(const integer &e)
{
  long p;
  <improve operands and compute p for add and sub 73>
  bool isexact;
  num_bf() = add(OP1→num_bf(), OP2→num_bf(), p, TO_NEAREST, isexact);
  <set error 2e 75>
}
void real_rep::improve_sub(const integer &e)
{
  long p;
  <improve operands and compute p for add and sub 73>
  bool isexact;
  num_bf() = sub(OP1→num_bf(), OP2→num_bf(), p, TO_NEAREST, isexact);
  <set error 2e 75>
}

```

```

73.<improve operands and compute p for add and sub 73> ≡
  OP1→guarantee_bound_two_to(e - 2);
  OP2→guarantee_bound_two_to(e - 2);
  integer log_epsilon = -1;
  if (¬isZero(OP1→num_bf()) ∧ ¬isZero(OP2→num_bf()))
    log_epsilon = e - Maximum(OP1→num_exp(), OP2→num_exp()) - 2;
  <convert log_epsilon to p 74>

```

This code is used in chunk 72.

```

74.<convert log_epsilon to p 74> ≡
  if (¬log_epsilon.islong()) {
    error_handler(1, "sorry: improve: computation□to□expensive");
    p = 10000; /* just to silence compiler warnings */
  }

```

```

}
else p = -log_epsilon.tolong();

```

This code is used in chunks 73, 77, 79, and 81.

```

75.<set error 2e 75> ≡
if (isexact ∧ OP1→exact() ∧ OP2→exact()) error_bf() = pZero;
else error_bf() = pow2(e);

```

This code is used in chunks 72, 77, and 79.

76. We deal with multiplication next. We have

$$\begin{aligned}
|\widehat{x \cdot y} - (x \cdot y)| &\leq |\widehat{x \cdot y} - (\hat{x} \cdot \hat{y})| + |\hat{x} \cdot (\hat{y} - y)| + |y \cdot (\hat{x} - x)| \\
&\leq |\widehat{x \cdot y}| \cdot \epsilon + |\hat{x}| \cdot qy + |y| \cdot qx
\end{aligned} \tag{5}$$

We first implement the part of computing errors. Since multiplications are expensive for operands with large precision, we first round them to 16 bit and then do the multiplications in precision 32. By choosing rounding mode “to  $\infty$ ” we avoid having to care for the inaccuracies.

```

<Functions for improving the approximation 65> +≡
void real_rep::compute_mul_error(bigfloat &operand_error)
{
  bigfloat x = OP1→num_bf();
  bigfloat y = OP2→num_bf();
  bigfloat error_x = mul(abs(round(x, 16, TO_INF)), OP2→error_bf(), 32, TO_INF);
  bigfloat error_y = mul(abs(round(y, 16, TO_INF)), OP1→error_bf(), 32, TO_INF);
  operand_error = add(error_x, error_y, 16, TO_INF);
}

```

77. Now we implement the improve function for multiplication. Again, our goal is to bound the first term of Relation 5 by  $q/2$  and the two other terms by  $q/4$ , or equivalently

$$\epsilon \leq \frac{q}{2|\widehat{x \cdot y}|}, qx \leq \frac{q}{4|y|}, qy \leq \frac{q}{4|\hat{x}|} \tag{6}$$

Since we don't know  $y$ , we bound it by  $y\_high = abs(y.num_bf()) + y.error_bf()$ .

```

<Functions for improving the approximation 65> +≡
void real_rep::improve_mul(const integer &e)
{
  bigfloat y_high = add(abs(OP2→num_bf()), OP2→error_bf(), 32, TO_P_INF);
  integer ex = e - log2(y_high) - 2;
  OP1→guarantee_bound_two_to(ex);
  integer log_epsilon = -1;
  if (¬isZero(OP1→num_bf())) {
    integer ey = e - OP1→num_exp() - 2;

```

```

    OP2→guarantee_bound_two_to(ey);
    log_epsilon = e - (OP1→num_exp() + OP2→num_exp()) - 1;
}
bool isexact;
long p;
⟨convert log_epsilon to p 74⟩
num_bf() = mul(OP1→num_bf(), OP2→num_bf(), p, TO_NEAREST, isexact);
⟨set error 2e 75⟩
}

```

78. Division is next. Remember that when the functions *compute\_div\_error* and *improve\_div* are called the sign of the divisor  $y$  is known to be nonzero, because this exception was handled in the division operator or by *init\_app\_bf*. Hence we can use  $y\_low = abs(y.num\_bf()) - y.error\_bf() > 0$  in our calculation. Let us consider the error estimate given for the division operator.

$$\begin{aligned}
 \left| \frac{\hat{x}}{y} - \frac{x}{y} \right| &\leq \left| \frac{\hat{x}}{y} \right| \cdot \epsilon + \frac{|\hat{x}| \cdot qy}{|y \cdot \hat{y}|} + \frac{qx}{|y|} \\
 &= \left| \frac{\hat{x}}{y} \right| \cdot \epsilon + \frac{1}{|y|} \cdot \left( \frac{|\hat{x}| \cdot qy}{|\hat{y}|} + qx \right)
 \end{aligned} \tag{7}$$

As before we avoid multiplication with operands having large precisions and also we avoid divisions completely.

```

⟨Functions for improving the approximation 65⟩ +≡
void real_rep::compute_div_error(bigfloat &operand_error)
{
    bigfloat y_low = sub(abs(OP2→num_bf()), OP2→error_bf(), 32, TO_N_INF);
    bigfloat x = OP1→num_bf();
    bigfloat z = mul(abs(round(x, 16)), OP2→error_bf(), 32, TO_INF);
    z = ldexp_bf(z, 1 - OP2→num_exp());
    operand_error = add(z, abs(OP1→error_bf()), 16, TO_INF);
    operand_error = ldexp_bf(operand_error, 1 - log2(y_low));
}

```

79. In the div version of *improve*, we again want to bound the last two terms of Relation 7 by  $q/4$  and the first one by  $q/2$ . This means

$$\epsilon \leq \frac{q}{2|\hat{z}|}, qx \leq \frac{q \cdot y}{4}, qy \leq \frac{q \cdot y \cdot |\hat{y}|}{4|\hat{x}|} \tag{8}$$

We proceed exactly as for the other cases of *improve*.

```

⟨Functions for improving the approximation 65⟩ +≡
void real_rep::improve_div(const integer &e)
{
    bigfloat y_low = sub(abs(OP2→num_bf()), OP2→error_bf(), 32, TO_N_INF);
    integer ex = e + log2(y_low) - 3;
}

```

```

    OP1→guarantee_bound_two_to(ex);
integer log_epsilon = -1;
if (¬isZero(OP1→num_bf())) {
    integer ey = ex + OP2→num_exp() - OP1→num_exp() - 3;
    OP2→guarantee_bound_two_to(ey);
    log_epsilon = e + OP2→num_exp() - OP1→num_exp() - 2;
}
bool isexact;
long p;
⟨convert log_epsilon to p 74⟩
num_bf() = div(OP1→num_bf(), OP2→num_bf(), p, TO_NEAREST, isexact);
⟨set error 2e 75⟩
}

```

80. We finally deal with the squareroot operation. We have

$$|\widehat{\sqrt{x}} - \sqrt{x}| \leq \widehat{\sqrt{x}} \cdot \epsilon + \frac{|\hat{x} - x|}{\sqrt{\hat{x}} + \sqrt{x}} \leq \widehat{\sqrt{x}} \cdot \epsilon + \frac{qx}{\sqrt{\hat{x}}} \quad (9)$$

We bound the logarithm of  $\sqrt{\hat{x}}$  by  $(x.\text{num\_exp}() + 1)/2$ .

```

⟨Functions for improving the approximation 65⟩ +≡
void real_rep::compute_sqrt_error(bigfloat &operand_error)
{
    integer log_sqrt = (num_exp() + 1) ≫ 1;
    if (¬OP1→exact()) operand_error = ldexp_bf(OP1→error_bf(), -log_sqrt);
}

```

81. Last we show how to improve the squareroot operation. To bound the two terms on the right hand side of Estimate 9 by  $q/2$ , we set  $qx = q \cdot \sqrt{x}/2$  and  $\epsilon = q/(2\sqrt{x})$ . Since we don't know  $\sqrt{x}$ , we bound its binary logarithm by  $\lfloor \log_2(x_{\text{low}})/2 \rfloor$  from below and by  $\lceil \log_2(x_{\text{high}})/2 \rceil$  from above.

```

⟨Functions for improving the approximation 65⟩ +≡
void real_rep::improve_sqrt(const integer &e)
{
    bool isexact;
    bigfloat x_low = sub(OP1→num_bf(), OP1→error_bf(), 32, TO_N_INF);
    integer ex = e + (log2(x_low) - 1)/2 - 1;
    OP1→guarantee_bound_two_to(ex);
    bigfloat x_high = add(OP1→num_bf(), OP1→error_bf(), 32, TO_P_INF);
    integer log_epsilon = e - (log2(x_high) + 1)/2 - 1;
    long p;
    ⟨convert log_epsilon to p 74⟩
    num_bf() = sqrt(OP1→num_bf(), p, TO_NEAREST, isexact, num_bf());
    if (isexact ∧ OP1→exact()) error_bf() = pZero;
    else error_bf() = pow2(e);
}

```

82. It remains to implement the functions *x.guarantee\_relative\_error* and *x.compute\_with\_precision(p)*. We begin with function *guarantee\_relative\_error*. It first performs a sign test. If the sign is nonzero, it calls *improve* with an appropriate error bound.

```

⟨ Functions for improving the approximation 65 ⟩ +≡
  void real::guarantee_relative_error(long p) const
  {
    if (is_double) return;
    if (PTR→sign_with_separation_bound(zero_integer) ≡ 0) return;
    PTR→init_app_bf();
    bigfloat num_low = abs(PTR→num_bf()) - PTR→error_bf();
    PTR→guarantee_bound_two_to(log2(num_low) - p);
    adjust();
  }

```

83. Function *compute\_with\_precision* uses the function *compute\_approximation* that in turn calls *compute\_op* for every unvisited node in the expression tree. To allow intelligent searching in the expression tree, we assign values of an enumeration type *statustype* to the member *status*. The meaning of these states is the following.

- *EMPTY* is the initial state.
- *VISITED* means that the computation for the node has been successfully finished.
- *CLEARED* is used by the function *clear\_visited\_marks* that we define in a moment.

```

⟨ Basic declarations 16 ⟩ +≡
  enum statustype { EMPTY = 0, CLEARED = 1, VISITED = 2 };

```

```

84.⟨ Functions for improving the approximation 65 ⟩ +≡
  void real::compute_with_precision(long p) const
  {
    if (is_double) return;
    PTR→init_app_bf();
    PTR→clear_visited_marks();
    PTR→compute_approximation(p);
    adjust();
  }

```

85. We need a function to clear the visited marks used by *compute\_approximation*.

```

⟨ Local functions of class real_rep 53 ⟩ +≡
  void clear_visited_marks();

```

**86.** This function is implemented as follows. We maintain the invariant that if any node is in the state *CLEARED*, then all children of that node are *CLEARED* too.

```
⟨ Sign functions 50 ⟩ +≡  
void real_rep::clear_visited_marks()  
{  
  if (status ≡ CLEARED) return;  
  if (OP1) OP1→clear_visited_marks();  
  if (OP2) OP2→clear_visited_marks();  
  status = CLEARED;  
}
```

**87.** Now we consider the function *compute\_approximation*(**long** *p*). When the function is called for the object *x\_rep*, we compute every operation in the expression tree of *x\_rep* with **bigfloat**s in precision *p*. The final result is stored as the new **bigfloat** approximation of *x\_rep*.

```
⟨ Functions for improving the approximation 65 ⟩ +≡  
void real_rep::compute_approximation(long p)  
{  
  if (status ≡ VISITED) return;  
  if (OP1) OP1→compute_approximation(p);  
  if (OP2) OP2→compute_approximation(p);  
  compute_op(p);  
  status = VISITED;  
}
```

**88.** References.

## References

- [BJMM93] M.O. Benouamer, P. Jaillon, D. Michelucci, and J-M. Moreau. A “lazy” solution to imprecision in computational geometry. In *Proceedings of the 5th Canadian Conference on Computational Geometry*, pages 73–78, 1993.
- [BMS94] C. Burnikel, K. Mehlhorn, and S. Schirra. How to compute the Voronoi diagram of line segments: Theoretical and experimental results. In *LNCS*, volume 855, pages 227–239. Springer-Verlag Berlin/New York, 1994. Proceedings of ESA’94.
- [FM91] S. Fortune and V. Milenkovic. Numerical stability of algorithms for line arrangements. In *Proceedings of the 7th ACM Symposium on Computational Geometry*, pages 334–341, 1991.
- [For92] S. Fortune. Numerical stability of algorithms for 2d Delaunay triangulations and Voronoi diagrams. In *Proceedings of the 8th Symposium on Computational Geometry*, pages 83–92, 1992.



- [Fv93] S. Fortune and C. van Wyk. Efficient exact arithmetic for computational geometry. *Proceedings of the 9th Symposium on Computational Geometry*, pages 163–171, 1993.
- [IEE87] IEEE standard 754-1985 for binary floating-point arithmetic. *reprinted in SIGPLAN 22*, 2:9–25, 1987.
- [KLN91] M. Karasick, D. Lieber, and L.R. Nackman. Efficient Delaunay triangulation using rational arithmetic. *ACM Transactions on Graphics*, 10(1):71–91, 1991.
- [Knu80] D.E. Knuth. *The Art of Computer Programming, Volume II*. Addison-Wesley, 1980.
- [LM92] Z. Li and V. Milenkovic. Constructing strongly convex hulls using exact or rounded arithmetic. *Algorithmica*, 8:345–364, 1992.
- [Mig92] M. Mignotte. *Mathematics for Computer Algebra*. Springer Verlag, 1992.
- [MN94a] K. Mehlhorn and S. Näher. Implementation of a sweep line algorithm for the straight line segment intersection problem. Technical Report MPI-I-94-160, Max-Planck-Institut für Informatik, Saarbrücken, 1994.
- [MN94b] K. Mehlhorn and St. Näher. The implementation of geometric algorithms. In *13th World Computer Congress IFIP94*, volume 1, pages 223–231. Elsevier Science B.V. North-Holland, Amsterdam, 1994.
- [MN95] K. Mehlhorn and St. Näher. LEDA: A library of efficient data types and algorithms. *CACM*, 38(1):96–102, 1995.
- [Näh95] St. Näher. LEDA Manual Version 3.2. Technical Report MPI-I-95-1-002, Max-Planck-Institut für Informatik, 1995.
- [Yap93] C.K. Yap. Towards exact geometric computation. In *Proceedings of the 5th Canadian Conference on Computational Geometry*, pages 405–419, 1993.
- [Yap94] C.K. Yap. *Fundamental Problems in Algorithmic Algebra*. Princeton University Press, 1994.
- [YD95] C. Yap and T. Dube. The exact computation paradigm. In *Computing in Euclidean Geometry II*. World Scientific Press, 1995.

## Index

*abs*: [12](#), [23](#), [30](#), [59](#), [66](#), [76](#), [77](#), [78](#), [79](#), [82](#).  
*absolute\_precision*: [60](#).  
*add*: [4](#), [27](#), [68](#), [71](#), [72](#), [76](#), [77](#), [78](#), [81](#).  
*ADDITION*: [16](#), [32](#), [54](#), [66](#), [68](#).  
*adjust*: [25](#), [26](#), [65](#), [82](#), [84](#).  
*adjust\_dbl*: [17](#), [20](#), [22](#), [24](#), [26](#), [28](#), [66](#), [68](#).  
**app-bf**: [16](#).  
*APP\_BF\_PTR*: [17](#), [20](#), [21](#), [23](#), [28](#), [54](#).  
*bf\_app*: [16](#).  
*BIG\_ENDIAN*: [39](#).  
**bigfloat**: [28](#).  
*BIGFLOAT*: [16](#), [20](#), [54](#), [62](#), [66](#), [68](#).  
*bound*: [60](#).  
*clear\_visited\_marks*: [83](#), [84](#), [85](#), [86](#).  
*CLEARED*: [83](#), [86](#).  
*compare*: [11](#).  
*compute\_approximation*: [17](#), [83](#), [84](#), [85](#), [87](#).  
*compute\_div\_error*: [68](#), [69](#), [78](#).  
*compute\_mul\_error*: [68](#), [69](#), [76](#).  
*compute\_op*: [17](#), [23](#), [67](#), [68](#), [83](#), [87](#).  
*compute\_op\_error*: [69](#), [70](#).  
*compute\_separation\_bound*: [53](#), [54](#), [56](#), [66](#).  
*compute\_sqrt\_error*: [68](#), [69](#), [80](#).  
*compute\_with\_precision*: [7](#), [17](#), [82](#), [83](#), [84](#).  
*CON*: [17](#), [20](#), [23](#), [31](#), [32](#), [41](#), [45](#), [49](#), [54](#), [62](#), [66](#), [68](#).  
**constructortype**: [16](#), [17](#).  
*correction*: [3](#), [23](#), [24](#), [38](#), [40](#), [44](#), [48](#), [63](#).  
*count*: [17](#), [19](#), [20](#), [21](#), [34](#), [45](#), [49](#), [62](#).  
*current\_precision\_is\_sufficient*: [66](#).  
*DEBUG\_REAL*: [21](#).  
*deg*: [52](#), [54](#), [66](#).  
*denominator*: [18](#).  
*dist*: [12](#), [30](#).  
*div*: [4](#), [68](#), [79](#).  
*DIVISION*: [16](#), [41](#), [54](#), [66](#), [68](#).  
**double**: [28](#).  
*DOUBLE*: [16](#), [20](#), [54](#), [66](#), [68](#).  
*e*: [17](#), [66](#), [72](#), [77](#), [79](#), [81](#).  
*e\_head*: [24](#).  
*EMPTY*: [83](#).  
*eps*: [3](#), [24](#), [38](#), [40](#), [44](#), [48](#).  
*err\_exp*: [17](#), [22](#), [59](#), [60](#), [66](#).  
*ERROR*: [17](#), [20](#), [23](#), [24](#), [28](#), [35](#), [38](#), [40](#), [42](#), [44](#), [45](#), [48](#), [49](#), [58](#), [61](#), [63](#).  
*ERROR\_BF*: [16](#), [17](#), [24](#), [59](#).  
*error\_bf*: [17](#), [24](#), [27](#), [28](#), [61](#), [66](#), [68](#), [71](#), [75](#), [76](#), [77](#), [78](#), [79](#), [80](#), [81](#), [82](#).  
*error\_handler*: [42](#), [45](#), [46](#), [68](#), [74](#).  
*error\_x*: [16](#), [76](#).  
*error\_y*: [76](#).  
*ex*: [77](#), [79](#), [81](#).  
*exact*: [17](#), [26](#), [59](#), [66](#), [68](#), [75](#), [80](#), [81](#).  
*ey*: [77](#), [79](#).  
*fabs*: [24](#), [35](#), [38](#), [40](#), [44](#), [63](#).  
*false*: [15](#), [17](#), [18](#), [19](#), [50](#), [60](#), [68](#).  
*get\_bigfloat\_error*: [7](#), [28](#).  
*get\_dblError*: [22](#).  
*get\_double\_error*: [7](#), [28](#).  
*get\_exponent*: [29](#), [54](#).  
*get\_precision*: [22](#), [54](#).  
*get\_significant*: [29](#).  
*guarantee\_bound\_two\_to*: [17](#), [60](#), [65](#), [66](#), [72](#), [73](#), [77](#), [79](#), [81](#), [82](#).  
*guarantee\_relative\_error*: [7](#), [82](#).  
*I*: [7](#).  
*improve*: [70](#), [79](#), [82](#).  
*improve\_add*: [66](#), [67](#), [72](#).  
*improve\_approximation\_to*: [7](#), [65](#).  
*improve\_div*: [66](#), [67](#), [78](#), [79](#).  
*improve\_mul*: [66](#), [67](#), [77](#).  
*improve\_op*: [67](#), [70](#).  
*improve\_operation*: [66](#).  
*improve\_sqrt*: [66](#), [67](#), [81](#).  
*improve\_sub*: [66](#), [67](#), [72](#).  
*in*: [11](#), [27](#).  
*Infinity*: [3](#), [23](#), [24](#), [35](#), [42](#), [44](#), [45](#), [58](#).  
*init\_app\_bf*: [17](#), [22](#), [23](#), [27](#), [28](#), [58](#), [65](#), [68](#), [78](#), [82](#), [84](#).  
*initialize\_separation\_bound*: [55](#), [56](#), [57](#).  
**int**: [50](#), [57](#).  
*is\_double*: [15](#), [18](#), [19](#), [26](#), [27](#), [28](#), [36](#), [37](#), [39](#), [42](#), [43](#), [46](#), [47](#), [49](#), [50](#), [63](#), [65](#), [82](#), [84](#).  
*is\_easy*: [63](#), [64](#).  
*isdouble*: [5](#), [18](#), [26](#).  
*isexact*: [68](#), [72](#), [75](#), [77](#), [79](#), [81](#).  
*islong*: [74](#).  
*isZero*: [17](#), [54](#), [59](#), [66](#), [73](#), [77](#), [79](#).

*k*: 7, 17, 29.  
*ldexp\_bf*: 29, 68, 78, 80.  
*LEDA\_MEMORY*: 17.  
*LITTLE\_ENDIAN*: 39, 43, 47.  
*log\_epsilon*: 73, 74, 77, 79, 81.  
*log\_sqrt*: 80.  
*log2*: 17, 27, 54, 66, 77, 78, 79, 81, 82.  
*low\_bound*: 27.  
*M*: 52.  
*Max*: 24.  
*MaxDbl*: 3, 17, 24.  
*MaxError*: 3, 17, 24, 35.  
*Maximum*: 29, 54, 73.  
*MinDbl*: 3, 17, 24, 38, 40, 44.  
*mul*: 4, 68, 76, 77, 78.  
*MULTIPLICATION*: 16, 32, 54, 66, 68.  
*n*: 12, 30.  
*n\_head*: 24.  
*n\_prefix*: 30.  
*NaN*: 35.  
*NEGATION*: 16, 49, 54, 66, 68.  
*nil*: 17, 18, 19, 20, 21, 45, 49, 62.  
*num*: 54.  
*NUM*: 17, 20, 23, 24, 28, 32, 33, 34,  
35, 38, 40, 41, 42, 43, 44, 45, 49,  
58, 61, 63.  
*num\_bf*: 17, 24, 26, 27, 28, 57, 58,  
59, 61, 66, 68, 72, 73, 76, 77, 78,  
79, 81, 82.  
*NUM\_BF*: 16, 17, 24, 54, 59.  
*NUM\_EXACT*: 16, 17, 24.  
*num\_exp*: 17, 22, 59, 73, 77, 78, 79, 80.  
*num\_low*: 82.  
*num\_ptr*: 43, 47.  
*NUM\_result*: 36, 37, 39, 43, 47.  
*numerator*: 18.  
*O*: 7.  
*op*: 70.  
*operand\_error*: 68, 71, 76, 78, 80.  
*operation\_error*: 68.  
**operator**: 7, 8, 9, 10, 17, 27, 32,  
41, 49, 64.  
*OP1*: 17, 18, 20, 21, 23, 34, 45, 49, 54,  
56, 62, 66, 68, 71, 72, 73, 75, 76,  
77, 78, 79, 80, 81, 86, 87.  
*OP2*: 17, 18, 20, 21, 23, 34, 45, 49,  
54, 56, 62, 66, 68, 71, 72, 73, 75,  
76, 77, 78, 79, 86, 87.  
*out*: 11, 27.  
*p*: 4, 17, 50, 65, 68, 72, 77, 79, 81,  
82, 84, 87.  
*powi*: 12, 30.  
*pow2*: 3, 75, 81.  
*prec*: 27.  
*Print*: 11.  
*ptr*: 54.  
*PTR*: 15, 18, 19, 21, 26, 27, 28, 32,  
33, 34, 35, 38, 40, 42, 44, 45, 48,  
49, 50, 63, 65, 82, 84.  
*pZero*: 16, 61, 75, 81.  
*q*: 7, 50, 57.  
*qx*: 40.  
*qy*: 40.  
*Read*: 11.  
**real**: 7, 17, 18, 19, 21.  
*REAL\_H*: 7.  
**real\_rep**: 7, 17, 20, 21.  
*relative\_precision*: 60.  
*round*: 76, 78.  
**sep\_bound**: 17, 52, 54, 56.  
*sep\_bound\_ptr*: 17, 20, 21, 54, 56, 66.  
*set\_glob\_prec*: 4.  
*sign*: 7, 11, 29, 30, 50, 54, 57, 58, 64.  
*sign\_with\_separation\_bound*: 17, 50,  
51, 57, 68, 82.  
*sq*: 12, 30.  
*sqrt*: 9, 30, 45, 47, 68, 81.  
*SQUAREROOT*: 16, 45, 54, 66, 68.  
*status*: 17, 20, 83, 86, 87.  
**statustype**: 83.  
*sub*: 4, 27, 68, 72, 78, 79, 81.  
*SUBTRACTION*: 16, 32, 54, 66, 68.  
*supported*: 17, 57, 60.  
*TO\_INF*: 71, 76, 78.  
*TO\_N\_INF*: 27, 78, 79, 81.  
*TO\_NEAREST*: 68, 72, 77, 79, 81.  
*TO\_P\_INF*: 27, 77, 81.  
*tobigfloat*: 7, 17, 22, 27, 28.  
*todouble*: 5, 7, 17, 18, 22, 24, 26, 28.  
*tolong*: 27, 74.  
*true*: 18, 19, 26, 50, 63.  
*twoMinDbl*: 3, 63.  
*Type\_Name*: 11.  
*upp\_bound*: 27.

*user\_bound\_is\_reached*: 60.  
*value*: 15, 18, 19, 26, 27, 28, 33, 36, 37,  
39, 42, 43, 45, 46, 47, 49, 50, 63, 64.  
*VISITED*: 83, 87.  
**void**: 23, 24, 26, 54, 56, 65, 66, 68, 72,  
76, 77, 78, 79, 80, 81, 82, 84, 86, 87.  
*x*: 5, 7, 10, 11, 12, 16, 18, 19, 20, 27,  
28, 29, 30, 32, 41, 45, 49, 64, 76, 78.  
*x\_error*: 63.  
*x\_high*: 81.  
*x\_low*: 81.  
*x\_num*: 27, 33, 38, 40, 41, 44.  
*x\_ptr*: 39.  
*x\_rep*: 18, 87.  
*y*: 10, 11, 12, 29, 30, 32, 41, 64, 76.  
*y\_error*: 63.  
*y\_high*: 77.  
*y\_low*: 78, 79.  
*y\_num*: 33, 38, 40, 41.  
*y\_ptr*: 39, 43.  
*z*: 30, 78.  
*z\_rep*: 32, 34, 35, 38, 40, 41, 44, 45,  
48, 49.  
*zero\_integer*: 6, 50, 54, 66, 68, 82.

## List of Refinements

- ⟨ Arithmetic operators 9 ⟩ Used in chunks 7 and 17.
- ⟨ Assignment operators 10 ⟩ Used in chunk 7.
- ⟨ Basic declarations 16, 52, 83 ⟩ Used in chunk 13.
- ⟨ Basic functions 23, 24, 26, 27, 28 ⟩ Used in chunk 13.
- ⟨ Comparison operators 8 ⟩ Used in chunks 7 and 17.
- ⟨ Constructors and destructors 18, 19, 20, 21 ⟩ Used in chunk 13.
- ⟨ Declaration of class **real\_rep** 17 ⟩ Used in chunk 13.
- ⟨ Declaration of mathematical functions in real.h 12 ⟩ Used in chunk 7.
- ⟨ Functions for improving the approximation 65, 66, 68, 72, 76, 77, 78, 79, 80, 81, 82, 84, 87 ⟩  
Used in chunk 13.
- ⟨ Global functions and constants 3, 5, 6, 29 ⟩ Used in chunk 13.
- ⟨ Implementation of mathematical functions in real.c 30 ⟩ Used in chunk 13.
- ⟨ LEDA functions 11 ⟩ Used in chunk 7.
- ⟨ Local functions of class **real\_rep** 53, 55, 67, 69, 85 ⟩ Used in chunk 17.
- ⟨ Operators 32, 41, 45, 49 ⟩ Used in chunk 13.
- ⟨ Sign functions 50, 54, 56, 57, 64, 86 ⟩ Used in chunk 13.
- ⟨ check for division by zero 42 ⟩ Used in chunk 41.
- ⟨ check for negative argument 46 ⟩ Used in chunk 45.
- ⟨ check if comparison is easy 63 ⟩ Used in chunk 64.
- ⟨ compute bigfloat approximations of increasing precision 60 ⟩ Used in chunk 57.
- ⟨ compute error for addition or subtraction 38 ⟩ Used in chunk 32.
- ⟨ compute error for division 44 ⟩ Used in chunk 41.
- ⟨ compute error for multiplication 40 ⟩ Used in chunk 32.
- ⟨ compute error for squareroot 48 ⟩ Used in chunk 45.
- ⟨ compute operand error for addition and subtraction 71 ⟩ Used in chunk 68.
- ⟨ convert log\_epsilon to p 74 ⟩ Used in chunks 73, 77, 79, and 81.
- ⟨ create nodes for two operands 33 ⟩ Used in chunks 32 and 41.
- ⟨ if addition is exact return double result 36 ⟩ Used in chunk 32.
- ⟨ if approximation is not valid set error to Infinity 35 ⟩ Used in chunks 32 and 41.
- ⟨ if current approximations determine the sign return it 58 ⟩ Used in chunk 57.
- ⟨ if division is exact return double result 43 ⟩ Used in chunk 41.
- ⟨ if multiplication is exact return double result 39 ⟩ Used in chunk 32.
- ⟨ if squareroot is exact return double result 47 ⟩ Used in chunk 45.
- ⟨ if subtraction is exact return double result 37 ⟩ Used in chunk 32.
- ⟨ improve operands and compute p for add and sub 73 ⟩ Used in chunk 72.
- ⟨ initialize *z\_rep* as result of a binary operation 34 ⟩ Used in chunks 32 and 41.
- ⟨ private data and functions of type real 15, 25 ⟩ Used in chunk 7.
- ⟨ **real.c** 13 ⟩
- ⟨ **real.h** 7 ⟩
- ⟨ remove operand pointers 62 ⟩ Used in chunks 61, 66, and 68.
- ⟨ set error  $2^e$  75 ⟩ Used in chunks 72, 77, and 79.
- ⟨ set number to zero 61 ⟩ Used in chunks 57 and 66.
- ⟨ sign is clear from the bigfloat approximation 59 ⟩ Used in chunks 57, 58, and 60.