

# Automatic Generation of Staged Geometric Predicates

Aleksandar Nanevski, Guy Blelloch and Robert Harper <sup>\*</sup>  
*Carnegie Mellon University* (`{aleks,blelloch,rwh}@cs.cmu.edu`)

**Abstract.** Algorithms in Computational Geometry and Computer Aided Design are often developed for the Real RAM model of computation, which assumes exactness of all the input arguments and operations. In practice, however, the exactness imposes tremendous limitations on the algorithms – even the basic operations become uncomputable, or prohibitively slow. In some important cases, however, the computations of interest are limited to determining the sign of polynomial expressions. In such circumstances, a faster approach is available: one can evaluate the polynomial in floating point first, together with some estimate of the rounding error, and fall back to exact arithmetic only if this error is too big to determine the sign reliably. A particularly efficient variation on this approach has been used by Shewchuk in his robust implementations of Orient and InSphere geometric predicates.

We extend Shewchuk’s method to arbitrary polynomial expressions. The expressions are given as programs in a suitable source language featuring basic arithmetic operations of addition, subtraction, multiplication and squaring, which are to be perceived by the programmer as exact. The source language also allows for anonymous functions; the use of such functions enables the common functional programming technique of *staging*. The method is presented formally through several judgments that govern the compilation of the source expression into target code, which is then easily transformed into SML or, in case of single-stage expressions, into C.

**Keywords:** robust predicates, floating-point filters, exact arithmetic, program transformation, computational geometry

---

<sup>\*</sup> This research was conducted within PSciCo project at Carnegie Mellon University. The PSciCo project is supported by National Science Foundation (NSF) under the title “Advanced Languages for Scientific Computation Environments” as part of the Experimental Software Systems program within CISE. The grant number is 9706572.



**Keywords:** robust predicates, floating-point filters, exact arithmetic, program transformation, computational geometry

## 1. Introduction

Algorithms in Computational Geometry and Computer Aided Design are often created for the Real RAM model of computation. The Real RAM model assumes exactness of all the arguments and operations involved in the calculations, thus making it easy to justify the algorithms and prove their correctness. Unfortunately, this very fact implies that the computations have to be done with unbounded or infinite precision, which can render the basic operations and predicates prohibitively slow or even uncomputable.

A practical and very useful compromise, when applicable, is to assume that the input arguments are of floating-point type. It is also very common that the required functionality involves computation of only the *sign* of a given *polynomial* expression. Such calculations are, for example, used in the geometric predicates for determining whether a point is in/out/on a given line, circle, plane, sphere, etc. These predicates are, in turn, fundamental building blocks of algorithms for some basic geometric structures such as convex hulls and Delaunay triangulations.

However, floating-point alone is not sufficient to guarantee that the evaluation of a polynomial expression will correctly obtain its sign. The perturbations accumulated during the computation due to rounding errors, if sufficiently large, can change the sign of the final result. This can present the program with an inconsistent view of the data set, and cause it to produce incoherent results, diverge, or even crash. On the other hand, when the input coordinates are floating-point numbers, the exact sign can always be computed, albeit slowly, by first converting the floating-point arguments into rational numbers, and then carrying out the prescribed operations in rational arithmetic.

One method that has been proposed as an efficiency improvement to the exact rational arithmetic involves the use of *floating-point filters*. A floating-point filter carries out the given computation in floating-point first, together with some sort of estimate of the rounding error (using either forward error analysis or interval arithmetic), and falls back to exact arithmetic only if the estimated error is too big to reliably determine the sign [2, 13, 17, 14, 21, 4]. Thus, it “filters out” the easy computations whose sign can be quickly determined and only leaves the



© 2002 Kluwer Academic Publishers. Printed in the Netherlands.

hard ones for the exact arithmetic. A particularly efficient variation of this approach was described by Jonathan Shewchuk in his PhD thesis [24]. Aside from performing the floating-point part of the computation as the first phase of the filter, it introduces additional filtration phases of ever-increasing precision. The phases are attempted in order, each phase building on the result from the previous one, until the correct sign is obtained. This reuse of previous results induces a small cost in storing these results in memory, but that should be negligible and may be amortized when more precise phases of computation are executed. Notice, however, that this need not be true for every particular application, and there has been interesting subsequent work trying to further optimize and avoid the reuse of intermediate results [11].

Developing robust geometric predicates via filtering, in general, can be very cumbersome and error prone, and typically the most efficient methods are the hardest to implement. The difficulties usually arise from the need to perform the error analysis of the algebraic expression being evaluated for the sign. The solution is to use a semantics preserving expression compiler [13, 5, 12] which automates the error analysis of a given algebraic expression and generates an equivalent program implementing the corresponding floating-point filter.

With Shewchuk's method of multi-phase filtering there are even further difficulties. In this approach, the phases of computation are not described in a uniform way. *They do not consist of always running the same program with higher and higher precision.* Rather, the phases are separate and fundamentally different blocks of computation, which very intricately dependent on each other's intermediate results. Hence, this approach would benefit the most from automation by an expression compiler.

We are also interested in designing predicates for functional languages, and developing geometric applications of functional languages [3]. In a functional language we can exploit some well-known programming techniques such as staging and specialization to speed up the computation. For example, consider filtering a set of points to see on what side of a plane defined by three points they lie. This is a typical test used, for example, in the Quickhull algorithm [1]. The test can be staged by first forming the plane and its normal and then checking the position of each point from the set. A staged test obviates the need to repeat the part of the computation pertinent to the normal whenever a new point is tested, and can potentially save a lot of work. Such staging of programs is certainly desirable in computational geometry, as testified for example by the development of LOOK [14], but is more naturally exploited in functional programming languages. Unfortunately, the expression compilers available to date do not support

staging. LOOK supports a limited form of staging by allowing partial results to be calculated with various levels of accuracy lazily (i.e. on demand). It, however, is unlikely to allow the same level of control as can be achieved by a compiler since it does not optimize across operations.

This paper reports on an expression compiler that addresses these shortcomings. The input to the compiler is a function written in an appropriate source language offering the basic arithmetic operations of addition, subtraction, multiplication and squaring, and allowing for nested anonymous functional expressions (*stages*). All the operations in the source language are perceived as exact. The output of the compiler is a program in the target language designed to be easily converted into Standard ML (SML) or, in the case of single-stage programs, to C. The resulting SML or C program will determine the *sign* of the source function at the given floating-point arguments, using a floating-point filter with several *phases*, when exact computation needs to be performed. In particular, in the case of Shewchuk’s basic geometric predicates, the expression compiler will generate code that, to a considerable extent, reproduces that of Shewchuk.

The rest of the text is organized as follows. Section 2 summarizes the main ideas behind floating-point filters and arbitrary precision floating-point arithmetic. The source and target languages are presented in Section 3, and the program transformation process is described in Section 4. Performance comparison with Shewchuk’s predicates is given in Section 5.

## 2. Background

From here on we assume floating-point arithmetic as prescribed by the IEEE standard and the to-nearest rounding mode with the round-to-even tie-breaking rule [15]. We also assume that no overflows or underflows occur.

### 2.1. ERROR ANALYSIS

One of the most important properties of a floating-point arithmetic is the correct rounding of the basic arithmetic operations. It requires that the computed result always look as if it were first computed exactly, and then rounded to the number of bits determined by the precision of the arithmetic. For example, let us denote by  $\otimes$  the “rounded”, i.e. floating-point version of the exact operation  $*$   $\in \{+, -, \times\}$ . Then, if  $x$  and  $y$  are floating-point numbers, and  $x \otimes y$  has a normalized mantissa (i.e. is not

a denormalized floating-point number), a consequence of the correct rounding is that

$$|x * y - x \otimes y| \leq \epsilon |x \otimes y| \quad \text{and} \quad |x * y - x \otimes y| \leq \epsilon |x * y|$$

The quantity  $\epsilon$  in the above inequality is called “machine epsilon”. If  $m$  is the precision of the arithmetic, i.e. the number of bits reserved for the normalized mantissa (without the hidden leading bit), then  $\epsilon = 2^{-(m+1)}$ . In the IEEE standard for double precision, for example,  $\epsilon = 2^{-53}$ . By abuse of notation, the above inequalities are often stated respectively as

$$x * y = (1 \pm \epsilon)(x \otimes y) = x \otimes y \pm \epsilon |x \otimes y| \quad (1)$$

and

$$x * y = x \otimes y \pm \epsilon |x * y|$$

The equation (1) provides a bound on absolute error of the expression when the expression consists of only a single floating point operation. Notice that the error is composed of two multiples,  $\epsilon$  and  $|x * y|$ , the first of which does not depend on the arguments  $x$  and  $y$ . The rounding error for a *composite* expression can also be split into two multiples, one of which does not depend on the arguments of the expression. This part of the error need not be computed at run-time when all the arguments of the expression are supplied, but can rather be completely obtained while preprocessing the expression. To this end, assume that the exact values  $X_i$  are approximated in floating-point as  $x_i$  with absolute error  $\delta_i p_i$ , i.e. that for  $i = 1, 2$  we have

$$X_i = x_i \pm \delta_i p_i$$

For example, if  $X_i$  was obtained from a single operation over floating-point numbers  $x$  and  $y$ , i.e. if  $X_i = x * y$ , then it would be  $x_i = x \otimes y$ ,  $\delta_i = \epsilon$  and  $p_i = |x_i|$ . Assume in addition that the quantities  $\delta_i$  do not depend on any run-time arguments and that the invariant  $|x_i| \leq p_i$  holds. This is clearly true in the base case when  $x_i$  is obtained from a single operation on exact arguments, as can be seen from (1). The quantities  $\delta_i$  are rational numbers, and the values  $p_i$  are floating-point. Diverging slightly from the customary nomenclature, we call these two multiples respectively the “relative error” and the “permanent” of the approximation  $x_i$ .

Using the inequalities for rounded floating-point arithmetic from above, we can derive

$$\begin{aligned} |(X_1 + X_2) - (x_1 \oplus x_2)| &= \\ &= |(X_1 + X_2) - (x_1 + x_2) + (x_1 + x_2) - (x_1 \oplus x_2)| \end{aligned}$$

$$\begin{aligned}
&\leq |(X_1 + X_2) - (x_1 + x_2)| + |(x_1 + x_2) - (x_1 \oplus x_2)| \\
&\leq (\delta_1 p_1 + \delta_2 p_2) + \epsilon |x_1 \oplus x_2| \\
&\leq \max(\delta_1, \delta_2)(p_1 + p_2) + \epsilon(p_1 \oplus p_2) \\
&\leq \max(\delta_1, \delta_2)(1 + \epsilon)(p_1 \oplus p_2) + \epsilon(p_1 \oplus p_2) \\
&= \left(\epsilon + \max(\delta_1, \delta_2)(1 + \epsilon)\right)(p_1 \oplus p_2)
\end{aligned}$$

The above inequality is, by abuse of notation, customarily written as

$$X_1 + X_2 = x_1 \oplus x_2 \pm \left(\epsilon + \max(\delta_1, \delta_2)(1 + \epsilon)\right)(p_1 \oplus p_2)$$

The relative error of the composite expression  $X_1 + X_2$  is then  $\epsilon + \max(\delta_1, \delta_2)(1 + \epsilon)$  and its permanent is  $p_1 \oplus p_2$ . Notice that the relative error again does not depend on the run-time arguments, and that the invariant  $|x_1 \oplus x_2| \leq p_1 \oplus p_2$  is preserved. Similar derivations produce

$$\begin{aligned}
X_1 - X_2 &= x_1 \ominus x_2 \pm \left(\epsilon + \max(\delta_1, \delta_2)(1 + \epsilon)\right)(p_1 \oplus p_2) \\
X_1 X_2 &= x_1 \otimes x_2 \\
&\quad \pm \left(\epsilon + (\delta_1 + \delta_2 + \delta_1 \delta_2)(1 + \epsilon)\right)(p_1 \otimes p_2) \\
X_1^2 &= x_1 \otimes x_1 \pm \left(\epsilon + (2\delta_1 + \delta_1^2)(1 + \epsilon)\right)(p_1 \otimes p_1)
\end{aligned} \tag{2}$$

The above formulas provide a quick test for the sign of  $X_i$ . Obviously,  $x_i$  and  $X_i$  have the same sign if  $|x_i| > \delta_i p_i$ . However, this test is not completely satisfactory since it contains exact multiplication of a rational number  $\delta_i$  and a floating-point number  $p_i$ . A simpler, although less tight test is

$$|x_i| > \lceil (1 + \epsilon)\delta_i \rceil_{fp} \otimes p_i \tag{3}$$

where  $\lceil Q \rceil_{fp}$  denotes the smallest floating-point value above the rational number  $Q$ . This is indeed the inequality we use in our expression compiler to test the sign of an evaluated expression.

## 2.2. EXPANSIONS

Another important feature of IEEE round-to-nearest arithmetic is that the roundoff error of the basic operations is always representable as a floating-point number and can be recovered from the result and the arguments of the operation. A way to recover the rounding error was first proposed by Dekker [10], and then later another one was proposed by Knuth [16]. Analysis of both methods can be found in [23]. We present Knuth's method below.

**THEOREM 1 (Knuth).** *In floating-point arithmetic with precision  $m \geq 3$ , if  $x = a \oplus b$  and  $c = x \ominus a$  then  $a + b = x + ((a \ominus (x \ominus c)) \oplus (b \ominus c))$ .*

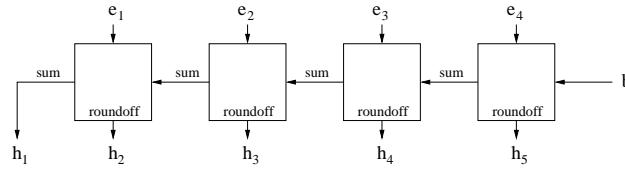


Figure 1. Adding a floating-point number to an expansion. The float  $b$  is added to the expansion  $e_1 + e_2 + e_3 + e_4$ , to produce a new expansion  $h_1 + \dots + h_5$ .

Knuth's theorem is significant because it provides a way to quickly perform *exact* addition of two floating-point numbers. First the addition  $x = a \oplus b$  is performed approximately, and then the roundoff error  $e = (a \ominus (x \ominus c)) \oplus (b \ominus c)$  is recovered. This takes only 6 floating-point operations, which is generally much faster than first converting  $a$  and  $b$  into rational numbers, and then adding them in rational arithmetic. The two values  $(x, e)$  put together represent the exact sum of  $a$  and  $b$ . One can view this pair as a sparse representation of the sum in a digit system with a radix  $2^{m+1}$ . Closing the set of sparse representations under addition leads to a very efficient data structure for exact computation. The values of this data type are lists of floating-point numbers sorted by magnitude and satisfying certain technical conditions about the alignment of their mantissas. These lists are called *expansions*, and each expansion represents the *exact* sum of its elements.

**DEFINITION 2 (Expansions).** *Expansion is a list of floating-point numbers  $\mathbf{x} = x_1 + \dots + x_{n-1} + x_n$ . The components of an expansion are ordered by decreasing magnitude and are nonoverlapping, i.e. the least-significant nonzero bit of  $x_i$  is more significant than the most significant nonzero bit of  $x_{i+1}$ . An expansion can contain components equal to zero, and these do not overlap with any other component.*

**EXAMPLE 1** In arithmetic with precision  $m = 4$ , the binary number

$$10101\ 11111\ 00\ 10010\ 00000\ 11$$

can be represented as the expansion

$$1.0110 \times 2^{23} - 1.0000 \times 2^{14} + 1.0010 \times 2^{11} + 1.1000 \times 2$$

For the sake of illustration, here we only picture the process of adding a floating-point number to an expansion (Figure 1) and of summing up two expansions (Figure 2). Quick algorithms for other basic arithmetic operations on this data type have been devised as well [20, 24, 7, 6, 9, 8].

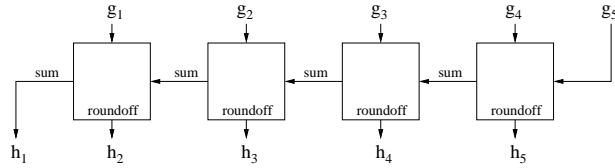


Figure 2. Summing two expansions. The components of  $e_1 + e_2 + e_3$  and  $f_1 + f_2$  are first merged by decreasing magnitude into the list  $[g_1, \dots, g_5]$ , which is then “normalized” into an expansion  $h_1 + \dots + h_5$ .

### 2.3. ADAPTIVE ARITHMETIC

Another consequence of Knuth’s theorem is a convenient ordering of operations that makes it possible to separate the computation into a sequence of filtering phases. Each phase is attempted after the previous one had failed to determine the sign reliably, and each computes with increasing precision, building on the result of the previous one. The following example is a bit contrived but illustrative nevertheless. Consider the expression  $X = (a_x - b_x)^2 + (a_y - b_y)^2$  where  $a_x, a_y, b_x$  and  $b_y$  are floating-point values. To find the sign of  $X$ , let  $v_x = a_x \ominus b_x$  and  $v_y = a_y \ominus b_y$ , and let  $e_x$  and  $e_y$  be the roundoffs from the two subtractions. Then

$$\begin{aligned} X &= (v_x + e_x)^2 + (v_y + e_y)^2 \\ &= (v_x^2 + v_y^2) + (2v_x e_y + 2v_y e_x) + (e_x^2 + e_y^2) \end{aligned}$$

In this sum, the summand  $v_x^2 + v_y^2$  is dominant, since  $|e_x| \leq \epsilon|v_x|$  and  $|e_y| \leq \epsilon|v_y|$  by (1). It is then a good heuristic to first compute  $v_x^2 + v_y^2$  and test it for sign before proceeding, because it is likely that  $v_x^2 + v_y^2$  will already have the same sign as  $X$ . The process can be sped up even more if this expression is first computed approximately to obtain  $X^A = (v_x \otimes v_x) \oplus (v_y \otimes v_y)$ . Then only if  $X^A$  has too big an error bound, as determined by the test (3), the computation of  $X^B = v_x^2 + v_y^2$  is undertaken exactly using the data type of expansions. As depicted in Figure 3, the computation of  $X^B$  does not need to start from scratch, but can reuse the approximate values of  $(v_x \otimes v_x)$  and  $(v_y \otimes v_y)$  and refine them with the recovered rounding errors to obtain the exact  $v_x^2$  and  $v_y^2$ . If  $X^B$  is also too crude an approximation of  $X$ , we can correct it by adding up the smaller terms  $(2v_x e_y + 2v_y e_x)$ , first approximately, and then correctly. Finally, if all of these approximations fail to give an answer, we can compute the exact result by adding the last summand  $e_x^2 + e_y^2$  to the expansion computed in the previous phase. Using this approach, we will compute the exact value only if absolutely necessary,



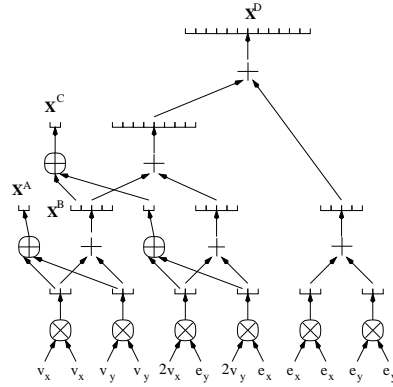


Figure 3. Evaluation of  $X$  in phases of increased precision with reuse of intermediate results.

and even then, the efforts spent on previous phases will not be wasted, but will rather be reused to obtain the exact result in an efficient way.

This idea to generalize floating-point filters into a hierarchy of adaptive precision filtering phases is due to Shewchuk. While the number and type of adaptive phases, strictly speaking, can vary with the expression, his experiments pointed to a scheme with four phases (depicted in Figure 3) as being optimal in practice for the basic geometric predicates that he considered. We adopt this scheme and present its formalization in Section 4. The arbitrary precision floating-point arithmetic and the data type of expansions is invented by Priest, later optimized by Shewchuk, and further on by Daumas. Detailed description and analysis of adaptive precision arithmetic and of the algorithms for basic operations can be found in [19, 20, 23, 24, 7, 6, 9, 8].

The whole method described above relies on the fact that the required floating-point operations will execute without any exceptions, i.e. that neither overflow nor underflow will occur during the computation. If exceptions do occur, the expansions holding the exact intermediate values may lose bits of precision and produce a distorted answer. A possible solution in such cases is to rerun the computation in some other, slower, form of exact arithmetic (for example in infinite precision rational numbers).

### 3. Source and target languages

The source language of the expression compiler is shown in Figure 4. Its syntax supports the basic arithmetic operations (including squaring), assignments and staged functional expressions. The arguments of

phrases	$\phi ::= x \mid c \mid e$
expressions	$e ::= \phi_1 + \phi_2 \mid \phi_1 - \phi_2 \mid \phi_1 \times \phi_2$ $\mid \sim\phi \mid \text{sq } \phi$
assignment lists	$\alpha ::= \text{val } x = e \mid \text{val } x = e \alpha$
programs	$\pi ::= \text{fn } [x_1, \dots, x_n] \Rightarrow \text{let } \alpha \text{ end}$ $\mid \text{fn } [x_1, \dots, x_n] \Rightarrow \text{let } \alpha \pi \text{ end}$

Figure 4. Source language

$$\text{orient2}(A, B, C) = \text{sgn} \begin{vmatrix} a_x - c_x & b_x - c_x \\ a_y - c_y & b_y - c_y \end{vmatrix}$$

```

fn [ax, ay, cx, cy] =>
let val acx = ax - cx
    val acy = ay - cy

    fn [ bx, by ] =>
let val d = acx × ( by - cy ) -
        acy × ( bx - cx )
end
end

```

Figure 5. Orient2D predicate: definition and implementation in the source language.

the functions should be perceived as floating-point values, while the intermediate results are assumed to be computed exactly. Squaring is included among the arithmetic operations because it can often be executed more quickly than the multiplication of two equal exact values, and has a better error bound. In addition, it provides the compiler with the knowledge that its result is non-negative, which can be used in some cases to optimize the code. In order to simplify the compilation process, the source language requires that all the assignments are non-trivial, i.e. it disallows assignments to variables of other variables or constants. The syntactic category  $\pi$  declares source programs with multiple stages of computation. A program defined in the source language is designed to compute the *sign* of the last expression in the assignment list of the program's last stage.

As an example of a program in the source language, consider the Orient2D geometric predicate and its implementation in Figure 5. Orient2D determines the position (in/out/on) of point  $B = (b_x, b_y)$  with

```

reals       $r ::= x \mid c \mid r_1 * r_2 \mid r_1 \otimes r_2 \mid \text{sq } r$ 
            $\mid \sim r \mid \text{abs } r \mid \text{double } r \mid \text{approx } r$ 
            $\mid \text{tail}_*(r_1, r_2, r_3) \mid \text{tail}_{\text{sq}}(r_1, r_2)$ 
assignment lists  $\lambda ::= \text{val } (x_1, \dots, x_n) = \text{lforce } x \lambda$ 
            $\mid \text{val } (x_1, \dots, x_n) = \text{rforce } x \lambda$ 
            $\mid \text{val } x = \text{susp } \lambda \text{ in}$ 
            $\quad ((x_1, \dots, x_n),$ 
            $\quad (x_1, \dots, x_m))$ 
            $\text{end } \lambda$ 
            $\mid \text{val } x = r \lambda \mid \text{empty}$ 
sign tests  $\sigma ::= \text{sign } r \mid \text{signtest } (r_1 \pm r_2)$ 
            $\text{with } \lambda \text{ in } \sigma \text{ end}$ 
functions  $\varphi ::= \text{fn } (x_1, \dots, x_n) =>$ 
            $\text{let } \lambda \text{ in } \sigma \text{ end}$ 
            $\mid \text{fn } (x_1, \dots, x_n) =>$ 
            $\text{let } \lambda \text{ in } \varphi \text{ end}$ 

```

Figure 6. Target language.

respect to the line from  $A = (a_x, a_y)$  to  $C = (c_x, c_y)$ . The presented implementation is staged in the coordinates of  $A$  and  $C$ . Once the predicate is applied to these two points, its result is a new function specialized to compute relative to the line  $\overline{AC}$ , without recomputing the intermediate results `acx` and `acy`.

The target language of the compilation is presented in Figure 6. It is designed to be easily converted to SML or C, so its semantics is best explained by referring to these languages. In the syntactic category of reals, the symbol `*` varies over the operations  $\{+, -, \times\}$ , and  $\sim$  is the unary negation operator. In the translation to SML or C, the values of the syntactic category of reals will be converted either into floating-point numbers or a library implementation of expansions. However, we chose not to make this distinction explicit and did not introduce separate types for floats and expansions in the target language. The reason is that we do not plan to do any programming in this language directly, but rather just use it for intermediate representation of programs before they are converted into SML or C.

The target language operations  $\oplus$ ,  $\ominus$  and  $\otimes$  will be interpreted as corresponding SML or C floating-point operations. They expect floating-point input, and produce floating-point output. The exact target-language operations `+`, `-` and `×` will be translated into the appropriate

exact operations on the type of expansions in SML or C. Constants are always floating-point values. The `tail` constructs compute the roundoffs from their corresponding floating-point operation. For example, `tail+(a, b, a ⊕ b)` will compute the roundoff from the addition  $a \oplus b$ , following Knuth’s theorem. The construct `double` is multiplication by 2 on expansions, and `approx` returns a floating-point number approximating the expansion argument with a relative error of  $2\epsilon$ .

To describe the role of the `susp`, `lforce` and `rforce` constructs, we need to make a clear distinction between stages and phases of computation in the target language. The source program contains nested functional expressions which we refer to as *stages*, and the computation advances through the stages when a partial input is supplied to the predicate. Once a source program is compiled, every stage gets transformed into a stage of the target language, which consists of four computational *phases* of increased precision. The first phase carries out the computation in floating-point. The other phases mix in elements of exact computation and are called only if higher precision is required. Because the results of the later phases may never actually be needed, these phases should not be carried out immediately as a partial instantiation reaches their stage. Rather, their corresponding computations should be *suspended* and *memoized* in order to be forced when and if required. Thus, one may say that the notion of stages refers to *partial evaluation* of code, while the notion of phases refers to *lazy evaluation* of code.

Going back to the target language, `susp` creates a piece of code, a *suspension*, to be evaluated when requested by `rforce` or `lforce`. As explained before, it provides a mechanism to pass intermediate results between different stages, and between different phases of the same stage. The output from a suspension contains two lists of intermediate values. The first list consists of intermediate values intended for some later phase of the current stage, and the second list consists of intermediate values intended for the following stage. The first list can be recovered by `lforce`-ing the suspension, and the second list by `rforce`-ing it (see Figure 7).

The `sign` function returns the sign of an expansion. The construct `signtest` first checks whether the magnitude  $|r_1|$  of the tested value is bigger than the magnitude  $|r_2|$  of the roundoff error. If so, it returns the sign of  $r_1$ . Otherwise, it cannot determine the sign of  $r_1$  with certainty, so it undertakes the computation of the next phase  $\lambda$ , followed by sign test  $\sigma$ . Values  $r_1$  and  $r_2$  are assumed to be floating-point.

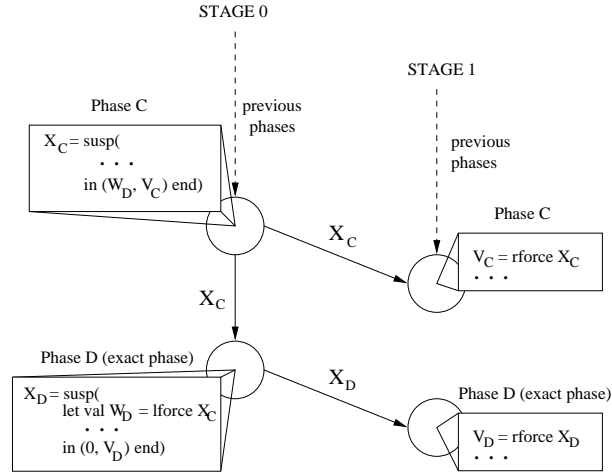


Figure 7. Passing intermediate results between phases and stages using `susp`, `lforce` and `rforce`.

#### 4. Compilation

To describe the compilation process, first notice that the source program, according to the grammar of the source language (Figure 4), can be viewed as a nonempty sequence of assignment lists, each representing a single stage of computation. Each of these stages is separately compiled into four target phases that are meant to perform the computation of the stage with increasing precision, as described in Section 2. At the end, these pieces of target code are pasted together in a target program, according to specific templates, so that sign checks are performed between subsequent phases, while respecting the staging specified in the source program.

The whole process is formalized using five judgments – four for compiling source stages into their target counterparts, and one judgment to compose all the obtained target stages and phases together into a target program. This section describes the compilation process in more detail, explains some decisions in designing the compilation judgments and illustrates representative rules of the judgments through several examples.

Before proceeding further, there is one technicality to notice. Namely, it can be assumed, without loss of generality, that the source program to be compiled is in a very specific format. First, we require that all of its assignments consist of a *single* binary operation acting on two other *variables* (rather than on two arbitrary source expressions), or of a single unary operation acting on another variable. The second

requirement is that the source program does not contain any floating-point constants – all the constants are replaced by fresh variables for error analysis purposes, and then put back into the target code at the end of the compilation. It is trivial to transform the source program so that it complies with these two prerequisites, so we do not present the formalization of this procedure. In our implementation it is carried out in parallel with the parsing of the source program.

In the following section, we illustrate the various phases of the process with the compilation of the source expression  $F$  given below.

```

F = fn [a, b] =>
    let val ab = a - b
        val ab2 = sq ab

        fn [c] =>
            let val d = c × ab2
            end
        end
    end

```

It consists of two stages: first expecting two arguments  $a$  and  $b$ , and the second expecting an argument  $c$ .

#### 4.1. COMPILING THE STAGES

This section illustrates the way we formalize the error analysis and the program transformation of the individual stages of the source program. The first phase, phase A, of the source program performs all the source operations approximately in floating-point. The transformation corresponding to the phase A follows closely the equations (2). We compute the relative error (which is a rational number known at compile-time), and emit target code for computing the permanent and for testing the reliability of the result sign. In order to describe the compilation for the phase A, we rely on the judgment

$$E_1 \vdash_A \alpha \rightsquigarrow \lambda; r_1, r_2 / E_2$$

This judgment relates a list  $\alpha$  of source assignments to the list  $\lambda$  of corresponding phase A target assignments. Expressions  $r_1$  and  $r_2$  are from the syntactic category of reals in the target language (Figure 6). The expression  $r_1$  is to be tested for sign at the end of the phase, and the expression  $r_2$  is an upper bound on the roundoff error. The assignments in  $\lambda$  will perform the phase A calculations and compute the appropriate permanent.

The contexts  $E_1$  and  $E_2$  deserve special attention. They are sets relating target language variables with their error estimates. The gram-

mar for their generation is presented below.

$$\begin{aligned} \text{contexts } E &::= \cdot \mid x : \tau, E \mid x \triangleright r, E \\ \text{errors } \tau &::= \mathcal{O}_A(\delta) \mid \mathcal{O}_B(\delta) \mid \mathcal{O}_C(\delta, \iota, \rho) \mid \mathcal{O}_D \mid \mathcal{P} \end{aligned}$$

Each variable in a context is bound in one of the four phases of the computation (A, B, C or D), and will have error estimates that are appropriate for that phase of the computation ( $\mathcal{O}_A(\delta)$ ,  $\mathcal{O}_B(\delta)$ ,  $\mathcal{O}_C(\delta, \iota, \rho)$  and  $\mathcal{O}_D$ ), where  $\delta$ ,  $\iota$  and  $\rho$  are rational numbers (Figure 7). For example, if the error relation  $x : \mathcal{O}_A(\delta) \in E$ , that means that the variable  $x$  which is bound in phase A, has been estimated by the compiler to have a relative error *bounded* from above by the rational number  $\delta$ . Similar meaning can be ascribed to the error assignments  $x : \mathcal{O}_B(\delta)$  for phase B. Phase C, on the other hand, is a mix of approximate and exact computations, and there are three rational values  $\delta$ ,  $\iota$  and  $\rho$  that govern phase C error estimations. We do not describe their meaning in this paper, but the formulas for their derivation can be found in the companion technical report [18]. Phase D is the exact phase, so there are no error estimates to associate with phase D variables. Finally, the temporary variables introduced to hold parts of the permanent are not analyzed for error. We still place them into the error contexts, just for clarity, but with the error tag  $\mathcal{P}$ . To reduce clutter, the error estimate of a variable  $x$  in a context  $E$  will be denoted simply as  $E(x)$ , as it will always be clear from the rule in which phase the variable is bound.

In addition to the error estimates, the contexts contain substitutions of variables by target language real expressions ( $x \triangleright r$ ). If some compilation rule needs to emit into target code a variable for which there is a substitution in the context, the substituting expression will be emitted instead. This serves two purposes. First, we can use it to express that certain variables in the code are just placeholders for floating-point constants – a situation occurring, as explained before, because of an assumed stricter form of the source programs. Second, it lets us optimize, in a single pass of the compiler, the code for computing the permanent of the expression – a process that will be illustrated below.

Now that we have laid out the structure of the contexts  $E_1$  and  $E_2$  in the judgment we are defining, we can describe their purpose. Simply, the compilation with the judgment starts with the context  $E_1$ , and ends with  $E_2$ . So,  $E_2$  is in fact  $E_1$  enlarged with the new variables, error estimates and substitutions introduced during the compilation. The context  $E_2$  is returned so that it can be threaded into other rules.

Going back to the analysis for the expression  $F$ , we illustrate how its phase A can be compiled using the above judgment. First of all, the expression  $F$  is specified as two stages: the one executing `val ab = a - b` and `val ab2 = sq ab`, and the other one executing `val d = c × ab2`.

The compilation for phase A starts by breaking down each stage of the source program into individual assignments. The rule is the following.

$$\frac{E_1 \vdash_A \text{val } \mathbf{x} = e \rightsquigarrow \lambda_H; s_1, s_2 / E' \quad E' \vdash_A \alpha \rightsquigarrow \lambda_T; r_1, r_2 / E_2}{E_1 \vdash_A \text{val } \mathbf{x} = e \alpha \rightsquigarrow \lambda_H \lambda_T; r_1, r_2 / E_2}$$

The rule “folds” the functionality of the compiler across the list of source assignments, carrying the context from one assignment to the next. Notice that the expressions  $s_1$  and  $s_2$  are never used – only the last expression in the assignment list is ever tested for sign. Now, to compile the assignment  $\text{val } \mathbf{ab} = \mathbf{a} - \mathbf{b}$ , we need a rule applicable to subtraction of input arguments. Input arguments are assumed to be error-less, so the following rule applies.

$$\frac{E(x_1^A) = E(x_2^A) = 0}{E \vdash_A \text{val } y = x_1 - x_2 \rightsquigarrow \text{val } y^A = x_1^A \ominus x_2^A; y^A, 0 / E, y^A : \mathcal{O}_A(\epsilon), y^P : \mathcal{P}, y^P \triangleright \text{abs}(y^A)}$$

When applied to the assignment  $\text{val } \mathbf{ab} = \mathbf{a} - \mathbf{b}$ , the meta variables  $y$ ,  $x_1$  and  $x_2$  are instantiated to  $\mathbf{ab}$ ,  $\mathbf{a}$  and  $\mathbf{b}$  respectively. The rule then emits the target code for the assignment to  $\mathbf{ab}^A$  (the superscripts  $A$  indicate that the target variable is bound in the phase A of the predicate). For the purpose of bookkeeping, the rule must also extend the context with information about the relative error and the permanent of  $\mathbf{ab}^A$ . The relative error of  $\mathbf{ab}^A$  is  $\epsilon$ , so the rule generates the error estimate  $\mathbf{ab}^A : \mathcal{O}_A(\epsilon)$ . Finally, the permanent  $\mathbf{ab}^P$  of  $\mathbf{ab}^A$  is equal to  $|\mathbf{ab}^A|$ , and so the substitution context is extended with  $\mathbf{ab}^P \triangleright \text{abs}(\mathbf{ab}^A)$ .

Next in the assignment list from our example is the assignment  $\text{val } \mathbf{ab2} = \text{sq } \mathbf{ab}$ . The squaring operation is handled by the rule

$$\frac{x_1^P \triangleright x_1^A \text{ or } x_1^P \triangleright \text{abs}(x_1^A) \in E}{E \vdash_A \text{val } y = \text{sq}(x_1) \rightsquigarrow \text{val } y^A = x_1^A \otimes x_1^A; y^A, \lceil \frac{(1+\epsilon)^2(2\delta_1 + \delta_1^2)}{1-\epsilon} \rceil_{fp} \otimes y^A / E, y^A : \mathcal{O}_A(\epsilon + (1+\epsilon)(2\delta_1 + \delta_1^2)), y^P : \mathcal{P}, y^P \triangleright y^A}$$

In the assignment to  $\mathbf{ab2}$ , the meta variables  $x_1$  and  $y$  of this rule are instantiated to  $\mathbf{ab}$  and  $\mathbf{ab2}$  respectively. Then the meta variable  $x_1^P$  becomes  $\mathbf{ab}^P$ . But  $\mathbf{ab}^P$  has already been introduced into the context with a substitution  $\mathbf{ab}^P \triangleright \text{abs}(\mathbf{ab}^A)$ . Thus, the premises of this rule are satisfied, and it can be applied. The meta variable  $\delta_1$  from the rule refers to the relative error of the variable  $x_1^A$  as read from the context, i.e.  $\mathcal{O}_A(\delta_1) = E(x_1^A)$ . In our example of assignment to  $\mathbf{ab2}$ , the variable



Table I. Compilation of the first stage of  $F$ .

source code	target code	context
phase A		
<code>val ab = a - b</code>	<code>val ab<sup>A</sup> = a<sup>A</sup> ⊖ b<sup>A</sup></code>	$\text{ab}^A : \mathcal{O}_A(\epsilon)$ $\text{ab}^P \triangleright \text{abs}(\text{ab}^A)$
<code>val ab2 = sq ab</code>	<code>val ab2<sup>A</sup> = ab<sup>A</sup> ⊗ ab<sup>A</sup></code>	$\text{ab2}^A : \mathcal{O}_A(3\epsilon + 3\epsilon^2 + \epsilon^3)$
phase B		
<code>val ab = a - b</code>	–	$\text{ab}^B : \mathcal{O}_B(\epsilon)$ $\text{ab}^B \triangleright \text{ab}^A$
<code>val ab2 = sq ab</code>	<code>val ab2<sup>B</sup> = sq ab<sup>A</sup></code>	$\text{ab2}^B : \mathcal{O}_B(2\epsilon + 3\epsilon^2 + \epsilon^3)$
phase C		
<code>val ab = a - b</code>	<code>val ab<sup>C</sup> = tail<sub>-</sub>(a<sup>A</sup>, b<sup>A</sup>, ab<sup>A</sup>)</code>	$\text{ab}^C : \mathcal{O}_C(0, \epsilon, 0)$
<code>val ab2 = sq ab</code>	<code>val ab2<sup>C</sup> = double(ab<sup>A</sup> ⊗ ab<sup>C</sup>)</code>	$\text{ab2}^C : \mathcal{O}_C(3\epsilon^2 + 3\epsilon^3,$ $2\epsilon \frac{1+\epsilon}{1-\epsilon}, \epsilon)$
phase D		
<code>val ab = a - b</code>	–	$\text{ab}^D \triangleright \text{ab}^C$
<code>val ab2 = sq ab</code>	<code>val ab2<sup>D</sup> = double(ab<sup>A</sup> × ab<sup>C</sup>)</code> <code>+ sq(ab<sup>C</sup>)</code>	–

$x_1^A$  is instantiated to  $\text{ab}^A$  and  $\delta_1$  is instantiated to  $\epsilon$ . The produced permanent for  $\text{ab2}^A$  is  $\text{ab2}^A$  itself, explaining why we avoided emitting any code for permanent computation so far. Any separate computation of the permanent for  $\text{ab2}$  would have been just a waste of effort, since it is already computed. For future use, however, this rule stores the substitution  $\text{ab2}^P \triangleright \text{ab2}^A$  into the context. The relative error for  $\text{ab2}$  is computed as  $\epsilon + (1 + \epsilon)(2\delta_1 + \delta_1^2) = (3\epsilon + 3\epsilon^2 + \epsilon^3)$  and is stored into the context.

That finishes the compilation of phase A of the first stage. The second stage contains only the assignment `val d = c × ab2`, and its phase A target code is obtained by the rule

$$\begin{array}{c}
 E(x_1^A) = 0 \quad x_2^P \triangleright x_2^A \text{ or } x_2^P \triangleright \text{abs}(x_2^A) \in E \\
 \hline
 E \vdash_A \text{val } y = x_1 \times x_2 \rightsquigarrow \text{val } y^A = x_1^A \otimes x_2^A; \\
 y^A, \lceil \frac{(1+\epsilon)^2 \delta_2}{1-\epsilon} \rceil_{fp} \otimes \text{abs}(y^A) / \\
 E, y^A : \mathcal{O}_A(\epsilon + (1 + \epsilon)\delta_2), y^P : \mathcal{P}, \\
 y^P \triangleright \text{abs}(y^A)
 \end{array}$$

The rule compiles the source assignment into `val dA = cA ⊗ ab2A` and expands the current context with the error estimate  $\text{d}^A : \mathcal{O}_A(4\epsilon + 6\epsilon^2 + 4\epsilon^3 + \epsilon^4)$  and the substitution  $\text{d}^P \triangleright \text{abs}(\text{d}^A)$ .

Table II. Compilation of the second stage of  $F$ . The source code for this stage consists of the single assignment `val d = c × ab2`.

target code	testing value	error estimate	context
phase A			
<code>val d<sup>A</sup> = c<sup>A</sup> × ab2<sup>A</sup></code>	$\mathbf{d}^A$	$\lceil \frac{(1+\epsilon)^2(3\epsilon+3\epsilon^2+\epsilon^3)}{1-\epsilon} \rceil_{fp}$ $\otimes \mathbf{abs}(\mathbf{d}^A)$	$\mathbf{d}^A : \mathcal{O}_A(4\epsilon + 6\epsilon^2 + 4\epsilon^3 + \epsilon^4)$ $\mathbf{d}^P \triangleright \mathbf{abs}(\mathbf{d}^A)$
phase B			
<code>val d<sup>B</sup> = c<sup>A</sup> × ab2<sup>B</sup></code>	<code>approx(d<sup>B</sup>)</code>	$\lceil \frac{(1+\epsilon)^2(2\epsilon+3\epsilon^2+\epsilon^3)}{1-2\epsilon} \rceil_{fp}$ $\otimes \mathbf{abs}(\mathbf{d}^A)$	$\mathbf{d}^B = \mathcal{O}_B(2\epsilon + 5\epsilon^2 + 4\epsilon^3 + \epsilon^4)$
phase C			
<code>val d<sup>C</sup> = c<sup>A</sup> × ab2<sup>C</sup></code>	$\mathbf{d}^C$	$\lceil \frac{5\epsilon^2+12\epsilon^3+6\epsilon^4-4\epsilon^5-3\epsilon^6}{(1-\epsilon)^2} \rceil_{fp}$ $\otimes \mathbf{abs}(\mathbf{d}^A)$	$\mathbf{d}^C : \mathcal{O}_C(\frac{5\epsilon^2+2\epsilon^3-3\epsilon^4}{1-\epsilon}, 2\epsilon(\frac{1+\epsilon}{1-\epsilon})^2, 2\epsilon + \epsilon^2)$
phase D			
<code>val d<sup>D</sup> = c<sup>A</sup> × ab2<sup>D</sup></code>	$\mathbf{d}^B + \mathbf{d}^D$	–	–

The remaining phases for  $F$  are obtained in a similar way. The reader is referred to [18] for complete definition. The steps in the derivation for the two stages, including the changes in the judgment contexts, are presented in Table I and Table II respectively. In addition to the target code and the contexts, Table II also shows, for each of the phases, the testing value and error estimate (recall that only the testing value and the error estimate of the last stage are actually emitted into the target code). As can be seen from Table II, the testing values for the four phases of the second stage are  $\mathbf{d}^A$ , `approx(dB)`,  $\mathbf{d}^C$  and  $\mathbf{d}^B + \mathbf{d}^D$ , respectively. In the first three phases, these will be checked against the rounding errors to determine if they have the correct sign. In phase D, the testing value is actually the exact value of the expression. The error estimates for the second stage are obtained from the corresponding rounding errors using (3), producing a quick floating-point test for the sign of the testing value. The error estimates are represented in the table in a symbolic form. It is important to notice that all of them are known at compile time, and are emitted into target code as floating-point constants<sup>1</sup>. So, for example,  $\lceil \frac{(1+\epsilon)^2(3\epsilon+3\epsilon^2+\epsilon^3)}{1-\epsilon} \rceil_{fp} = 3.33067\text{e-}16$  and  $\lceil \frac{(1+\epsilon)^2(2\epsilon+3\epsilon^2+\epsilon^3)}{1-2\epsilon} \rceil_{fp} = 2.22045\text{e-}16$ .

<sup>1</sup> In the actual SML and C implementations, these values are calculated in an initialization routine, rather than placed in the code as decimal constants, in order to avoid rounding errors in the decimal-to-binary conversion.

```

fn [aA, bA] =>
let val abA = aA ⊖ bA
    val ab2A = abA ⊗ abA
    val yB =
      susp
        val ab2B = sq abA
        in ((), ab2) end
    val yC =
      susp
        val abC =
          tail_(aA, bA, abA)
        val ab2C =
          double(abA ⊗ abC)
        in ((abC, (ab2C)) end
    val yD =
      susp
        val (abC) = lforce yC
        val ab2D =
          double(abA × abC)
            + sq abC
        in ((), ab2D) end
in
fn [cA] =>
let val dA = cA ⊗ ab2A
in
    signtest (dA ± (3.33067e-16
                  ⊗ abs(dA)))
  with
    val (ab2B) = rforce yB
    val dB = cA × ab2B
    val yBX = approx(dB)
  in signtest
    (yBX ± (2.22045e-16
              ⊗ abs(dA)))
  with
    val (ab2C) =
      rforce (yC)
    val dC = cA ⊗ ab2C
  in signtest
    (dC ± (2.22045e-16
            ⊗ yBX ⊕
            6.16298e-32
            ⊗ abs(dA)))
  with
    val (ab2D) =
      rforce yD
    val dD = cA × ab2D
  in sign(dB + dD) end
end

```

Figure 8. Target code for the example expression  $F$ .

## 4.2. COMPILING THE PROGRAM

Once all the stages of the source program have been compiled, they need to be pasted together into a target program, in such a way that the phases can “communicate” their intermediate results. For illustration, the target code resulting from the compilation of the expression  $F$  is presented in Figure 8.

The translation is done through a new judgment

$$E_1 \vdash_P \pi; x_B, x_C, x_D \rightsquigarrow \varphi / V_B, V_C, V_D$$

which takes a source program  $\pi$  and compiles it into a target program  $\varphi$ . This judgment works in a bottom-up manner – the later stages are pasted in first (recall that a source program is a “list” of stages; the judgment first processes the tail of the list, and then pastes in the head stages). Thus, it is possible that the target term  $\varphi$  will not have all of its variables bound – some of them might have been introduced in one of the previous stages, and thus will be compiled and bound by the judgment only later. The meta variables  $x_B$ ,  $x_C$  and  $x_D$  hold object-code variables, freshly allocated in the previous stage to hold that stage’s suspensions, and then passed to the current stage to be

`rforce`'d if needed. The variables  $V_B$ ,  $V_C$  and  $V_D$  hold the object-code variables that the mentioned suspensions should populate with intermediate values. They are passed back to the previous stage so that the stage can be correctly constructed.

To determine which object-code variables will be passed via suspensions to a particular phase, we use the following function.

$$\text{fv}(\lambda, S) = (S \cup \text{free variables of } \lambda) \setminus \text{bound variables of } \lambda$$

For example, if  $\lambda_D$  is the assignment list for the exact phase (phase D) of the *last* stage in a program  $\pi$ , its free variables will be  $V_D = \text{fv}(\lambda_D, \emptyset)$ . Some of these free variables will be bound in the  $\lambda_A$ ,  $\lambda_B$  or  $\lambda_C$  list of the same stage, but some will have to be passed by a suspension from the exact phase of the previous stage (see Figure 6). The variables to be placed in this suspension are therefore all characterized by the fact that they are introduced in the exact phase of some previous stage. Thus, their set is  $V_D \cap \text{dom}_D E$ , where  $\text{dom}_D E$  is the set of variables from the context  $E$  that have phase D error estimates.

We can similarly determine the suspensions for the phase C of the last stage. Since phase C needs to bind some of the variables from  $\lambda_D$ , we don't just consider the free variables of  $\lambda_C$ , but rather set  $V_C = \text{fv}(\lambda_C, V_D)$ . As before, some of these variables will be bound in  $\lambda_A$  and  $\lambda_B$ , but those that are not will need to be passed via suspension from the phase C of the preceding stage. These variables are in the set  $V_C \cap \text{dom}_C E$ , where  $\text{dom}_C E$  is, analogously to the phase D case, the set of object-code variables from context  $E$  bound in some of the previous C phases.

In a similar way, the phase B will request the set  $V_B \cap \text{dom}_B E$  where  $V_B = \text{fv}(\lambda_B, V_C)$  passed as a suspension from phase B of the preceding stage. Finally, phase A doesn't require any variable passing, since the computations of this phase are always carried out immediately in each stage, and are never suspended.

The above discussion motivates the following rule of the  $\vdash_P$  judgment. The rule applies only if  $\pi$  is a single-stage program, and since the judgment is recursively applied, it serves to compile the *last* stage of the source program.

$$\frac{\begin{array}{l} E, x_i^A : \mathcal{O}_A(0) \vdash_A \alpha \rightsquigarrow \lambda_A; r_1^A, r_2^A / E_1 \quad E_1 \vdash_B \alpha \rightsquigarrow \lambda_B; r_1^B, r_2^B / E_2 \\ E_2 \vdash_C \alpha \rightsquigarrow \lambda_C; r_1^C, r_2^C / E_3 \quad E_3 \vdash_D \alpha \rightsquigarrow \lambda_D; r_1^D / E_4 \end{array}}{E \vdash_P \text{fn } [x_1, \dots, x_n] \Rightarrow \text{let } \alpha \text{ end}; x_B, x_C, x_D \rightsquigarrow \Phi / V_B \cap \text{dom}_B E, V_C \cap \text{dom}_C E, V_D \cap \text{dom}_D E}$$

where  $\Phi$  is defined as

```

fn [x1, ..., xn] =>
let λA
in signtest (r1A ± r2A) with
  val (VB ∩ domB E) =
    rforce(xB)
  λB
  val yBX = r1B
in signtest (yBX ± r2B) with
  val (VC ∩ domC E) =
    rforce(xC)
  λC
in signtest
  (r1C ± ⌊ $\frac{2\epsilon(1+\epsilon)^2}{1-\epsilon}$ ⌋fp ⊗ yBX
  ⊕ r2C)
  with
  val (VD ∩ domD E) =
    rforce(xD)
  λD
  in sign (r1D) end
end

```

Similar analysis of variable passing can be performed if the stage considered is not the last one. Then one only needs to take into account that some object-code variables might be requested from the subsequent stage, and factor them in when creating the suspensions. The rule that handles this case is

$$\frac{
\begin{array}{l}
E, x_i^A : O(0) \vdash \alpha \rightsquigarrow \lambda_A; r_1^A, r_2^A / E_1 \quad E_1 \vdash_B \alpha \rightsquigarrow \lambda_B; r_1^B, r_2^B / E_2 \\
E_2 \vdash_C \alpha \rightsquigarrow \lambda_C; r_1^C, r_2^C / E_3 \quad E_3 \vdash_D \alpha \rightsquigarrow \lambda_D; r_1^D / E_4 \\
E_4 \vdash_P \pi; y_B, y_C, y_D \rightsquigarrow \varphi / U_B, U_C, U_D
\end{array}
}{
\begin{array}{l}
E \vdash_P \text{fn } [x_1, \dots, x_n] \Rightarrow \text{let } \alpha \pi \text{ end}; x_B, x_C, x_D \\
\rightsquigarrow \Phi / V_B \cap \text{dom}_B E, V_C \cap \text{dom}_C E, V_D \cap \text{dom}_D E
\end{array}
}$$

where  $V_D = \text{fv}(\lambda_D, U_D)$ ,  $V_C = \text{fv}(\lambda_C, U_C \cup V_D)$ ,  $V_B = \text{fv}(\lambda_B, U_B \cup V_C)$ , and the program  $\Phi$  is defined as follows.

```

fn [x1, ..., xn] =>
let λA
  val yB =
    susp
    val (VB ∩ domB E) =
      rforce xB
    λB
    in (VD ∩ domB E, UB)
  end
  val yC =
    susp
    val (VC ∩ domC E) =
      rforce xC
    λC
    in (VD ∩ domC E, UC)
  end
  val yD =
    susp
    val (VD ∩ domD E) =
      rforce xD
    val (VD ∩ domB E) =
      lforce yB
    val (VD ∩ domC E) =
      lforce yC
    λD
    in (( ), UD) end
  in ((), UD) end
end

```

Finally, if  $\pi$  is a source program, then as described before, it can be assumed that all its assignment expressions consist of a single operation acting only on *variables*, and that its constants  $c_i$  are replaced by free variables  $y_i$ . The target program  $\varphi$  for  $\pi$  is obtained through

the judgment after all these new variables are placed into context with relative error 0 together with their substitutions with constants.

$$E, y_i^A : \mathcal{O}_A(0), y_i^A \triangleright c_i \vdash_P \pi; x_B, x_C, x_D \rightsquigarrow \varphi/V_B, V_C, V_D$$

Notice how the pieces of target code shown in Tables I and II, which represent various stages and phases of computation, are pasted together into the target program in Figure 8. For clarity, the empty suspensions and forcings have been deleted from this target program.

## 5. Performance

We have already mentioned that our automatically generated code for 2- and 3-dimensional `Orient`, `InCircle` and `InSphere` predicates to a large extent resembles that of Shewchuk [24]. Of course, this similarity is hard to quantify, if for no other reason than because our predicates are generated in our target language, while Shewchuk’s predicates are in C. Nevertheless, we wanted to measure the extent to which the logical and mathematical differences in the code influence the efficiency of our predicates. For that purpose we translated (automatically) the generated predicates from the target language into C and compared the translations against Shewchuk’s C implementations. All the results are obtained on a Pentium II on 266 MHz and 96 Mb of RAM.

The first test consisted of running the compared predicates on a common set of input entries. Each set had 1000 entries, and each entry was a list of point coordinates, in cardinality and dimension appropriate for the particular predicate. The coordinates of the points were drawn with a uniform random distribution from the set of floating-point numbers with exponents between  $-63$  and  $63$ . The summary of the results is represented in Table III. As can be seen, our C predicates are of comparable speed with Shewchuk’s, except in the case of `InSphere` where Shewchuk’s hand-tuned version is about 2.4 times faster. The `InSphere` predicate is the most complex of all and it is only natural that it can benefit the most from optimizations.

One of the most visible differences between our `InSphere` predicate and Shewchuk’s is the number of variables declared in the program. Our version of `InSphere` declares a new `double` array (which can be of considerable size) for every local variable in the target code intended to hold an exact value of an expansion type. However, a lot of this memory can actually be reused, because only a minor portion of the exact values needs to be accessible throughout the run of the program. The reuse will improve the cache management of the automatically

Table III. Performance comparison with Shewchuk’s predicates. The presented results are times for an average run of a predicate on random inputs.

	Shewchuk’s version	Automatically generated version	Ratio
Orient2D	0.208 ms	0.249 ms	1.197
Orient3D	0.707 ms	0.772 ms	1.092
InCircle	6.440 ms	5.600 ms	0.870
InSphere	16.430 ms	39.220 ms	2.387

Table IV. Performance comparison with Shewchuk’s predicates for 2d divide-and-conquer Delaunay triangulation.

	Shewchuk’s version	Automatically generated version	Ratio
uniform random	1187.1 ms	1410.3 ms	1.19
tilted grid	2060.4 ms	3677.5 ms	1.78
co-circular	1190.2 ms	1578.3 ms	1.33

generated programs and certainly increase their time efficiency. However, it is important to notice that this problem is not inherent to the automatically generated predicates, but is due to the naive translation from our target language into C. A better translator could probably decrease these differences considerably.

For the second test we modified *Triangle*, Shewchuk’s 2d Delaunay triangulator [22] to use the automatically generated predicates. The testing included triangulations of three different sets of 50,000 sample points: uniformly random *in* a unit square, tilted grid and uniformly random *on* a unit circle. The summary of the results is represented in Table IV. As can be seen, our predicates are a bit slower in the degenerate cases of tilted grid and co-circular points. Triangulation of such point-distributions often requires the higher phases of the filter, which can most benefit from a more economical memory management.

## 6. Future Work

The most immediate extensions of the compiler should focus on exploiting the paradigm of staging even better. Staging of expressions prevents recomputing already obtained intermediate results, and each stage in the source program translates into four phases of the target

program, with four approximations of different precision to the given intermediate result. However, currently there is no feedback from the exact phase to the inexact phases. If a computation ever carries out its exact phase, it will obtain the exact values that the intermediate results from the previous stages were approximating, and could potentially use this exact value to increase the accuracy of the intermediate approximations. It would be interesting and useful to devise a scheme that would exploit both the adaptive precision arithmetic and the staging in this broader manner.

A longer term goal could be to exploit the structure of the computation to obtain better error bounds. Priest has derived sufficient conditions which guarantee that the result from a certain floating-point operation will actually be computed exactly, i.e. will not incur any roundoff error [20]. While putting this idea in practice will likely require a non-trivial amount of theorem proving, it might still be feasible, since geometric predicates are typically short expressions, and the time for their compilation is not really crucial.

Finally, one may wonder how to extend the source language with higher-order functions and recursion. Adding higher-order functions for the sake of structuring the code will most likely require that every single intermediate variable in the program be replaced with a tuple containing that variable's phase A value and a suspension for the other three phases. This is required since now functions in the language can test signs of arbitrary values, even those produced by other functions, so the values have to be equipped with means to compute themselves exactly. But this is likely to be too slow, defeating the whole purpose of the expression compiler. On the other hand, adding recursive functions is even less realistic. Performing error analysis for recursive functions is hard – it is one of the main goals of the whole mathematical field of numerical analysis.

## 7. Conclusion

This paper presents an expression compiler for automatic generation of functions for testing the sign of a given arithmetic expression over floating-point constants and variables. In addition to the basic operations of addition, subtraction, multiplication, squaring and negation, our expressions can contain nested anonymous functions and thus exploit the optimization technique of staging, that is well-known in functional programming. The output of the compiler is a target program in a suitably designed intermediate language, which can be easily converted to SML or, in case of single-stage programs, to C.



Our compiler is an extension of the idea of Shewchuk [24]. Shewchuk used his method to develop quick robust implementations for the Orient and InSphere geometric predicates. We generalize this approach to arbitrary expressions. In particular, when applied to source expressions for these geometric predicates, our compiler generates code that, to a large extent, resembles that of Shewchuk. The idea behind the approach is to split the computation into several phases of increasing precision (but decreasing speed), each of which builds upon the result of the previous phase, while using forward error analysis to achieve reliable sign tests.

There remain, however, two caveats when generating predicates with this general approach – the produced code works correctly (1) only if no overflow or underflow happen, and (2) only in round-to-nearest, tie-to-even floating-point arithmetic complying with the IEEE standard. If overflow or underflow happens in the course of the run of some predicate, the expansions holding exact intermediate results may lose bits of information and distort the final outcome. Thus, we need to recognize such situations and, in those supposedly rare cases, rerun the computation in another form of exact arithmetic (say in infinite precision rational numbers). Unfortunately, even though the IEEE standard prescribes flags that can be read to check for overflow and underflow, the Standard Basis Library of ML does not provide any functions for their testing.

As concerning the second requirement, the IEEE standard is implemented on most modern processors. Unfortunately, on the Intel x86 family this is not a default setup. This family uses internal floating-point registers that are larger than 64-bits reserved for values of floating-point type. This property can occasionally make them round incorrectly in the to-nearest mode (for an example, see [20] page 103) and thus destroys the soundness of the language semantics. This default can be changed by setting a processor flag, but again, the Standard Basis Library does not provide any means for it. We believe that these two described insufficiencies can easily be remedied, and should be if SML is to become a language with serious applications in numerical analysis and scientific computing.

## References

1. Barber, C. B., D. P. Dobkin, and H. Huhdanpaa: 1996, ‘The Quickhull Algorithm for Convex Hulls’. *ACM Transactions on Mathematical Software* **22**(4), 469–483.
2. Benouamer, M. O., P. Jaillon, D. Michelucci, and J. M. Moreau: 1993, ‘A lazy exact arithmetic’. In: E. E. Swartzlander, M. J. Irwin, and J. Jullien (eds.):

- Proceedings of the 11th IEEE Symposium on Computer Arithmetic*. Windsor, Canada, pp. 242–249.
3. Blleloch, G., H. Burch, K. Crary, R. Harper, G. Miller, and N. Walkington: 2001, ‘Persistent Triangulations’. *Journal of Functional Programming*. To appear.
  4. Brönnimann, H., C. Burnikel, and S. Pion: 2001, ‘Interval Arithmetic Yields Efficient Dynamic Filters for Computational Geometry’. *Discrete Applied Mathematics* **109**(1–2), 25–47.
  5. Burnikel, C., S. Funke, and M. Seel: 2001, ‘Exact Geometric Computation using Cascading’. *International Journal of Computational Geometry and Applications* **11**(3), 245–266.
  6. Daumas, M.: 1999, ‘Multiplications of floating point expansions’. In: I. Koren and P. Kornerup (eds.): *Proceedings of the 14th Symposium on Computer Arithmetic*. Adelaide, Australia, pp. 250–257.
  7. Daumas, M. and C. Finot: 1998, ‘Division of floating point expansions’. In: *IMACS-GAMM International Symposium on Scientific Computing, Computer Arithmetic and Validated Numerics*. Budapest, Hungaria, pp. 45–46.
  8. Daumas, M. and C. Finot: 1999a, ‘Algorithm, Proof and Performances of a new Division of Floating Point Expansions’. Technical Report 3771, INRIA, Le Chesnay, France.
  9. Daumas, M. and C. Finot: 1999b, ‘Division of floating point expansions with an application to the computation of a determinant’. *Journal of Universal Computer Science* **5**(6), 323–338.
  10. Dekker, T. J.: 1971, ‘A Floating Point Technique for Extending the Available Precision’. *Numerische Mathematik* **18**(3), 224–242.
  11. Devillers, O. and S. Pion: 2002, ‘Efficient Exact Geometric Predicates for Delaunay Triangulations’. Technical Report 4351, INRIA.
  12. Fabri, A., B. Gärtner, S. Hert, S. Hirsch, M. Hoffmann, L. Kettner, S. Pion, M. Teillaud, R. Veltkamp, and M. Yvinec: 2002, ‘The CGAL Manual’. Release 2.4.
  13. Fortune, S. and C. J. V. Wyk: 1993, ‘Efficient Exact Arithmetic for Computational Geometry’. In: *Ninth Annual Symposium on Computational Geometry*. pp. 163–172.
  14. Funke, S. and K. Mehlhorn: 2000, ‘LOOK – A Lazy Object-Oriented Kernel for Geometric Computation’. In: *Proceedings of the 16th Symposium on Computational Geometry*. pp. 156–165.
  15. IEEE: 1985, ‘IEEE Standard for Binary Floating-Point Arithmetic’. *ACM SIGPLAN Notices* **22**(2), 9–25.
  16. Knuth, D. E.: 1981, *The Art of Computer Programming: Seminumerical Algorithms*. Addison-Wesley.
  17. Michelucci, D. and J. M. Moreau: 1997, ‘Lazy arithmetic’. *IEEE Transactions on Computers* **46**(9).
  18. Nanevski, A., G. Blleloch, and R. Harper: 2001, ‘Automatic Generation of Staged Geometric Predicates’. Technical Report CMU-CS-01-141, Carnegie Mellon University.
  19. Priest, D. M.: 1991, ‘Algorithms for arbitrary precision floating point arithmetic’. In: P. Kornerup and D. W. Matula (eds.): *Proceedings of the 10th IEEE Symposium on Computer Arithmetic (Arith-10)*. Grenoble, France, pp. 132–144.

20. Priest, D. M.: 1992, 'On Properties of Floating Point Arithmetics: Numerical Stability and the Cost of Accurate Computations'. Ph.D. thesis, University of California, Berkeley.
21. Seshia, S. A., G. E. Blelloch, and R. W. Harper: 2000, 'A Performance Comparison of Interval Arithmetic and Error Analysis in Geometric Predicates'. Technical Report CMU-CS-00-172, School of Computer Science, Carnegie Mellon University.
22. Shewchuk, J. R.: 1996, 'Triangle: Engineering a 2D Quality Mesh Generator and Delaunay Triangulator'. In: *First Workshop on Applied Computational Geometry*. Philadelphia, Pennsylvania, pp. 124–133. <http://www.cs.cmu.edu/~quake/triangle.html>.
23. Shewchuk, J. R.: 1997a, 'Adaptive Precision Floating-Point Arithmetic and Fast Robust Geometric Predicates'. *Discrete and Computational Geometry* **18**, 305–363.
24. Shewchuk, J. R.: 1997b, 'Delaunay Refinement Mesh Generation'. Ph.D. thesis, Carnegie Mellon University.

