

Strongly Typed Numerical Computations[†]

Matthieu Martel

Laboratoire de Mathématiques et Physique (LAMPS)
Université de Perpignan Via Domitia, France
`matthieu.martel@univ-perp.fr`

Abstract. It is well-known that numerical computations may sometimes lead to wrong results because of roundoff errors. We propose an ML-like type system (strong, implicit, polymorphic) for numerical computations in finite precision, in which the type of an expression carries information on its accuracy. We use dependent types and a type inference which, from the user point of view, acts like ML type inference. Basically, our type system accepts expressions for which it may ensure a certain accuracy on the result of the evaluation and it rejects expressions for which a minimal accuracy on the result of the evaluation cannot be inferred. The soundness of the type system is ensured by a subject reduction theorem and we show that our type system is able to type implementations of usual simple numerical algorithms.

1 Introduction

It is well-known that numerical computations may sometimes lead to wrong results because of the accumulation of roundoff errors [8]. Recently, much work has been done to detect these accuracy errors in finite precision computations [1], by static [6, 9, 18] or dynamic [7] analysis, to find the least data formats needed to ensure a certain accuracy (precision tuning) [11, 12, 17] and to optimize the accuracy by program transformation [5, 14]. All these techniques are used late in the software development cycle, once the programs are entirely written.

In this article, we aim at exploring a different direction. We aim at detecting and correcting numerical accuracy errors at software development time, i.e. during the programming phase. From a software engineering point of view, the advantages of our approach are many since it is well-known that late bug detection is time and money consuming. We also aim at using intensively used techniques recognized for their ability to discard run-time errors. This choice is motivated by efficiency reasons as well as for end-user adoption reasons.

We propose an ML-like type system (strong, implicit, polymorphic [15]) for numerical computations in which the type of an arithmetic expression carries

[†]This work is supported by the Office for Naval Research Global under Grant NICOP N62909-18-1-2068 (Tycoon project).
<https://www.onr.navy.mil/en/Science-Technology/ONR-Global>

information on its accuracy. We use dependent types [16] and a type inference which, from the user point of view, acts like ML [13] type inference [15] even if it slightly differs in its implementation. While type systems have been widely used to prevent a large variety of software bugs, to our knowledge, no type system has been targeted to address numerical accuracy issues in finite precision computations. Basically, our type system accepts expressions for which it may ensure a certain accuracy on the result of the evaluation and it rejects expressions for which a minimal accuracy on the result of the evaluation cannot be inferred.

Let us insist on the fact that we use a dependent type system. Consequently, the type corresponding to a function of some argument x depends on the type of x itself. The soundness of our type system relies on a subject reduction theorem introduced in Section 4. Based on an instrumented operational semantics computing both the finite precision and exact results of a numerical computation, this theorem shows that the error on the result of the evaluation of some expression e is less than the error predicted by the type of e . Obviously, as any non-trivial type system, our type system is not complete and rejects certain programs that would not produce unbounded numerical errors. Our type system has been implemented in a prototype language `Num1` and we show that, in practice, our type system is expressive enough to type implementations of usual simple numerical algorithms [2] such as the ones of Section 5. Let us also mention that our type system represents a new application of dependent type theory motivated by applicative needs. Indeed, dependent types arise naturally in our context since accuracy depends on values.

This article is organized as follows. Section 2 introduces informally our type system and shows how it is used in our implementation of a ML-like programming language, `Num1`. The formal definition of the types and of the inference rules are given in Section 3. A soundness theorem is given in Section 4. Section 5 presents experimental results and Section 6 concludes.

2 Programming with Types for Numerical Accuracy

In this section, we present informally how our type system works throughout a programming sequence in our language, `Num1`. First of all, we use real numbers $r\{s, u, p\}$ where r is the value itself, and $\{s, u, p\}$ the format of r . The format of a real number is made of a sign $s \in \text{Sign}$ and integers $u, p \in \text{Int}$ such that u is the unit in the first place of r , written $\text{ufp}(r)$ and p the precision (i.e. the number of digits of the number). We have $\text{Sign} = \{0, \oplus, \ominus, \top\}$ and $\text{sign}(r) = 0$ if $r = 0$, $\text{sign}(r) = \oplus$ if $r > 0$ and $\text{sign}(r) = \ominus$ if $r < 0$. The set Sign is equipped with the partial order relation $\prec \subseteq \text{Sign} \times \text{Sign}$ defined by $0 \prec \oplus$, $0 \prec \ominus$, $\oplus \prec \top$ and $\ominus \prec \top$. The ufp of a number x is

$$\text{ufp}(x) = \min \{i \in \mathbb{N} : 2^{i+1} > x\} = \lfloor \log_2(x) \rfloor . \quad (1)$$

The term p defines the precision of r . Let $\varepsilon(r)$ be the absolute error on r , we assume that $\varepsilon(r) < 2^{u-p+1}$. The errors on the numerical constants arising in programs are specified by the user or determined by default by the system.

Format	Name	p	e bits	e_{min}	e_{max}
Binary16	Half precision	11	5	-14	+15
Binary32	Single precision	24	8	-126	+127
Binary64	Double precision	53	11	-1122	+1223
Binary128	Quadruple precision	113	15	-16382	+16383

Fig. 1. Basic binary IEEE754 formats.

The errors on the computed values can be inferred by propagation of the initial errors. Similarly to Equation (1), we also define the *unit* in the *last place* (ulp) used later in this article. The ulp of a number of precision p is defined by

$$\text{ulp}(x) = \text{ufp}(x) - p + 1 . \quad (2)$$

For example, the type of 1.234 is `real{+,0,53}` since `ufp(1.234) = 0` and since we assume that, by default, the real numbers have the same precision as in the IEEE754 double precision floating-point format [1] (see Figure 1). Other formats may be specified by the programmer, as in the example below. Let us also mention that our type system is independent of a given computer arithmetic. The interpreter only needs to implement the formats given by the type system, using floating-point numbers, fixed-point numbers [10], multiple precision numbers¹, etc in order to ensure that the finite precision operations are computed exactly. The special case of IEEE754 floating-point arithmetic, which introduces additional errors due to the roundoff on results of operations can also be treated by modifying slightly the equations of Section 3.

```
> 1.234 ;; (* precision of 53 bits by default *)
- : real{+,0,53} = 1.2340000000000000

> 1.234{4};; (* precision of 4 bits specified by the user *)
- : real{+,0,4} = 1.2
```

Notice that, in `Num1`, the type information is used by the pretty printer to display only the correct digits of a number and a bound on the roundoff error.

Note that accuracy is not a property of a number but a number that states how closely a particular floating-point number matches some ideal true value. For example, using the basis $\beta = 10$ for the sake of simplicity, the floating-point value 3.149 represents π with an accuracy of 3. It itself has a precision of 4. It represents the real number 3.14903 with an accuracy of 4. As in ML, our type system admits parameterized types [15].

```
> let f = fun x -> x + 1.0 ;;
val f : real{'a','b','c'} -> real{<expr>,<expr>,<expr>} = <fun>

> verbose true ;;
- : unit = ()

> f ;;
- : real{'a','b','c'} -> real{(SignPlus 'a 'b 1 0),((max 'b 0) +_ (sigma+ 'a 1)),
(((max 'b 0) +_ (sigma+ 'a 1)) -_ (max ('b -_ 'c) -53))-_ (iota ('b -_ 'c) -53))} = <fun>
```

¹<https://gmplib.org/>

In the example above, the type of `f` is a function of an argument whose parameterized type is `real{'a','b','c'}`, where `'a'`, `'b'` and `'c'` are three type variables. The return type of the function `f` is `Real{e0,e1,e2}` where `e0`, `e1` and `e2` are arithmetic expressions containing the variables `'a'`, `'b'` and `'c'`. By default these expressions are not displayed by the system (just like higher order values are not explicitly displayed in ML implementations) but we may enforce the system to print them. In `Num1`, we write `+`, `-`, `*` and `/` the floating-point operators. The integer operators are written `+`, `-`, `*` and `/`. The expressions arising in the type of `f` are explained in Section 3. As shown below, various applications of `f` yield results of various types, depending on the type of the argument.

```
> f 1.234 ;;
- : real{+,1,53} = 2.2340000000000000
```

```
> f 1.234{4} ;;
- : real{+,1,5} = 2.2
```

If the interpreter detects that the result of some computation has no significant digit, then an error is raised. For example, it is well-known that in IEEE754 double precision $(10^{16} + 1) - 10^{16} = 0$. Our type system rejects this computation.

```
> (1.0e15 + 1.0) - 1.0e15 ;;
- : real{+,50,54} = 1.0
```

```
> (1.0e16 + 1.0) - 1.0e16 ;;
Error: The computed value has no significant digit. Its ufp is 0 but the ulp of the certified value is 1
```

Last but not least, our type system accepts recursive functions. For example, we have:

```
> let rec g x = if x < 1.0 then x else g (x * 0.07) ;;
val g : real{+,0,53} -> real{+,0,53} = <fun>
```

```
> g 1.0 ;;
- : real{+,0,53} = 0.0700000000000000
```

```
> g 2.0 ;;
Error: This expression has type real{+,1,53} but an expression was expected of type real{+,0,53}
```

In the above session, the type system unifies the return type of the function with the type of the conditional. The types of the `then` and `else` branches also need to be unified. Then the return type is `real{+,0,53}` which corresponds to the type of the value `1.0` used in the `then` branch. The type system also unifies the return type with the type of the argument since the function is recursive. Finally, we obtain that the type of `g` is `real{+,0,53} -> real{+,0,53}`. As a consequence, we cannot call `g` with an argument whose `ufp` is greater than `ufp(1.0) = 0`. To overcome this limitation, we introduce new comparison operations for real numbers. While the standard comparison operator `<` has type `'a -> 'a -> bool`, the operator `<{s,u,p}` has type `real{s,u,p} -> real{s,u,p} -> bool`. In other words, the compared value are cast in the format `{s,u,p}` before performing the comparison. Now we can write the code:

```
> let rec g x = if x <{* ,10,15} 1.0 then x else g (x * 0.07) ;;
val g : real{* ,10,15} -> real{* ,10,15} = <fun>
```

```
> g 2.0 ;;
- : real{* ,10,15} = 0.1
```

```

> g 456.7 ;;
- : real{*,10,15} = 0.1

> g 4567.8 ;;
Error: This expression has type real{+,12,53} but an expression was expected of
type real{*,10,15}

```

Interestingly, unstable functions (for which the initial errors grow with the number of iterations) are not typable. This is a desirable property of our system.

```

> let rec h n = if (n=0) then 1.0 else 3.33 * (h (n -_ 1)) ;;
Error: This expression has type real{+,-1,-1} but an expression was expected of
type real{+,-3,-1}

```

Stable computations should be always accepted by our type system. Obviously, this is not the case and, as any non-trivial type system, our type system rejects some correct programs. The challenge is then to accept enough programs to be useful from an end-user point of view. We end this section by showing another example representative of what our type system accepts. More examples are given later in this article, in Section 5. The example below deals with the implementation of the Taylor series $\frac{1}{1-x} = \sum_{n \geq 0}^{+\infty} x^n$. The implementation gives rise to a simple recursion, as shown in the programming session below.

```

> let rec taylor x{*,-1,25} xn i n = if (i > n) then 0.0{*,10,20}
                                     else xn + (taylor x (x * xn) (i +_ 1) n) ;;
val taylor : real{*,-1,25} -> real{*,10,20} -> int -> int -> real{*,10,20} = <fun>

> taylor 0.2 1.0 0 5;;
- : real{*,10,20} = 1.2499 +/- 0.0009765625

```

Obviously, our type system computes the propagation of the errors due to finite precision but does not take care of the method error intrinsic to the implemented algorithm (the Taylor series instead of the exact formula $\frac{1}{1-x}$ in our case.) All the programming sessions introduced above as well as the additional examples of Section 5 are fully interactive in our system, Num1, i.e. the type judgments are obtained instantaneously (about 0.01 second in average following our measurements) including the most complicated ones.

3 The Type System

In this section, we introduce the formal definition of our type system for numerical accuracy. First, in Section 3.1, we define the syntax of expressions and types and we introduce a set of inference rules. Then we define in Section 3.2 the types of the primitives for the operators among real numbers (addition, product, etc.) These types are crucial in our system since they encode the propagation of the numerical accuracy information.

3.1 Expressions, Types and Inference Rules

In this section, we introduce the expressions, types and typing rules for our language. For the sake of simplicity, the syntax introduced hereafter uses notations

$$\begin{array}{c}
\frac{}{\Gamma \vdash i : \text{int}} \text{(INT)} \qquad \frac{}{\Gamma \vdash \mathbf{b} : \text{bool}} \text{(BOOL)} \\
\frac{\text{sign}(x) \prec s \quad \text{ufp}(x) \leq u}{\Gamma \vdash \mathbf{r}\{s, u, p\} : \text{real}\{s, u, p\}} \text{(REAL)} \qquad \frac{\Gamma(\text{id}) = t}{\Gamma \vdash \text{id} : t} \text{(ID)} \\
\frac{\Gamma \vdash e_0 : \text{bool} \quad \Gamma \vdash e_1 : t_1 \quad \Gamma \vdash e_2 : t_2 \quad t = t_1 \sqcup t_2}{\Gamma \vdash \text{if } e_0 \text{ then } e_1 \text{ else } e_2 : t} \text{(COND)} \\
\frac{\Gamma, x : t_1 \vdash e : t_2}{\Gamma \vdash \lambda x. e : \Pi x : t_1. t_2} \text{(ABS)} \qquad \frac{\Gamma, x : t_1, f : \Pi y : t_1. t_2 \vdash e : t_2 \quad y \text{ not free in } t_2}{\Gamma \vdash \text{rec } f x. e : \Pi x : t_1. t_2} \text{(REC)} \\
\frac{\Gamma \vdash e_1 : \Pi x : t_0. t_1 \quad \Gamma \vdash e_2 : t_2 \quad t_2 \sqsubseteq t_0}{\Gamma \vdash e_1 e_2 : t_2[x \mapsto e_1]} \text{(APP)}
\end{array}$$

Fig. 2. Typing rules for our language.

à la lambda calculus instead of the ML-like syntax employed in Section 2. In our system, expressions and types are mutually dependent. They are defined inductively using the grammar of Equation (3).

$$\begin{aligned}
\text{Expr } \ni e ::= & \mathbf{r}\{s, u, p\} \in \text{Real}_{u,p} \mid i \in \text{Int} \mid \mathbf{b} \in \text{Bool} \mid \text{id} \in \text{Id} \\
& \mid \text{if } e_0 \text{ then } e_1 \text{ else } e_2 \mid \lambda x. e \mid e_0 e_1 \mid \text{rec } f x. e \mid t \\
\text{Typ } \ni t ::= & \mid \text{int} \mid \text{bool} \mid \text{real}\{i_0, i_1, i_2\} \mid \alpha \mid \Pi x : e_0. e_1 \\
\text{LExp } \ni i ::= & \mid \text{int} \mid \text{op} \in \text{Id}_i \mid \alpha \mid i_0 i_1
\end{aligned} \tag{3}$$

In Equation (3), the e terms correspond to expressions. Constants are integers $i \in \text{Int}$, booleans $\mathbf{b} \in \text{Bool}$ and real numbers $\mathbf{r}\{s, u, p\}$ where \mathbf{r} is the value itself, $s \in \text{Sign}$ is the sign and $u, p \in \text{Int}$ the ufp (see Equation (1)) and precision of \mathbf{r} . We have $\text{Sign} = \{0, \oplus, \ominus, \top\}$ and $\text{sign}(r) = 0$ if $r = 0$, $\text{sign}(r) = \oplus$ if $r > 0$ and $\text{sign}(r) = \ominus$ if $r < 0$. The set Sign is equipped with the partial order relation $\prec \subseteq \text{Sign} \times \text{Sign}$ defined by $0 \prec \oplus$, $0 \prec \ominus$, $\oplus \prec \top$ and $\ominus \prec \top$. The term p defines the precision of r . Let $\varepsilon(r)$ be the absolute error on r , we assume that

$$\varepsilon(r) < 2^{u-p+1} . \tag{4}$$

The errors on the numerical constants arising in programs are specified by the user or determined by default by the system. The errors on the computed values can be inferred by propagation of the initial errors.

In Equation (3), identifiers belong to the set Id and we assume a set of pre-defined identifiers $+$, $-$, \times , \leq , $=$, \dots related to primitives for the logical and arithmetic operations. We write $+$, $-$, \times and \div the operations on real numbers and $+$, $-$, \times and \div the operations among integers. The language also admits conditionals, functions $\lambda x. e$, applications $e_0 e_1$ and recursive functions $\text{rec } f x. e$ where f is the name of the function, x the parameter and e the body. The language of expressions also includes type expressions t defined by the second production of the grammar of Equation (3).

The definition of expressions and type is mutually recursive. Type variables are denoted α, β, \dots and $\Pi x : e_0. e_1$ is used to introduce dependent types [16].

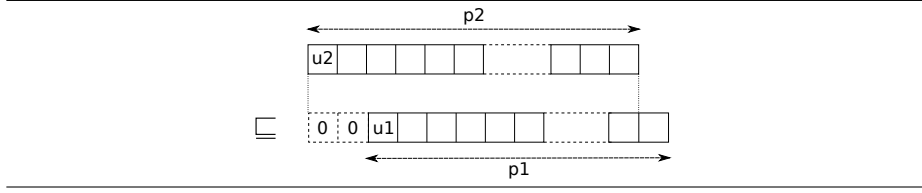


Fig. 3. The sub-typing relation \sqsubseteq of Equation (6).

Let us notice that our language does not explicitly contain function types $t_0 \rightarrow t_1$ since they are encoded by means of dependent types. Let \equiv denote the syntactic equivalence, we have

$$t_0 \rightarrow t_1 \equiv \Pi x : t_0. t_1 \quad \text{with } x \text{ not free in } t_1 . \quad (5)$$

For convenience, we also write $\lambda x_0. x_1 \dots x_n. e$ instead of $\lambda x_0. \lambda x_1 \dots \lambda x_n. e$ and $\Pi x_0 : t_0. x_1 : t_1 \dots x_n : t_n. e$ instead of $\Pi x_0 : t_0. \Pi x_1 : t_1 \dots \Pi x_n : t_n. e$.

The types of constants are `int`, `bool` and `real` $\{i_0, i_1, i_2\}$ where i_0 , i_1 and i_2 are integer expressions denoting the format of the real number. Integer expressions of `lExpr` \subseteq `Expr` are a subset of expressions made of integer numbers, integer primitives of `ldl` \subseteq `ld` (such as $+$, $-$, \times , etc.), type variables and applications. Note that this definition restricts significantly the set of expressions which may be written inside `real` types.

The typing rules for our system are given in Figure 2. These rules are mostly classical. The type judgment $\Gamma \vdash e : t$ means that in the type environment Γ , the expression e has type t . A type environment $\Gamma : \text{ld} \rightarrow \text{Typ}$ map identifiers to types. We write $\Gamma x : t$ the environment Γ in which the variable x has type t . The typing rules (INT) and (BOOL) are trivial. Rule (REAL) states that the type of a real number `r` $\{s, u, p\}$ is `real` $\{s, u, p\}$ assuming that the actual sign of `r` is less than s and that the `ulp` of `r` is less than u . Following Rule (ID), an identifier `id` has type t if $\Gamma(\text{id}) = t$. Rules (COND), (ABS) and (REC) are standard rules for conditionals and abstractions respectively. The rule for application, (APP), requires that the first expression e_1 has type $\Pi x : t_0. t_1$ (which is equivalent to $t_0 \rightarrow t_1$ if x is not free in t_1) and that the argument e_2 has some type $t_2 \sqsubseteq t_0$. The sub-typing relation \sqsubseteq is introduced for real numbers. Intuitively, we want to allow the argument of some function to have a smaller `ulp` than what we would require if we used $t_0 = t_2$ in Rule (APP), provided that the precision p remains as good with t_2 as with t_0 . This relaxation allows to type more terms without invalidating the type judgments. Formally, the relation \sqsubseteq is defined by

$$\text{real}\{s_1, u_1, p_1\} \sqsubseteq \text{real}\{s_2, u_2, p_2\} \iff s_1 \sqsubseteq s_2 \wedge u_2 \geq u_1 \wedge p_2 \leq u_2 - u_1 + p_1 . \quad (6)$$

In other words, the sub-typing relation of Equation (6) states that it is always correct to add zeros before the first significant digit of a number, as illustrated in Figure 3.

3.2 Types of Primitives

In this section, we introduce the types of the primitives of our language. As mentioned earlier, the arithmetic and logic operators are viewed as functional constants of the language. The type of a primitive for an arithmetic operation among integers $*_- \in \{+_, -_, \times_, \div_-\}$ is

$$t_{*_} = \Pi x : \text{int}. y : \text{int}. \text{int} . \quad (7)$$

The type of comparison operators $\bowtie \in \{=, \neq, <, >, \leq, \geq\}$ are polymorphic with the restriction that they reject the type $\text{real}\{\mathbf{s}, \mathbf{u}, \mathbf{p}\}$ which necessitates special comparison operators:

$$t_{\bowtie} = \Pi x : \alpha. y : \alpha. \text{bool} \quad \alpha \neq \text{real}\{\mathbf{s}, \mathbf{u}, \mathbf{p}\} . \quad (8)$$

For real numbers, we use comparisons at a given accuracy defined by the operators $\bowtie_{\{u,p\}} \in \{\<_{\{u,p\}}, >_{\{u,p\}}\}$. We have

$$t_{\bowtie_{\{u,p\}}} = \Pi \mathbf{s} : \text{int}, \mathbf{u} : \text{int}, \mathbf{p} : \text{int}. \text{real}\{\mathbf{s}, \mathbf{u}, \mathbf{p} + 1\} \rightarrow \text{real}\{\mathbf{s}, \mathbf{u}, \mathbf{p} + 1\} \rightarrow \text{bool} .$$

Notice that the operands of a comparison $\bowtie_{\{u,p\}}$ must have $p+1$ bits of accuracy. This is to avoid unstable tests, as detailed in the proof of Lemma 3 in Section 4. An unstable test is a comparison between two approximate values such that the result of the comparison is altered by the approximation error. For instance, if we reuse an example of Section 2, in IEEE754 double precision, the condition $10^{16} + 1 = 10^{16}$ evaluates to **true**. We need to avoid such situations in our language in order to preserve our subject reduction theorem (we need the control-flow be the same in the finite precision and exact semantics). Let us also note that our language does not provide an equality relation $=_{\{u,p\}}$ for **real** values. Again this is to avoid unstable tests. Given values x and y of type $\text{real}\{\mathbf{s}, \mathbf{u}, \mathbf{p}\}$, the programmer is invited to use $|x - y| < 2^{u-p+1}$ instead of $x = y$ in order to get rid of the perturbations of the finite precision arithmetic.

The types of primitives for real arithmetic operators are fundamental in our system since they encode the propagation of the numerical accuracy information. They are defined in figures 4 and 5. The type t_* of some operation $* \in \{+, -, \times, \div\}$ is a pi-type with takes six arguments $\mathbf{s}_1, \mathbf{u}_1, \mathbf{p}_1, \mathbf{s}_2, \mathbf{u}_2$ and \mathbf{p}_2 of type **int** corresponding to the sign, **ufp** and precision of the two operands of $*$ and which produces a type $\text{real}\{\mathbf{s}_1, \mathbf{u}_1, \mathbf{p}_1\} \rightarrow \text{real}\{\mathbf{s}_2, \mathbf{u}_2, \mathbf{p}_2\} \rightarrow \text{real}\{\mathcal{S}_*(\mathbf{s}_1, \mathbf{s}_2), \mathcal{U}_*(\mathbf{s}_1, \mathbf{u}_1, \mathbf{s}_2, \mathbf{u}_2), \mathcal{P}_*(\mathbf{u}_1, \mathbf{p}_1, \mathbf{u}_2, \mathbf{p}_2)\}$ where \mathcal{S}_* , \mathcal{U}_* and \mathcal{P}_* are functions which compute the sign, **ufp** and precision of the result of the operation $*$ in function of $\mathbf{s}_1, \mathbf{u}_1, \mathbf{p}_1, \mathbf{s}_2, \mathbf{u}_2$ and \mathbf{p}_2 . These functions extend the functions used in [12].

The functions \mathcal{S}_* determine the sign of the result of an operation in function of the signs of the operands and, for additions and subtractions, in function of the **ufp** of the operands. The functions \mathcal{U}_* compute the **ufp** of the result. Notice that \mathcal{U}_+ and \mathcal{U}_- use the functions σ_+ and σ_- , respectively. These functions are defined in the bottom right corner of Figure 5 to increment the **ufp** of the result of some addition or subtraction in the relevant cases only. For example if a and b are two positive real numbers then $\text{ufp}(a + b)$ is possibly $\max(\text{ufp}(a), \text{ufp}(b)) + 1$

$$\begin{aligned}
\mathbf{t}_* &= \mathbb{I} \mathbf{s}_1 : \mathbf{int}, \mathbf{u}_1 : \mathbf{int}, \mathbf{p}_1 : \mathbf{int}, \mathbf{s}_2 : \mathbf{int}, \mathbf{u}_2 : \mathbf{int}, \mathbf{p}_2 : \mathbf{int}. \\
&\quad \mathbf{real}\{\mathbf{s}_1, \mathbf{u}_1, \mathbf{p}_1\} \rightarrow \mathbf{real}\{\mathbf{s}_2, \mathbf{u}_2, \mathbf{p}_2\} \\
&\quad \rightarrow \mathbf{real}\{\mathcal{S}_*(\mathbf{s}_1, \mathbf{u}_1, \mathbf{s}_2, \mathbf{u}_2), \mathcal{U}_*(\mathbf{s}_1, \mathbf{u}_1, \mathbf{s}_2, \mathbf{u}_2), \mathcal{P}_*(\mathbf{s}_1, \mathbf{u}_1, \mathbf{p}_1, \mathbf{s}_2, \mathbf{u}_2, \mathbf{p}_2)\} \\
\mathcal{U}_+(\mathbf{s}_1, \mathbf{u}_1, \mathbf{s}_2, \mathbf{u}_2) &= \max(\mathbf{u}_1, \mathbf{u}_2) + \sigma_+(\mathbf{s}_1, \mathbf{s}_2) \\
\mathcal{P}_+(\mathbf{s}_1, \mathbf{u}_1, \mathbf{p}_1, \mathbf{s}_2, \mathbf{u}_2, \mathbf{p}_2) &= \max(\mathbf{u}_1, \mathbf{u}_2) + \sigma_+(\mathbf{s}_1, \mathbf{s}_2) - \\
&\quad \max(\mathbf{u}_1 - \mathbf{p}_1, \mathbf{u}_2 - \mathbf{p}_2) - \iota(\mathbf{u}_1 - \mathbf{p}_1, \mathbf{u}_2 - \mathbf{p}_2) \\
\mathcal{U}_-(\mathbf{s}_1, \mathbf{u}_1, \mathbf{s}_2, \mathbf{u}_2) &= \max(\mathbf{u}_1, \mathbf{u}_2) + \sigma_-(\mathbf{s}_1, \mathbf{s}_2) \\
\mathcal{P}_-(\mathbf{s}_1, \mathbf{u}_1, \mathbf{p}_1, \mathbf{s}_2, \mathbf{u}_2, \mathbf{p}_2) &= \max(\mathbf{u}_1, \mathbf{u}_2) + \sigma_-(\mathbf{s}_1, \mathbf{s}_2) - \\
&\quad \max(\mathbf{u}_1 - \mathbf{p}_1, \mathbf{u}_2 - \mathbf{p}_2) - \iota(\mathbf{u}_1 - \mathbf{p}_1, \mathbf{u}_2 - \mathbf{p}_2) \\
\mathcal{U}_\times(\mathbf{s}_1, \mathbf{u}_1, \mathbf{s}_2, \mathbf{u}_2) &= \mathbf{u}_1 + \mathbf{u}_2 + 1 \\
\mathcal{P}_\times(\mathbf{s}_1, \mathbf{u}_1, \mathbf{p}_1, \mathbf{s}_2, \mathbf{u}_2, \mathbf{p}_2) &= \mathbf{u}_1 + \mathbf{u}_2 + 1 - \\
&\quad \max(\mathbf{u}_1 + \mathbf{u}_2 + 1 - \mathbf{p}_1, \mathbf{u}_1 + \mathbf{u}_2 + 1 - \mathbf{p}_2) - \iota(\mathbf{p}_1, \mathbf{p}_2) \\
\mathcal{U}_\div(\mathbf{s}_1, \mathbf{u}_1, \mathbf{s}_2, \mathbf{u}_2) &= \mathbf{u}_1 - \mathbf{u}_2 + 1 \\
\mathcal{P}_\div(\mathbf{s}_1, \mathbf{u}_1, \mathbf{p}_1, \mathbf{s}_2, \mathbf{u}_2, \mathbf{p}_2) &= \mathcal{P}_\times(\mathbf{u}_1, \mathbf{p}_1, \mathbf{u}_2, \mathbf{p}_2) - 1 \quad \iota(x, y) = \begin{cases} 1 & \text{if } x = y, \\ 0 & \text{otherwise.} \end{cases}
\end{aligned}$$

Fig. 4. Types of the primitives corresponding to the elementary arithmetic operations $*$ $\in \{+, -, \times, \div\}$. The functions \mathcal{S}_* and σ_* are defined in Figure 5.

but if $a > 0$ and $b < 0$ then $\text{ufp}(a + b)$ is not greater than $\max(\text{ufp}(a), \text{ufp}(b))$. The functions \mathcal{P}_* compute the precision of the result. Basically, they compute the number of bits between the ufp and the ulp of the result.

We end this section by exhibiting some properties of the functions \mathcal{P}_* . Let $\varepsilon(x)$ denote the error on $x \in \text{Real}_{u,p}$. We have $\varepsilon(x) < 2^{u-p+1} = \text{ulp}(x)$. Let us start with addition. Lemma 1 relates the accuracy of the operands to the accuracy of the result of an addition between two values x and y . Lemma 2 is similar to Lemma 1 for product.

Lemma 1. *Let x and y be two values such that $\varepsilon(x) < 2^{u_1-p_1+1}$ and $\varepsilon(y) < 2^{u_2-p_2+1}$. Let $z = x + y$, $u = \mathcal{U}_+(s_1, u_1, s_2, u_2)$ and $p = \mathcal{P}_+(s_1, u_1, p_1, s_2, u_2, p_2)$. Then $\varepsilon(z) < 2^{u-p+1}$.*

Proof. The errors on addition may be bounded by $e_+ = \varepsilon(x) + \varepsilon(y)$. Then the most significant bit of the error has weight $\text{ufp}(e_+)$ and the accuracy of the result is $p = \text{ufp}(x + y) - \text{ufp}(e_+)$. Let $u = \text{ufp}(x + y) = \max(u_1, u_2) + \sigma_+(s_1, s_2) = \mathcal{U}_+(s_1, u_1, s_2, u_2)$. We need to over-approximate e_+ in order to ensure p . We have $\varepsilon(x) < 2^{u_1-p_1+1}$ and $\varepsilon(y) < 2^{u_2-p_2+1}$ and, consequently, $e_+ < 2^{u_1-p_1+1} + 2^{u_2-p_2+1}$. We introduce the function $\iota(x, y)$ also defined in Figure 4 and which is equal to 1 if $x = y$ and 0 otherwise. We have

$$\begin{aligned}
\text{ufp}(e_+) &< \max(u_1 - p_1 + 1, u_2 - p_2 + 1) + \iota(u_1 - p_1, u_2 - p_2) \\
&\leq \max(u_1 - p_1, u_2 - p_2) + \iota(u_1 - p_1, u_2 - p_2)
\end{aligned}$$

Let us write $p = \max(u_1 - p_1, u_2 - p_2) - \iota(u_1 - p_1, u_2 - p_2) = \mathcal{P}_+(s_1, u_1, p_1, s_2, u_2, p_2)$. We conclude that $u = \mathcal{U}_+(s_1, u_1, s_2, u_2)$, $p = \mathcal{P}_+(s_1, u_1, p_1, s_2, u_2, p_2)$ and $\varepsilon(z) < 2^{u-p+1}$. \square

Lemma 2. *Let x and y be two values such that $\varepsilon(x) < 2^{u_1-p_1+1}$ and $\varepsilon(y) < 2^{u_2-p_2+1}$. Let $z = x \times y$, $u = \mathcal{U}_\times(s_1, u_1, s_2, u_2)$ and $p = \mathcal{P}_\times(s_1, u_1, p_1, s_2, u_2, p_2)$. Then $\varepsilon(z) < 2^{u-p+1}$.*

\mathcal{S}_+					\mathcal{S}_\times and \mathcal{S}_\div				
$s_1 \setminus s_2$	0	+	-	\top	$s_1 \setminus s_2$	0	+	-	\top
0	0	+	-	\top	0	0	0	0	0
+	+	+	+ if $u_1 < u_2$ - if $u_2 < u_1$ \top otherwise	\top	+	0	+	-	\top
-	-	+ if $u_2 < u_1$ - if $u_1 < u_2$ \top otherwise	-	\top	-	0	-	+	\top
\top	\top	\top	\top	\top	\top	0	\top	\top	\top

\mathcal{S}_-					σ_+				
$s_1 \setminus s_2$	0	+	-	\top	$s_1 \setminus s_2$	0	+	-	\top
0	0	-	+	\top	0	0	0	0	0
+	+	- if $u_1 < u_2$ + if $u_2 < u_1$ \top otherwise	+	\top	+	0	1	0	1
-	-	-	- if $u_2 < u_1$ + if $u_1 < u_2$ \top otherwise	\top	-	0	0	1	1
\top	\top	\top	\top	\top	\top	0	1	1	1

σ_-				
$s_1 \setminus s_2$	0	+	-	\top
0	0	0	0	0
+	0	0	1	1
-	-	0	1	0
\top	\top	0	1	1

Fig. 5. Operators used in the types of the primitives of Figure 4.

Proof. For product, we have $p = \text{ufp}(x \times y) - \text{ufp}(e_\times)$ with $e_\times = x \cdot \varepsilon(y) + y \cdot \varepsilon(x) + \varepsilon(x) \cdot \varepsilon(y)$. Let $u = u_1 + u_2 + 1 = \mathcal{U}_\times(s_1, u_1, s_2, u_2)$. We have, by definition of ufp , $2^{u_1} \leq x < 2^{u_1+1}$ and $2^{u_2} \leq y < 2^{u_2+1}$. Then e_\times may be bound by

$$\begin{aligned} e_\times &< 2^{u_1+1} \cdot 2^{u_2-p_2+1} + 2^{p_2+1} \cdot 2^{u_1-p_1+1} + 2^{u_1-p_1+1} \cdot 2^{u_2-p_2+1} \\ &= 2^{u_1+u_2-p_2+2} + 2^{u_1+u_2-p_1+2} + 2^{u_1+u_2-p_1-p_2+2} \end{aligned} \quad (9)$$

Since $u_1 + u_2 - p_1 - p_2 + 2 < u_1 + u_2 - p_1 + 2$ and $u_1 + u_2 - p_1 - p_2 + 2 < u_1 + u_2 - p_2 + 2$, we may get rid of the last term of Equation (9) and we obtain that

$$\begin{aligned} \text{ufp}(e_\times) &< \max(u_1 + u_2 - p_1 + 2, u_1 + u_2 - p_2 + 2) + \iota(p_1, p_2) \\ &\leq \max(u_1 + u_2 - p_1 + 1, u_1 + u_2 - p_2 + 1) + \iota(p_1, p_2) \end{aligned}$$

Let us write $p = \max(u_1 + u_2 - p_1 + 1, u_1 + u_2 - p_2 + 1) - \iota(p_1, p_2) = \mathcal{P}_\times(s_1, u_1, p_1 s_2, u_2, p_2)$. Then $u = \mathcal{U}_\times(s_1, u_1, s_2, u_2)$, $p = \mathcal{P}_\times(s_1, u_1, p_1 s_2, u_2, p_2)$ and $\varepsilon(z) < 2^{u-p+1}$. \square

Note that, by reasoning on the exponents of the values, the constraints resulting from a product become linear. The equations for subtraction and division are almost identical to the equations for addition and product, respectively. Note that the result of a division has one less bit than the result of a product. This is due to the fact that, even if the operands are finite numbers, the result of the division may be irrational and possibly needs to be truncated. We conclude this section with the following theorem which summarize the properties of the types of the result of the four elementary operations.

$$\begin{array}{c}
\frac{|r - v_{\mathbb{F}}| < 2^{u-p+1} \quad \text{ufp}(r) \leq u \quad \text{sign}(v_{\mathbb{F}}) < s}{r\{s, u, p\} \rightarrow_{\mathbb{F}} v_{\mathbb{F}}} \quad (\text{FVal}) \qquad \frac{v_{\mathbb{R}} = r}{r\{s, u, p\} \rightarrow_{\mathbb{R}} v_{\mathbb{R}}} \quad (\text{RVal}) \\
\\
\frac{e_0 \rightarrow e'_0}{e_0 * e_1 \rightarrow e'_0 * e_1} \quad (\text{Op1}) \qquad \frac{e_1 \rightarrow e'_1}{v * e_1 \rightarrow v * e'_1} \quad (\text{Op2}) \qquad * \in \{+, -, \times, \div, +-, --, \times-, \div-\} \\
\\
\frac{v = v_0 * v_1}{v_0 * v_1 \rightarrow v} \quad (\text{Op}) \qquad * \in \{+, -, \times, \div, +-, --, \times-, \div-\} \qquad \text{rec } f \ x.e \rightarrow \lambda x.e(\text{rec } f \ x.e/f) \quad (\text{REC}) \\
\\
\frac{e_0 \rightarrow e'_0}{e_0 \bowtie e_1 \rightarrow e'_0 \bowtie e_1} \quad (\text{Cmp1}) \qquad \frac{e_1 \rightarrow e'_1}{v \bowtie e_1 \rightarrow v \bowtie e'_1} \quad (\text{Cmp2}) \qquad \bowtie \in \{<_{\{u,p\}}, >_{\{u,p\}}, <, >\} \\
\\
\frac{b = (v_0^{\mathbb{F}} - v_1^{\mathbb{F}} \bowtie 2^{u-p+1})}{v_0 \bowtie_{\{u,p\}} v_1 \rightarrow_{\mathbb{F}} b} \quad (\text{FCmp}) \qquad \frac{b = (v_0 \bowtie v_1)}{v_0 \bowtie_{\{u,p\}} v_1 \rightarrow_{\mathbb{R}} b} \quad (\text{RCmp}) \qquad \bowtie \in \{<_{\{u,p\}}, >_{\{u,p\}}\} \\
\\
\frac{b = v_0 \bowtie v_1}{v_0 \bowtie v_1 \rightarrow b} \quad (\text{ACmp}) \qquad \bowtie \in \{=, \neq, <, >, \leq, \geq\} \\
\\
\frac{e_1 \rightarrow e'_1}{e_0 \ e_1 \rightarrow e_0 \ e'_1} \quad (\text{App1}) \qquad \frac{e_0 \rightarrow e'_0}{e_0 \ v \rightarrow e'_0 \ v} \quad (\text{App2}) \qquad (\lambda x.e) \ v \rightarrow e(v/x) \quad (\text{Red}) \\
\\
\frac{e_0 \rightarrow e'_0}{\text{if } e_0 \text{ then } e_1 \text{ else } e_2 \rightarrow \text{if } e'_0 \text{ then } e_1 \text{ else } e_2} \quad (\text{Cond}) \\
\\
\frac{v = \text{true}}{\text{if } v \text{ then } e_1 \text{ else } e_2 \rightarrow e_1} \quad (\text{CondTrue}) \qquad \frac{v = \text{false}}{\text{if } v \text{ then } e_1 \text{ else } e_2 \rightarrow e_2} \quad (\text{CondFalse})
\end{array}$$

Fig. 6. Operational semantics for our language.

Theorem 1. *Let x and y be two values such that $\varepsilon(x) < 2^{u_1-p_1+1}$ and $\varepsilon(y) < 2^{u_2-p_2+1}$ and let $*$ $\in \{+, -, \times, \div\}$ be an elementary operation. Let $z = x * y$, $u = \mathcal{U}_*(s_1, u_1, s_2, u_2)$ and $p = \mathcal{P}_*(s_1, u_1, p_1, s_2, u_2, p_2)$. Then $\varepsilon(z) < 2^{u+p-1}$.*

Proof. The cases of addition and product correspond to Lemma 1 and Lemma 2, respectively. The cases of subtraction and division are similar. \square

4 Soundness of the Type System

In this section, we introduce a subject reduction theorem proving the consistency of our type system. We use two operational semantics $\rightarrow_{\mathbb{F}}$ and $\rightarrow_{\mathbb{R}}$ for the finite precision and exact arithmetics, respectively. The exact semantics is used for proofs. Obviously, in practice, only the finite precision semantics is implemented. We write \rightarrow whenever a reduction rule holds for either $\rightarrow_{\mathbb{F}}$ or $\rightarrow_{\mathbb{R}}$ (in this case, we assume that the same semantics $\rightarrow_{\mathbb{F}}$ or $\rightarrow_{\mathbb{R}}$ is used in the lower and upper parts of the same sequent). Both semantics are displayed in Figure 6. They concern the subset of the language of Equation (3) which do not deal with types.

$$\begin{aligned}
\text{EvalExpr } \ni e ::= & r\{s, u, p\} \in \text{Real}_{u,p} \mid \mathbf{i} \in \text{Int} \mid \mathbf{b} \in \text{Bool} \mid \mathbf{id} \in \text{Id} \\
& \mid \text{if } e_0 \text{ then } e_1 \text{ else } e_2 \mid \lambda x.e \mid e_0 \ e_1 \mid \text{rec } f \ x.e \mid e_0 * e_1 \quad . \quad (10)
\end{aligned}$$

In Equation (10), $*$ denotes an arithmetic operator $* \in \{+, -, \times, \div, +_-, -_-, \times_-, \div_-\}$. In Figure 6, Rule (FVAL) of $\rightarrow_{\mathbb{F}}$ transforms a syntactic element describing a real number $\mathbf{r}\{\mathbf{s}, \mathbf{u}, \mathbf{p}\}$ in a certain format into a value $v_{\mathbb{F}}$. The finite precision value $v_{\mathbb{F}}$ is an approximation of \mathbf{r} with an error less than the ulp of $\mathbf{r}\{\mathbf{s}, \mathbf{u}, \mathbf{p}\}$. In the semantics $\rightarrow_{\mathbb{R}}$, the real number $\mathbf{r}\{\mathbf{s}, \mathbf{u}, \mathbf{p}\}$ simply produces the value \mathbf{r} without any approximation by Rule (RVal). Rules (Op1) and (Op2) evaluate the operands of some binary operation and Rule (Op) performs an operation $* \in \{+, -, \times, \div, +_-, -_-, \times_-, \div_-\}$ between two values v_0 and v_1 .

Rules (Cmp1), (Cmp2) and (ACmp) deal with comparisons. They are similar to Rules (Op1), (Op2) and (Op) described earlier. Note that the operators $<$, $>$, $=$, \neq concerned by Rule (ACmp) are polymorphic except that they do not accept arguments of type `real`. Rules (FCmp) and (RCmp) are for the comparison of `real` values. Rule (FCmp) is designed to avoid unstable tests by requiring that the distance between the two compared values is greater than the ulp of the format in which the comparison is done. With this requirement, a condition cannot be invalidated by the roundoff errors. Let us also note that, with this definition, $x <_{u,p} y \not\Rightarrow y >_{u,p} x$ or $x >_{u,p} y \not\Rightarrow y <_{u,p} x$. For the semantics $\rightarrow_{\mathbb{R}}$, Rule (RCmp) simply compares the exact values.

The other rules are standard and are identical in $\rightarrow_{\mathbb{F}}$ and $\rightarrow_{\mathbb{R}}$. Rules (App1), (App2) and (Red) are for applications and Rule (Rec) is for recursive functions. We write $e\langle v/x \rangle$ the term e in which v has been substituted to the free occurrences of x . Rules (Cond), (CondTrue) and (CondFalse) are for conditionals.

The rest of this section is dedicated to our subject reduction theorem. First of all, we need to relate the traces of $\rightarrow_{\mathbb{F}}$ and $\rightarrow_{\mathbb{R}}$. We introduce new judgments

$$\Gamma \models (e_{\mathbb{F}}, e_{\mathbb{R}}) : t . \quad (11)$$

Intuitively, Equation (11) means that expression $e_{\mathbb{F}}$ simulates $e_{\mathbb{R}}$ up to accuracy t . In this case, $e_{\mathbb{F}}$ is syntactically equivalent to $e_{\mathbb{R}}$ up to the values which, in $e_{\mathbb{F}}$, are approximations of the values of $e_{\mathbb{R}}$. The quantification of the approximation is given by type t .

Formally, \models is defined in Figure 7. These rules are similar to the typing rules of Figure 2 excepted that they operate on pairs $(e_{\mathbb{F}}, e_{\mathbb{R}})$. They are also designed for the language of Equation (10) and, consequently, deal with the elementary arithmetic operations $+$, $-$, \times and \div as well as the comparison operators. The difference between the rules of Figure 2 and Figure 7 is in Rule (VReal) which states that a `real` value $v_{\mathbb{R}}$ is correctly simulated by a value $v_{\mathbb{F}}$ up to accuracy `real` $\{\mathbf{s}, \mathbf{u}, \mathbf{p}\}$ if $|v_{\mathbb{R}} - v_{\mathbb{F}}| < 2^{u-p+1}$. It is easy to show, by examination of the rules of Figure 2 and Figure 7 that

$$\Gamma \models (e_{\mathbb{F}}, e_{\mathbb{R}}) : t \implies \Gamma \vdash e_{\mathbb{F}} : t . \quad (12)$$

We introduce now Lemma 3 which asserts the soundness of the type system for one reduction step. Basically, this lemma states that types are preserved by reduction and that concerning the values of type `real`, the distance between the finite precision value and the exact value is less than the ulp given by the type.

Lemma 3 (Weak subject reduction). *If $\Gamma \models (e_{\mathbb{F}}, e_{\mathbb{R}}) : t$ and if $e_{\mathbb{F}} \rightarrow_{\mathbb{F}} e'_{\mathbb{F}}$ and $e_{\mathbb{R}} \rightarrow_{\mathbb{R}} e'_{\mathbb{R}}$ then $\Gamma \models (e'_{\mathbb{F}}, e'_{\mathbb{R}}) : t$.*

$$\begin{array}{c}
\frac{}{\Gamma \models (i, i) : \text{int}} \quad (\text{INT}) \qquad \frac{}{\Gamma \models (b, b) : \text{bool}} \quad (\text{BOOL}) \qquad \frac{\Gamma(\text{id}) = t}{\Gamma \models (\text{id}, \text{id}) : t} \quad (\text{ID}) \\
\\
\frac{\text{sign}(r) < s \quad \text{ufp}(r) \leq u}{\Gamma \models (r\{s, u, p\}, r\{s, u, p\}) : \text{real}\{s, u, p\}} \quad (\text{SREAL}) \qquad \frac{|v_{\mathbb{R}} - v_{\mathbb{F}}| < 2^{u-p+1}}{\Gamma \models (v_{\mathbb{F}}, v_{\mathbb{R}}) : \text{real}\{s, u, p\}} \quad (\text{VREAL}) \\
\\
\frac{\Gamma \models (e_{1\mathbb{F}}, e_{1\mathbb{R}}) : \text{real}\{s_1, u_1, p_1\} \quad \Gamma \models (e_{2\mathbb{F}}, e_{2\mathbb{R}}) : \text{real}\{s_1, u_1, p_1\} \quad * \in \{+, -, \times, \div\}}{\Gamma \models (e_{1\mathbb{F}} * e_{2\mathbb{F}}, e_{1\mathbb{R}} * e_{2\mathbb{R}}) : \text{real}\{\mathcal{S}_*(s_1, u_1, s_2, u_2), \mathcal{U}_*(s_1, u_1, s_2, u_2), \mathcal{P}_*(s_1, u_1, p_1, s_2, u_2, p_2)\}} \quad (\text{ROP}) \\
\\
\frac{\Gamma \models (e_{1\mathbb{F}}, e_{1\mathbb{R}}) : \text{real}\{s_1, u, p + 1\} \quad \Gamma \models (e_{2\mathbb{F}}, e_{2\mathbb{R}}) : \text{real}\{s_1, u, p + 1\} \quad * \in \{<, >\}}{\Gamma \models (e_{1\mathbb{F}} \bowtie_{u,p} e_{2\mathbb{F}}, e_{1\mathbb{R}} \bowtie_{u,p} e_{2\mathbb{R}}) : \text{bool}} \quad (\text{RCMP}) \\
\\
\frac{\Gamma \models (e_{1\mathbb{F}}, e_{1\mathbb{R}}) : \text{int} \quad \Gamma \models (e_{2\mathbb{F}}, e_{2\mathbb{R}}) : \text{int} \quad * \in \{+-, -\times, \times-, \div-\}}{\Gamma \models (e_{1\mathbb{F}} *_- e_{2\mathbb{F}}, e_{1\mathbb{R}} *_- e_{2\mathbb{R}}) : \text{int}} \quad (\text{INTOP}) \\
\\
\frac{\Gamma \models (e_{1\mathbb{F}}, e_{1\mathbb{R}}) : t \quad \Gamma \models (e_{2\mathbb{F}}, e_{2\mathbb{R}}) : t \quad t \neq \text{real}\{s, u, p\} \quad \bowtie \in \{=, \neq, <, >, \leq, \geq\}}{\Gamma \models (e_{1\mathbb{F}} \bowtie e_{2\mathbb{F}}, e_{1\mathbb{R}} \bowtie e_{2\mathbb{R}}) : \text{bool}} \quad (\text{ACMP}) \\
\\
\frac{\Gamma \models (e_{0\mathbb{F}}, e_{0\mathbb{R}}) : \text{bool} \quad \Gamma \models (e_{1\mathbb{F}}, e_{1\mathbb{R}}) : t_1 \quad \Gamma \models (e_{2\mathbb{F}}, e_{2\mathbb{R}}) : t_2 \quad t = t_1 \sqcup t_2}{\Gamma \models (\text{if } e_{0\mathbb{F}} \text{ then } e_{1\mathbb{F}} \text{ else } e_{2\mathbb{F}}, \text{if } e_{0\mathbb{R}} \text{ then } e_{1\mathbb{R}} \text{ else } e_{2\mathbb{R}}) : t} \quad (\text{COND}) \\
\\
\frac{\Gamma, x : t_1 \models (e_{\mathbb{F}}, e_{\mathbb{R}}) : t_2}{\Gamma \models (\lambda x. e_{\mathbb{F}}, \lambda x. e_{\mathbb{R}}) : \Pi x : t_1. t_2} \quad (\text{ABS}) \qquad \frac{\Gamma, x : t_1, f : \Pi y : t_1. t_2 \models (e_{\mathbb{F}}, e_{\mathbb{R}}) : t_2}{\Gamma \models (\text{rec } f x. e_{\mathbb{F}}, \text{rec } f x. e_{\mathbb{R}}) : \Pi x : t_1. t_2} \quad (\text{REC}) \\
\\
\frac{\Gamma \models (e_{1\mathbb{F}}, e_{1\mathbb{R}}) : \Pi x : t_0. t_1 \quad \Gamma \models (e_{2\mathbb{F}}, e_{2\mathbb{R}}) : t_2 \quad t_2 \sqsubseteq t_0}{\Gamma \models (e_{1\mathbb{F}} e_{2\mathbb{F}}, e_{1\mathbb{R}} e_{2\mathbb{R}}) : t_1[x \mapsto e_2]} \quad (\text{APP})
\end{array}$$

Fig. 7. Simulation relation \models used in our subject reduction theorem.

Proof. By induction on the structure of expressions and case examination on the possible transition rules of Figure 6.

- If $e_{\mathbb{F}} \equiv e_{\mathbb{R}} \equiv r\{s, u, p\}$ then $\Gamma \models (r\{s, u, p\}, r\{s, u, p\}) : \text{real}\{s, u, p\}$ and, from the reduction rules (FVal) and (RVal) of Figure 6, $r\{s, u, p\} \rightarrow_{\mathbb{F}} v_{\mathbb{F}}$ and $r\{s, u, p\} \rightarrow_{\mathbb{R}} v_{\mathbb{R}}$ with $|v_{\mathbb{F}} - v_{\mathbb{R}}| < 2^{u-p+1}$. So $\Gamma \models (v_{\mathbb{F}}, v_{\mathbb{R}}) : \text{real}\{s, u, p\}$.
- If $e_{\mathbb{F}} \equiv e_{0\mathbb{F}} * e_{1\mathbb{F}}$ and $e_{\mathbb{R}} \equiv e_{0\mathbb{R}} * e_{1\mathbb{R}}$ then several cases must be distinguished.
 - If $e_{\mathbb{F}} \equiv v_{0\mathbb{F}} * v_{1\mathbb{F}}$ and $e_{\mathbb{R}} \equiv v_{0\mathbb{R}} * v_{1\mathbb{R}}$ then, by induction hypothesis, $\Gamma \models (v_{0\mathbb{F}}, v_{0\mathbb{R}}) : \text{real}\{s_0, u_0, p_0\}$, $\Gamma \models (v_{1\mathbb{F}}, v_{1\mathbb{R}}) : \text{real}\{s_1, u_1, p_1\}$ and, consequently, from Rule (VREAL),

$$|v_{0\mathbb{R}} - v_{0\mathbb{F}}| < 2^{u_0-p_0+1} \quad \text{and} \quad |v_{1\mathbb{R}} - v_{1\mathbb{F}}| < 2^{u_1-p_1+1} \quad . \quad (13)$$

Following Figure 4, the type t of e is

$$\begin{aligned}
t &= (\Pi s_1 : \text{int}, u_1 : \text{int}, p_1 : \text{int}, s_2 : \text{int}, u_2 : \text{int}, p_2 : \text{int}. \\
&\quad \text{real}\{s_1, u_1, p_1\} \rightarrow \text{real}\{s_2, u_2, p_2\} \rightarrow \\
&\quad \rightarrow \text{real}\{\mathcal{S}_*(s_1, u_1, s_2, u_2), \mathcal{U}_*(s_1, u_1, s_2, u_2), \mathcal{P}_*(s_1, u_1, p_1, s_2, u_2, p_2)\} \\
&\quad) s_1 u_1 p_1 s_2 u_2 p_2 , \\
&= \text{real}\{\mathcal{S}_*(s_1, u_1, s_2, u_2), \mathcal{U}_*(s_1, u_1, s_2, u_2), \mathcal{P}_*(s_1, u_1, p_1, s_2, u_2, p_2)\} \\
&= \text{real}\{s, u, p\}
\end{aligned}$$

By Rule (Op), $e \rightarrow_{\mathbb{F}} v_{\mathbb{F}}$ and $e \rightarrow_{\mathbb{R}} v_{\mathbb{R}}$ and, by Theorem 1, with the assumptions of Equation (13), we know that $|v_{\mathbb{R}} - v_{\mathbb{F}}| < 2^{u-p+1}$. Consequently, $\Gamma \models (v_{\mathbb{F}}, v_{\mathbb{R}}) : \mathbf{real}\{\mathbf{s}, \mathbf{u}, \mathbf{p}\}$.

- If $e_{\mathbb{F}} \equiv v_{0\mathbb{F}} * v_{1\mathbb{F}}$ and $e_{\mathbb{R}} \equiv v_{0\mathbb{R}} * v_{1\mathbb{R}}$ with $\Gamma \models (v_0, v_1) : \mathbf{int}$ then, by Rule (Op), $e \rightarrow (v, v)$ and, by Equation (7), $\Gamma \vdash v : \mathbf{int}$. If $e \equiv e_0 * e_1$ then, by Rule (Op1), $e \rightarrow e_0 * e_1'$ and we conclude by induction hypothesis. The case $e \equiv e_0 * v_1$ is similar to the former one.
- If $e_{\mathbb{F}} \equiv e_{0\mathbb{F}} \bowtie_{u,p} e_{1\mathbb{F}}$ and $e_{\mathbb{R}} \equiv e_{0\mathbb{R}} \bowtie_{u,p} e_{1\mathbb{R}}$ then several cases have to be examined.
 - If $e_{\mathbb{F}} \equiv v_{0\mathbb{F}} \bowtie_{u,p} v_{1\mathbb{F}}$ and $e_{\mathbb{R}} \equiv v_{0\mathbb{R}} \bowtie_{u,p} v_{1\mathbb{R}}$ then by rules (FCmp) and (RCmp) $e_{\mathbb{F}} \rightarrow_{\mathbb{F}} b_{\mathbb{F}}$, $e_{\mathbb{R}} \rightarrow_{\mathbb{R}} b_{\mathbb{R}}$ with $b_{\mathbb{F}} = v_{0\mathbb{F}} - v_{1\mathbb{F}} \bowtie_{\{u,p\}} 2^{u-p+1}$ and $b_{\mathbb{R}} = v_{0\mathbb{R}} - v_{1\mathbb{R}} \bowtie_{\{u,p\}} 0$. By rule (RCmp) of Figure 7, $\Gamma \models (v_{0\mathbb{F}}, v_{1\mathbb{F}}) : \mathbf{real}\{\mathbf{s}, \mathbf{u}, \mathbf{p}\}$ and $\Gamma \models (v_{0\mathbb{R}}, v_{1\mathbb{R}}) : \mathbf{real}\{\mathbf{s}, \mathbf{u}, \mathbf{p}\}$. Consequently, $|v_{0\mathbb{R}} - v_{0\mathbb{F}}| < 2^{u-p+1}$ and $|v_{1\mathbb{R}} - v_{1\mathbb{F}}| < 2^{u-p+1}$. By combining thz former equations, we obtain that $|(v_{0\mathbb{R}} - v_{1\mathbb{R}}) - (v_{0\mathbb{F}} - v_{1\mathbb{F}})| < 2^{u-p}$. Consequently, $b_{\mathbb{F}} = b_{\mathbb{R}}$ and we conclude that $\Gamma \models (b_{\mathbb{F}}, b_{\mathbb{R}}) : \mathbf{bool}$.
 - The other cases for $e_{\mathbb{F}} \equiv e_{0\mathbb{F}} \bowtie_{u,p} e_{1\mathbb{F}}$ are similar to the cases $e_{\mathbb{F}} \equiv v_{0\mathbb{F}} * v_{1\mathbb{F}}$ examined previously.
- The other cases simply follow the structure of the terms, by application of the induction hypothesis. \square

Let $\rightarrow_{\mathbb{F}}^*$ (resp. $\rightarrow_{\mathbb{R}}^*$) denote the reflexive transitive closure of $\rightarrow_{\mathbb{F}}$ (resp. $\rightarrow_{\mathbb{R}}$). Theorem 2 expresses the soundness of our type system for sequences of reduction of arbitrary length.

Theorem 2 (Subject reduction). *If $\Gamma \models (e_{\mathbb{F}}, e_{\mathbb{R}}) : t$ and if $e_{\mathbb{F}} \rightarrow_{\mathbb{F}}^* e'_{\mathbb{F}}$ and $e_{\mathbb{R}} \rightarrow_{\mathbb{R}}^* e'_{\mathbb{R}}$ then $\Gamma \models (e'_{\mathbb{F}}, e'_{\mathbb{R}}) : t$.*

Proof. By induction on the length of the reduction sequence, using Lemma 3. \square

Theorem 2 asserts the soundness of our type system. It states that the evaluation of an expression of type $\mathbf{real}\{\mathbf{s}, \mathbf{u}, \mathbf{p}\}$ yields a result of accuracy 2^{u-p+1} .

5 Experiments

In this section, we report some experiments showing how our type system behaves in practice. Section 5.1 presents Num1 implementations of usual mathematical formulas while Section 5.2 introduce a larger example demonstrating the expressive power of our type system.

5.1 Usual Mathematic Formulas

Our first examples concern usual mathematic formulas, to compute the volume of geometrical objects or formulas related to polynomials. These examples aim at showing that usual mathematical formulas are typable in our system. We start with the volume of the sphere and of the cone.

```

> let sphere r = (4.0 / 3.0) * 3.1415926{+,1,20} * r * r * r ;;
val sphere : real{'a','b','c'} -> real{<expr>,<expr>,<expr>} = <fun>

> sphere 1.0 ;;
- : real{+,7,20} = 4.188

> let cone r h = (3.1415926{+,1,20} * r * r * h) / 3.0 ;;
val cone : real{'a','b','c'} -> real{'a','b','c'}
      -> real{<expr>,<expr>,<expr>} = <fun>

> cone 1.0 1.0 ;;
- : real{+,4,20} = 1.0472

```

We repeatedly define the function `sphere` with more precision in order to show the impact on the accuracy of the results. Note that the results now have 15 digits instead of the former 5 digits.

```

> let sphere r = (4.0 / 3.0) * 3.1415926535897932{+,1,53} * r * r * r ;;
val sphere : real{'a','b','c'} -> real{<expr>,<expr>,<expr>} = <fun>

> sphere 1.0 ;;
- : real{+,7,52} = 4.1887902047863

```

The next examples concern polynomials. We start with the computation of the discriminant of a second degree polynomial.

```

> let discriminant a b c = b * b - 4.0 * a * c ;;
val discriminant : real{'a','b','c'} -> real{'d','e','f'} -> real{'g','h','i'}
      -> real{<expr>,<expr>,<expr>} = <fun>

> discriminant 2.0 -11.0 15.0 ;;
- : real{+,8,52} = 1.0000000000000

```

Our last example concerning usual formulas is the Taylor series development of the sine function. In the code below, observe that the accuracy of the result is correlated to the accuracy of the argument. As mentioned in Section 2, error methods are neglected, only the errors due to the finite precision are calculated (indeed, $\sin \frac{\pi}{8} = 0.382683432\dots$).

```

let sin x = x - ((x * x * x) / 3.0) + ((x * x * x * x * x) / 120.0) ;;
val sin : real{'a','b','c'} -> real{<expr>,<expr>,<expr>} = <fun>

> sin (3.14{1,6} / 8.0) ;;
- : real{*,0,6} = 0.3

> sin (3.14159{1,18} / 8.0) ;;
- : real{*,0,18} = 0.37259

```

5.2 Newton-Raphson Method

In this section, we introduce a larger example to compute the zero of a function using the Newton-Raphson method. This example, which involves several higher order functions, shows the expressiveness of our type system. In the programming session below, we first define a higher order function `deriv` which takes as argument a function and computes its numerical derivative at a given point. Then we define a function `g` and compute the value of its derivative at point 2.0. Next, by partial application, we build a function computing the derivative

of `g` at any point. Finally, we define a function `newton` which searches the zero of a function. The `newton` function is also an higher order function taking as argument the function for which a zero has to be found and its derivative.

```
> let deriv f x h = ((f (x + h)) - (f x)) / h ;;
val deriv : (real{<expr>,<expr>,<expr>} -> real{'a','b','c'})
            -> real{<expr>,<expr>,<expr>} -> real{'d','e','f'}
            -> real{<expr>,<expr>,<expr>} = <fun>

> let g x = (x*x) - (5.0*x) + 6.0 ;;
val g : real{'a','b','c'} -> real{<expr>,<expr>,<expr>} = <fun>

> deriv g 2.0 0.01 ;;
- : real{*,5,51} = -0.990000000000000

> let gprime x = deriv g x 0.01 ;;
val gprime : real{<expr>,<expr>,<expr>} -> real{<expr>,<expr>,<expr>} = <fun>

> let rec newton x xold f fprime = if ((x-xold)<0.01{*,10,20}) then x
                                   else newton (x-((f x)/(fprime x))) x f fprime ;;
val newton : real{*,10,21} -> real{0,10,20} -> (real{*,10,21} -> real{'a','b','c'})
            -> (real{*,10,21} -> real{'d','e','f'}) -> real{*,10,21} = <fun>

> newton 9.0 0.0 g gprime ;;
- : real{*,10,21} = 5.771
```

We call the `newton` function with our function `g` and its derivative computed by partial application of the `deriv` function. We obtain a root of our polynomial `g` with a guaranteed accuracy. Note that while Newton-Raphson method converges quadratically in the reals, numerical errors may perturb the process [4].

6 Conclusion

In this article, we have introduced a dependent type system able to infer the accuracy of numerical computations. Our type system allows one to type non-trivial programs corresponding to implementations of classical numerical analysis methods. Unstable computations are rejected by the type system. The consistency of typed programs is ensured by a subject reduction theorem. To our knowledge, this is the first type system dedicated to numerical accuracy. We believe that this approach has many advantages going from early debugging to compiler optimizations. Indeed, we believe that the usual type `float` proposed by usual ML implementations, and which is a simple clone of the type `int`, is too poor for numerical computations. We also believe that this approach is a credible alternative to static analysis techniques for numerical precision [6, 9, 18]. For the developer, our type system introduces few changes in the programming style, limited to giving the accuracy of the inputs of the accuracy of comparisons to allow the typing of certain recursive functions.

A first perspective to the present work is the implementation of a compiler for Num1. We aim at using the type information to select the most appropriate formats (the IEEE754 formats of Figure 1, multiple precisions numbers of the GMP library when needed or requested by the user or fixed-point numbers.) At longer term, we also aim at introducing safe compile-time optimizations based on type preservation: an expression may be safely (from the accuracy point of view)

substituted to another expression as long as both expressions are mathematically equivalent and that the new expression has a greater type than the older one in the sense of Equation (6). Finally, a second perspective is to integrate our type system into other applicative languages. In particular, it would be of great interest to have such a type system inside a language used to build critical embedded systems such as the synchronous language `Lustre` [3]. In this context numerical accuracy requirements are strong and difficult to obtain. Our type system could be integrated naturally inside `Lustre` or similar languages.

References

1. ANSI/IEEE: IEEE Standard for Binary Floating-point Arithmetic (2008)
2. Atkinson, K.: An Introduction to Numerical Analysis, 2nd Edition. Wiley (1989)
3. Caspi, P., Pilaud, D., Halbwegs, N., Plaice, J.: Lustre: A declarative language for programming synchronous systems. In: POPL. pp. 178–188. ACM Press (1987)
4. Damouche, N., Martel, M., Chapoutot, A.: Impact of accuracy optimization on the convergence of numerical iterative methods. In: LOPSTR’15. LNCS, vol. LNCS 9527, pp. 1–18. Springer (2015)
5. Damouche, N., Martel, M., Chapoutot, A.: Improving the numerical accuracy of programs by automatic transformation. STTT 19(4), 427–448 (2017)
6. Darulova, E., Kuncak, V.: Sound compilation of reals. In: POPL’14. pp. 235–248. ACM (2014)
7. Denis, C., de Oliveira Castro, P., Petit, E.: Verificarlo: Checking floating point accuracy through monte carlo arithmetic. In: ARITH’16. pp. 55–62. IEEE (2016)
8. Franco, A.D., Guo, H., Rubio-González, C.: A comprehensive study of real-world numerical bug characteristics. In: ASE. pp. 509–519. IEEE (2017)
9. Goubault, E.: Static analysis by abstract interpretation of numerical programs and systems, and FLUCTUAT. In: SAS. LNCS, vol. 7935, pp. 1–3. Springer (2013)
10. Graphics, M.: Algorithmic C Datatypes, software version 2.6 edn. (2011), <http://www.mentor.com/esl/catapult/algorithmic>
11. Lam, M.O., Hollingsworth, J.K., de Supinski, B.R., LeGendre, M.P.: Automatically adapting programs for mixed-precision floating-point computation. In: Supercomputing, ICS’13. pp. 369–378. ACM (2013)
12. Martel, M.: Floating-point format inference in mixed-precision. In: NFM. LNCS, vol. 10227, pp. 230–246 (2017)
13. Milner, R., Harper, R., MacQueen, D., Tofte, M.: The Definition of Standard ML. MIT Press (1997)
14. Panchekha, P., Sanchez-Stern, A., Wilcox, J.R., Tatlock, Z.: Automatically improving accuracy for floating point expressions. In: PLDI. pp. 1–11. ACM (2015)
15. Pierce, B.C.: Types and programming languages. MIT Press (2002)
16. Pierce, B.C. (ed.): Advanced Topics in Types and Programming Languages. MIT Press (2004)
17. Rubio-Gonzalez, C., Nguyen, C., Nguyen, H.D., Demmel, J., Kahan, W., Sen, K., Bailey, D.H., Iancu, C., Hough, D.: Precimonious: tuning assistant for floating-point precision. In: HPCNSA. pp. 27:1–27:12. ACM (2013)
18. Solovyev, A., Jacobsen, C., Rakamaric, Z., Gopalakrishnan, G.: Rigorous estimation of floating-point round-off errors with symbolic taylor expansions. In: FM. LNCS, vol. 9109, pp. 532–550. Springer (2015)