

# End-To-End Formal Verification of a Fast and Accurate Floating-Point Approximation

Florian Faissole ✉ 

Mitsubishi Electric R&D Centre Europe, 35700 Rennes, France

Paul Geneau de Lamarlière ✉ 

Mitsubishi Electric R&D Centre Europe, 35700 Rennes, France

Université Paris-Saclay, CNRS, ENS Paris-Saclay, Inria, LMF, 91190 Gif-sur-Yvette, France

Guillaume Melquiond ✉ 

Université Paris-Saclay, CNRS, ENS Paris-Saclay, Inria, LMF, 91190 Gif-sur-Yvette, France

---

## Abstract

Designing an efficient yet accurate floating-point approximation of a mathematical function is an intricate and error-prone process. This warrants the use of formal methods, especially formal proof, to achieve some degree of confidence in the implementation. Unfortunately, the lack of automation or its poor interplay with the more manual parts of the proof makes it way too costly in practice. This article revisits the issue by proposing a methodology and some dedicated automation, and applies them to the use case of a faithful *binary64* approximation of exponential. The peculiarity of this use case is that the target of the formal verification is not a simple modeling of an external code; it is an actual floating-point function defined in the logic of the Coq proof assistant, which is thus usable inside proofs once its correctness has been fully verified. This function presents all the attributes of a state-of-the-art implementation: bit-level manipulations, large tables of constants, obscure floating-point transformations, exceptional values, etc. This function has been integrated into the proof strategies of the CoqInterval library, bringing a 20× speedup with respect to the previous implementation.

**2012 ACM Subject Classification** Software and its engineering → Formal software verification; Theory of computation → Interactive proof systems; Theory of computation → Automated reasoning; Mathematics of computing → Mathematical software performance; Mathematics of computing → Interval arithmetic

**Keywords and phrases** Program verification, floating-point arithmetic, formal proof, automated reasoning, mathematical library

**Digital Object Identifier** 10.4230/LIPIcs.ITP.2024.14

**Supplementary Material** *Software*: <https://gitlab.inria.fr/coqinterval/interval.git> [6]  
archived at [swh:1:dir:78da3e6e98b7ef018180119255ce1e10a048cc88](https://swh.1:dir:78da3e6e98b7ef018180119255ce1e10a048cc88)

**Funding** *Guillaume Melquiond*: This work was partly supported by the NuSCAP project (ANR-20-CE48-0014) of the French national research agency (ANR).

## 1 Introduction

A mathematical library is a set of floating-point functions that are designed to approximate mathematical functions. They are used in various domains, ranging from engineering to scientific computing and experimental mathematics. For such applications, these functions are required to be both accurate and fast to compute. To meet those requirements, the code of such a floating-point function is usually intricate and its correctness is far from trivial [14]. This warrants verifying the latter formally, which can be long and tedious [8, 9].



© Florian Faissole, Paul Geneau de Lamarlière, and Guillaume Melquiond;  
licensed under Creative Commons License CC-BY 4.0

15th International Conference on Interactive Theorem Proving (ITP 2024).

Editors: Yves Bertot, Temur Kutsia, and Michael Norrish; Article No. 14; pp. 14:1–14:18

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

Formally verified mathematical libraries can even be used for formal proofs. An example of such usage is the `CoqInterval` library,<sup>1</sup> which provides a set of strategies for the Coq proof assistant that automate the verification of enclosures of real-valued expressions. It is based on a formalization of rigorous polynomial approximations that are computed using an interval arithmetic with floating-point bounds [11]. Originally, the floating-point computations were performed one bit at a time in the logic of the Coq system. Later, support for 63-bit integers was added to Coq to speed up computation [5]. Even then, a formal verification of the following approximation of Siegfried Rump’s integral – an example known to cause computer algebra systems to struggle due to the large number of oscillations of the integrand – would still take minutes to complete:

$$\int_0^8 \sin(x + \exp x) dx = 0.3474 \pm 10^{-6}.$$

Indeed, proving this approximation requires computing polynomial approximations of the integrand on numerous subintervals of  $[0; 8]$ , which itself requires computing enclosures of the sine and exponential functions on many inputs.

To improve performance, recent work has added support for performing hardware floating-point computations inside Coq proofs [12]. These built-in operations are axiomatized with a purely computational specification, which has been formally proved to comply with the IEEE 754 standard thanks to the `Flocq` library [3]. This makes it possible to both trust the specification and use it for formal reasoning. Thanks to this new feature, the time needed to verify the above approximation is reduced to just a few seconds.

This remains much slower than what could be achieved by state-of-the-art libraries [17, §12]. Part of the reason is that we are performing computations inside the logic system, but also that the code itself uses hardware floating-point numbers in a very naive way. To make `CoqInterval` even more suitable for this kind of numerically intensive proofs, we would like to improve on this last point by providing it with a mathematical library that is specialized for hardware floating-point numbers and that is also formally verified. The work presented in this article focuses on the implementation of the exponential function.

## 1.1 The original `CoqInterval` implementation

Prior to this work, `CoqInterval` would use a single algorithm for the exponential function, but instantiated twice: once for floating-point numbers computed in hardware and once for floating-point numbers slowly emulated in the logic of Coq.<sup>2</sup> This was made possible thanks to a suitable abstraction of floating-point arithmetic [12]. Having a single algorithm, and thus a single proof of correctness, made the large formalization effort that went into adding hardware computations to `CoqInterval` much less tedious.

`CoqInterval`’s original algorithm for computing an enclosure of  $\exp x$  goes as follows. First, using the following mathematical identities, an argument reduction brings the input  $x$  into the interval  $[-2^{-8}; 0]$ :

$$\begin{aligned} \exp x &= (\exp(-x))^{-1} \\ \exp x &= (\exp(x/2))^2 \end{aligned} \tag{1}$$

<sup>1</sup> <https://coqinterval.gitlabpages.inria.fr/>

<sup>2</sup> While slow, this emulation of floating-point arithmetic is still useful for proofs that require more than the 53 bits of precision provided by the *binary64* format.

■ **Listing 1** Guaranteed approximation of exponential in OCaml. The output is a pair of floating-point numbers that enclose  $\exp x$ . Symbols `invLog2_64`, `log2_64h`, `log2_64l`, `cst`, `p1`, `p2`, etc, are predefined floating-point literals.

```
let iexp x =
  if x < -0x1.74385446d71c4p9 then (0., 0x1.p-1074) else
  if x > 0x1.62e42fefa39efp9 then (0x1.fffffffffff2ap1023, infinity) else
  let k' = x *. invLog2_64 +. 0x1.8p52 in
  let k = k' -. 0x1.8p52 in
  let t = (x -. k *. log2_64h) -. k *. log2_64l in
  let y = t *. (p1 +. t *. (p2 +. t *. (... + t *. p5))) in
  let ki = int_of_float k' - 0x18000000000000 in
  let p0 = cst.(ki land 63) in
  let d = 0x1.25p-57 in
  let lb = p0 +. (p0 *. y -. d) in
  let ub = p0 +. (p0 *. y +. d) in
  next_down (ldexp lb (ki asr 6)), next_up (ldexp ub (ki asr 6))
```

Second, the alternating series  $\exp(-x) = \sum(-x)^n/n!$  is computed using interval arithmetic to a high enough order. Thanks to the use of interval arithmetic and an alternating series, the algorithm is guaranteed to compute an enclosure of the real number  $\exp x$ . Third, the argument reduction is reversed to reconstruct the final interval result.

By suitably choosing the order of truncation of the series, one can obtain arbitrarily tight enclosures of  $\exp x$ , assuming that the precision of the floating-point arithmetic used to compute the interval bounds can be made accordingly large. This property is invaluable when used in conjunction with the original multi-precision floating-point arithmetic of CoqInterval. But for hardware floating-point numbers and their fixed precision of 53 bits, the property is pointless. The inadequacies of the implementation of `exp` then become prominent. First, Equation (1) means that computing an enclosure of  $\exp x$  is not constant time, but proportional to the magnitude of  $x$ . Second, an alternating series is the worst way of approximating a value, as part of the computations performed at order  $i$  are immediately canceled by those at order  $i + 1$  and thus have been performed in vain. Third, while interval arithmetic is correct by construction, hence very proof-friendly, it performs twice as many floating-point operations as needed.

## 1.2 The whole new algorithm

An approximation of the exponential function, as found in usual mathematical libraries, does not suffer from these defects, as it is generally implemented along the following guidelines [14, §6.2]. It would first perform a constant-time argument reduction using the following mathematical identity:

$$\exp x = \exp(x - k \cdot \ln 2) \cdot 2^k \quad \text{with } k = \lceil x / \ln 2 \rceil \in \mathbb{Z}.$$

where  $\lceil \cdot \rceil$  denotes the nearest integer. Then, a low-degree polynomial approximation of  $\exp$  around 0 would be evaluated. Finally, the result reconstruction is trivial, as it is a multiplication by a power of two. The whole algorithm amounts to just a few tens of operations; it is thus extremely fast.

Listing 1 shows the implementation we have devised, represented as an OCaml function for readability. (Its translation to Coq's  $\lambda$ -calculus is straightforward.) Given a finite floating-point number  $x$ , the code computes a pair of floating-point numbers enclosing  $\exp x$ . In particular, it uses the functions `next_down` and `next_up` to compute the predecessor and successor of a floating-point number.

While the code seems to contain useless, if not adverse, floating-point computations, this is not the case. For example, `k` looks like it could be directly computed as `x *. invLog2_64` by canceling `0x1.8p52 -. 0x1.8p52`. This optimization would completely break the function, causing it to no longer approximate the exponential, not even roughly. Similarly, `t` should not be rewritten as `x -. k *. (log2_64h +. log2_64l)`, and `d` should not be moved outside of the parentheses in the computations of `lb` and `ub`. So, not only does floating-point arithmetic ignore the usual algebraic laws of associativity and distributivity, but floating-point experts actively rely on the lack of these laws to compute more accurate approximations.

### 1.3 Challenges of the formal verification

Contrarily to the previous algorithm, the adequacy of this new implementation of the exponential no longer derives from the use of interval arithmetic, so the proof is no longer straightforward. But at the same time, the proof effort needs to be sufficiently light so that it is worth replacing a feature that is already good enough for most use cases.

There have been several attempts at formally verifying this kind of state-of-the-art implementation using the Coq proof assistant, but they all have suffered from various shortcomings. It might have been that the floating-point arithmetic was modeled without any exceptional value [3, §6.2.3]. Indeed, when a computer-assisted proof is meant to complement a pen-and-paper proof, it is acceptable that it only focuses on the most intricate parts of the proof, which the absence of exceptional behavior is hardly ever. But, since this idealized arithmetic does not match the behavior of hardware floating-point numbers, it cannot be used here. Some later attempt solved the issue of the exceptional values [7], but it was still targeting the verification of some code meant to run outside of Coq and thus did not need to cover all of its facets. On the contrary, the algorithm shown in Listing 1 will effectively be run when checking subsequent Coq proofs, so absolutely no shortcuts can be taken.

### 1.4 Contributions

This article proposes a fully proven, fast, and accurate implementation of the exponential function for `CoqInterval`. The intricacy of this implementation corresponds to what is typically found in the state of the art: tables of precomputed values, mixed floating-point integer operations, etc. The proof covers all aspects of the correctness: the argument reduction, the polynomial approximation, and the reconstruction. In addition, this article describes our methodology for formally verifying floating-point approximations of mathematical functions. In particular, we will present the automated strategies that were added to make this verification as painless as possible.

Section 2 reminds both the arithmetic language and the notion of well-behaved expression that were introduced in a previous work [7]. Section 3 explains how some strategies of `CoqInterval` have been improved to automatically verify properties involving tight bounds on rounding errors. Section 4 details the new features added to the arithmetic language and associated tools to tackle the algorithm of Listing 1: hardware floating-point numbers, conversions, macro-operations, array accesses, etc. Section 5 describes the methodology used to formally verify the correctness of exponential, as well as some unusual properties of floating-point arithmetic we ended up with. Section 6 explains how this work relates to some other works. Finally, Section 7 concludes with some benchmarks and some perspectives.

## 2 Preliminaries

This work is partly built on top of a framework for modeling floating-point expressions [7]. In particular, that framework provides some facilities to automate the proof of the absence of exceptional behaviors, thus making it possible for the user to focus on a modeling of floating-point expressions as real numbers. This section reminds the features of that framework that are the most relevant to the presented work. Section 2.1 focuses on the expressions and their various interpretations, while Section 2.2 shows how one can jump between interpretations to ease the proof process.

In the following, the unary operator  $\circ(\cdot)$  designates a rounding operator from  $\mathbb{R}$  to  $\mathbb{R}$ ; it returns the real number the nearest to the input that fits in the target floating-point format (with unlimited range) [3, §3.2.2]. This theoretical operator is at the core of the IEEE-754 standard for floating-point arithmetic.

### 2.1 Arithmetic expressions

An arithmetic expression  $e$  is represented as the value of an inductive type corresponding to a typed abstract syntax tree, namely an expression tree [7]. The nodes of an expression tree correspond to arithmetic expressions, including floating-point operations, integer operations, and some functions such as `nearbyint`.

The expression  $e$  can then be interpreted in several ways, two of which are relevant here. First, it can be interpreted as the floating-point number  $\llbracket e \rrbracket_{\text{flt}}$  that would be obtained according to the IEEE-754 standard. Second,  $e$  can be interpreted as the value  $\llbracket e \rrbracket_{\text{rnd}}$  obtained by performing all the operations on real numbers and rounding their results. For example, in the case of the floating-point addition, we have  $\llbracket u + v \rrbracket_{\text{rnd}} = \circ(\llbracket u \rrbracket_{\text{rnd}} + \llbracket v \rrbracket_{\text{rnd}})$ . In the case of integer operations,  $\llbracket e \rrbracket_{\text{flt}}$  performs computations modulo a power of two, while  $\llbracket e \rrbracket_{\text{rnd}}$  performs operations on unbounded integers, *e.g.*,  $\llbracket u + v \rrbracket_{\text{rnd}} = \llbracket u \rrbracket_{\text{rnd}} + \llbracket v \rrbracket_{\text{rnd}}$ . The first interpretation corresponds to the value actually computed by an implementation, and therefore the value on which we need to prove a correctness theorem. The second interpretation, however, is the one that is the more amenable to formal reasoning, as it is not susceptible to exceptional behaviors such as overflows.

There are two features of expression trees that are of interest to us. The first one is the support for let-binding operators, with binders represented by their de Bruijn indices, to express sharing between sub-expressions and to guide proofs. The second one is the availability of exact arithmetic operations, as they are commonly encountered in implementation of mathematical functions. As far as  $\llbracket \cdot \rrbracket_{\text{flt}}$  is concerned, there is no difference in interpretation between exact and inexact operations over floating-point numbers; they are performed as mandated by the IEEE-754 standard. For  $\llbracket \cdot \rrbracket_{\text{rnd}}$ , exact arithmetic operations, however, are not rounded, which makes formal proofs, both manual and automatic, much simpler. This raises the concern of whether such a proof about  $\llbracket e \rrbracket_{\text{rnd}}$  is meaningful, which Theorem 1 below will tackle.

Let us illustrate these two features on the example of the argument reduction of the Cody-Waite exponential [4], variants of which are still widely used in modern implementations, including in the code shown in Listing 1:

$$\begin{aligned} k &\leftarrow \text{nearbyint}(x \cdot C), \\ t &\leftarrow x - k \cdot c_1 - k \cdot c_2, \end{aligned}$$

with  $c_1 + c_2 \simeq 1/C$  and  $c_2 \ll c_1$ . Below is the formulation of this argument reduction as an expression tree.

```

Let (NearbyInt (Op MUL (Var 0) (BinFl C)))
(Op SUB
  (OpExact SUB (Var 1) (OpExact MUL (Var 0) (BinFl c1)))
  (Op MUL (Var 0) (BinFl c2)))

```

Notice that both floating-point operations in  $x - k \cdot c_1$  are annotated as exact operations by using the `OpExact` constructor. All the other operations are marked as potentially inexact (`Op` constructor). This gives the following value for  $\llbracket t \rrbracket_{\text{rnd}}$ :

$$\circ(x - k \cdot c_1 - \circ(k \cdot c_2)) \quad \text{with } k = \lceil \circ(x \cdot C) \rceil.$$

## 2.2 Relation between interpretations

As mentioned earlier, a correctness statement is about  $\llbracket e \rrbracket_{\text{flt}}$ , while a user only wants to have to deal with  $\llbracket e \rrbracket_{\text{rnd}}$ , as it is free of exceptional behaviors and contains fewer rounding operations. In order to bridge the gap between both interpretations, a predicate `WB` (for *well-behaved*) is defined recursively over expressions. For example, the proposition `WB(Op DIV u v)` is defined as

$$\text{WB}(u) \wedge \text{WB}(v) \wedge \llbracket v \rrbracket_{\text{rnd}} \neq 0 \wedge |\circ(\llbracket u \rrbracket_{\text{rnd}} / \llbracket v \rrbracket_{\text{rnd}})| \leq \Omega$$

with  $\Omega$  the value of the largest finite floating-point number. In other words, for the floating-point division  $u/v$  to be well-behaved, it is sufficient that  $u$  and  $v$  are well-behaved, that the interpretation of  $v$  as a real number is non-zero, and that the division over real numbers, once rounded, does not overflow the floating-point format. The predicate `WB` is defined in a similar way for the other inexact operations over floating-point numbers. For exact operations, the formula contains an additional conjunct that states that the result is exactly representable. For example, the proposition `WB(OpExact ADD u v)` is defined as

$$\text{WB}(u) \wedge \text{WB}(v) \wedge \circ(\llbracket u \rrbracket_{\text{rnd}} + \llbracket v \rrbracket_{\text{rnd}}) = \llbracket u \rrbracket_{\text{rnd}} + \llbracket v \rrbracket_{\text{rnd}} \wedge |\llbracket u \rrbracket_{\text{rnd}} + \llbracket v \rrbracket_{\text{rnd}}| \leq \Omega.$$

The key result is that, if an expression  $e$  is well-behaved, then  $\llbracket e \rrbracket_{\text{flt}}$  is a finite floating-point number and it represents the real number  $\llbracket e \rrbracket_{\text{rnd}}$ . This is expressed by the following theorem:

► **Theorem 1.** *Given an expression  $e$ ,  $\text{WB}(e) \Rightarrow \llbracket e \rrbracket_{\text{flt}} \text{ finite} \wedge \llbracket e \rrbracket_{\text{flt}} = \llbracket e \rrbracket_{\text{rnd}}$ .*

When applying Theorem 1, the user is left with a subgoal `WB(e)`, which is painful to prove by hand. So, to ease the proof process, the framework proposes a proof strategy called `simplify_wb`, which tackles this subgoal by applying a procedure similar to `CoqInterval`'s `interval` strategy to every conjunct of `WB(e)` individually. In practice, one can expect all the conjuncts related to the absence of exceptional behaviors to be automatically proved. Conjuncts related to exact operations, *e.g.*,  $\circ(\llbracket u \rrbracket_{\text{rnd}} + \llbracket v \rrbracket_{\text{rnd}}) = \llbracket u \rrbracket_{\text{rnd}} + \llbracket v \rrbracket_{\text{rnd}}$ , are however out of the scope of `CoqInterval`. So, the user will have to prove them either manually or using a dedicated tool like `Gappa` [3, §4.3].

Note that, in order for `simplify_wb` to make use of the `interval` strategy of `CoqInterval`, the latter had to be enhanced with some support for rounding operators, as they appear in almost all the conjuncts of `WB(e)`. This support was based on the so-called *standard model* of floating-point arithmetic. (Section 3.2 will propose a better approach.) For instance, given an enclosure  $u \in [\underline{u}; \bar{u}]$ , `CoqInterval` would compute an enclosure of  $\circ(u)$  as follows, assuming a *binary64* format:

$$\circ(u) \in [\underline{u} - \varepsilon; \bar{u} + \varepsilon] \quad \text{with } \varepsilon = \max(2^{-1075}, -2^{-53} \cdot \underline{u}, +2^{-53} \cdot \bar{u}). \quad (2)$$

### 3 Automated tools and rounding errors

Consider the code of Listing 1. To verify its correctness, we need to prove the following bound on the absolute error between the exponential function and its floating-point degree-5 polynomial approximation:

$$\forall t \in \mathbb{R}, |t| \leq 355 \cdot 2^{-16} \Rightarrow |1 + y - \exp t| \leq 11 \cdot 2^{-62} \quad \text{with } y = \circ(t \cdot \circ(p_1 + \circ(t \cdot \circ(p_2 + \dots))))). \quad (3)$$

The traditional methodology to prove such a bound is as follows [3, §6.2.3]. One would split the expression  $1 + y - \exp t$  into two parts  $e_1 + e_2$ , with  $e_1 = y - t \cdot (p_1 + t \cdot (p_2 + \dots))$  and  $e_2 = (1 + t \cdot (p_1 + t \cdot (p_2 + \dots))) - \exp t$ . On one hand, expression  $e_1$ , which contains only arithmetic operations and rounding operators, can be bounded using the dedicated Gappa tool [3, §4.3]. On the other hand, expression  $e_2$ , which contains no rounding operator, can be bounded using the rigorous polynomial approximations of CoqInterval [11]. Combining the proofs of both bounds gives the final result.

Since support for rounding operators has been added to CoqInterval so that `simplify_wb` could automatically prove conjuncts of WB [7], it should now be possible to perform this kind of proof directly, without any need for such algebraic manipulations nor the use of an external tool. Unfortunately, several issues arise when using the `interval` strategy on Equation (3). Indeed, it is slightly more involved than the usual conjuncts of WB.

First of all, both sides of the subtraction are strongly correlated, since the left-hand side  $1 + y$  was chosen among the best possible floating-point approximations of the right-hand side  $\exp t$ . This means that naive interval arithmetic, as used in Equation (2) to define the enclosure of a rounding operator, will cause an overestimation of the final enclosure that is so large that it becomes useless for proving anything interesting. On this example, the strategy would only be able to prove that the error is bounded by  $10^{-2}$ , very far from the expected bound of  $11 \cdot 2^{-62}$ . So, the first step is to define rigorous polynomial approximations for rounding operators (Section 3.1).

This is not sufficient though, as the strategy would only succeed in proving a bound of  $24 \cdot 2^{-62}$ , which is already quite good, but not sufficient to prove the correctness of the code of Listing 1. This overestimation is a consequence of using the standard model of floating-point arithmetic to derive Equation (2), as it is a bit too naive. So, the second step is to prove tighter bounds on rounding errors (Section 3.2).

#### 3.1 Rigorous polynomial approximations

The correlation issue of naive interval arithmetic is well-known, and it is independent of rounding errors. In fact, even the interval evaluation of  $(a + x) \cdot (b - x)$  would suffer from it, as  $a + x$  and  $b - x$  vary in opposite directions with respect to  $x$ . A first solution to this issue is to split the domain of  $x$  into smaller sub-intervals and to take the union of the enclosures of the whole expression on all these sub-intervals. This approach is very simple proof-wise, but it scales poorly computation-wise, so it should only be used as a last resort. A second solution is to compute enclosures whose bounds symbolically depend on  $x$  rather than being just numerical values. This approach scales better, but it requires a much larger formalization effort.

CoqInterval provides both approaches [11]. In particular, the second approach is implemented using rigorous polynomial approximations. Instead of just computing a single interval  $[\underline{e}; \bar{e}]$  that encloses an expression  $e(x)$  for any  $x \in X$ , it computes a polynomial  $P$  and an interval  $\Delta$  such that, for any  $x \in X$ , we have  $e(x) - P(x) \in \Delta$ , which we denote by  $e \in (P, \Delta)_X$ . Those polynomial enclosures can then be composed. For example, if we have



$f \in (P_f, \Delta_f)_X$  and  $g \in (P_g, \Delta_g)_X$ , we also have  $f + g \in (P_f + P_g, \Delta_f + \Delta_g)_X$ . This makes it possible to compute the polynomial enclosure of an arbitrary expression, by induction on its structure.

Therefore, to benefit from the rigorous polynomial approximations of CoqInterval, we need to be able to compute a polynomial enclosure of  $\circ(u)$ , given an enclosure  $u \in (P, \Delta)_X$ . To do so, we rewrite  $\circ(u(x))$  into the sum  $[\circ(u(x)) - u(x)] + u(x)$ . For the left-hand side, we use the degree-0 enclosure  $\circ(u) - u \in (0, [-\varepsilon; \varepsilon])_X$  with  $\varepsilon$  computed as in Equation (2). Then, by adding the original enclosure  $(P, \Delta)_X$ , we get a polynomial enclosure of  $\circ(u)$ .

This change to CoqInterval was straightforward, but it has shifted the perspective on rounding operators in the library. Indeed, the original implementation, which was designed for `simplify_wb`, computed an enclosure of  $\circ(u)$  given an enclosure of  $u$ . Then, the user could ask for an enclosure of  $\circ(u) - u$ , which would be correct but overestimated. The new implementation computes a tight enclosure of  $\circ(u) - u$  from an enclosure of  $u$ , from which it derives an enclosure of  $\circ(u)$ . This change has been propagated up to the surface language, that is, CoqInterval now recognizes the expression  $\circ(u) - u$  as an atomic error for an expression  $u$  rather than a subtraction between two sub-expressions involving  $u$ .

### 3.2 Tighter error bounds

By adding support for rounding operators, CoqInterval is now able to automatically prove Equation (3), but only if the rightmost bound is changed to  $24 \cdot 2^{-62}$ . It fails for any tighter bound, especially for  $11 \cdot 2^{-62}$ , which we need to prove the correctness of the implementation of Listing 1. As mentioned earlier, the issue comes from the simplicity of the standard model of floating-point arithmetic, which states that the absolute error between  $\circ(u)$  and  $u$  is bounded by  $2^{-53} \cdot |u|$ , assuming that  $u$  is in the normal range. While this is sensibly true for values of  $u$  slightly larger than a power of two, this is off by a factor two for values of  $u$  that are slightly smaller than a power of two. A better model of floating-point errors is to bound the absolute error between  $\circ(u)$  and  $u$  by  $\frac{1}{2}\text{ulp}(u)$ , where  $\text{ulp}$  denotes the *unit in the last place*, which is the distance between  $|u|$  and its successor. In other words, given an enclosure  $u \in [\underline{u}; \bar{u}]$ , we have the following enclosure of the absolute error:

$$\circ(u) - u \in \left[-\frac{\varepsilon}{2}; \frac{\varepsilon}{2}\right] \quad \text{with } \varepsilon = \text{ulp}(\max(-\underline{u}, \bar{u})).$$

Not only is this new enclosure tighter, but it also makes the implementation and its proof more generic, as it separates the concerns about the target format and the rounding direction. Regarding the target format, one just has to choose the corresponding definition for  $\text{ulp}$ . As a consequence, CoqInterval now supports not only the floating-point formats of Flocq, but also its fixed-point formats. As for the rounding direction, it is a matter of choosing the enclosing interval:  $[-\frac{\varepsilon}{2}; \frac{\varepsilon}{2}]$  for rounding to nearest,  $[0; \varepsilon]$  for rounding toward  $+\infty$ , and so on.

Thanks to these improvements, the `interval` strategy can now directly prove Equation (3). This proof only takes a tenth of a second using degree-10 polynomials (default degree for `interval`). Note that the use of polynomial approximations, rather than the use of more naive variants of interval arithmetic, is critical for this proof, as can be experienced by reducing the degree. With degree 3, it takes about one second; with degree 2, it takes about one minute; and with degree 1, it does not seem to terminate.

## 4 Supported formats and expressions

Since the goal of our work is to formally prove the function shown in Listing 1, we need several new features that were missing from the earlier work on the Cody-Waite algorithm [7]. First of all, since our implementation relies on hardware support for both integer and floating-point



numbers, we need an interpretation of the expressions from Section 2 into the corresponding types (Section 4.1). Since the implementation also uses an array of pre-calculated values to reduce the degree of the polynomial approximation, the grammar of expressions has been extended with array accesses (Section 4.2). Finally, as mentioned earlier, the algorithm takes advantage of the inaccuracies and “flaws” of floating-point arithmetic to implement optimized versions of `nearbyint` and `int_of_float`. Hence, to ease the proof of this algorithm, we have added support for these optimized operations (Section 4.3).

## 4.1 Hardware operations

Similarly to  $\llbracket e \rrbracket_{\text{flt}}$ , we would like to define another interpretation  $\llbracket e \rrbracket_{\text{prim}}$  which represents the computation of  $e$  using the hardware types provided by Coq’s standard library. The `PrimFloat` module offers support for hardware floating-point numbers [12], while the `PrimInt63` module offers support for OCaml’s 63-bit integers [5]. Both modules provide constants, basic operations (+, −, ×, /, etc.), comparisons (=, <, ≤), conversions, and some miscellaneous functions (e.g., floating-point predecessor and successor functions). They also provide axiomatized specifications for these hardware operations.

Since hardware floating-point numbers are just an instance of Flocq’s generic floating-point numbers, we have derived the following variant of Theorem 1:

► **Theorem 2.** *Given an expression  $e$ ,  $\text{WB}(e) \Rightarrow \llbracket e \rrbracket_{\text{prim}} \text{ finite} \wedge \llbracket e \rrbracket_{\text{prim}} = \llbracket e \rrbracket_{\text{rnd}}$ .*

There are two things to note about the definition of  $\llbracket e \rrbracket_{\text{prim}}$ . First, not all operations can be performed directly on hardware types. Fused multiply-add (FMA), for example, is not yet provided by the `PrimFloat` module. Therefore, to complete the definition of  $\llbracket e \rrbracket_{\text{prim}}$ , we emulate these missing operations using the Flocq library (i.e., convert the operands to the formalized Flocq types, compute the result using Flocq’s operations, and convert it back to the hardware type).

Second, we have made our integers 32-bit wide, so that we can use Coq’s 63-bit hardware integers to compute  $\llbracket e \rrbracket_{\text{prim}}$  while maintaining our ability to export verified algorithms as C programs. For 32-bit integer expressions,  $\text{WB}(e)$  hence requires  $\llbracket e \rrbracket_{\text{rnd}}$  to remain inside the  $[-2^{31}; 2^{31} - 1]$  range.

## 4.2 Array accesses

The implementation of Listing 1 starts with an argument reduction that is very similar to Cody & Waite’s, but based on a slightly different identity:

$$\exp(x) = \exp\left(x - k \cdot \frac{\ln 2}{64}\right) \cdot 2^{k/64} \quad \text{with} \quad k = \left\lceil x \cdot \frac{64}{\ln 2} \right\rceil. \quad (4)$$

This makes the reduced argument much smaller, but it also means that the reconstruction is not a simple multiplication by an integer power of 2 anymore. To multiply by  $2^{k/64}$  for some integer  $k$ , we first compute the Euclidean division of  $k$  by 64, in other words find  $k_q$  and  $k_r$  such that  $k = k_q \cdot 64 + k_r$  and  $0 \leq k_r \leq 63$ . Since there are only finitely many different values of  $k_r$ , we pre-compute the floating-point number closest to  $2^{k_r/64}$  for each value of  $k_r$  and store the results in a table `cst`. Therefore, to multiply by  $2^{k/64}$ , we first multiply by `cst.[ $k_r$ ]` and then by  $2^{k_q}$  (see Listing 1).

We have defined `cst` using the Coq standard library `PArray`, which provides persistent arrays [5]. To ease proofs, we have added a constructor `ArrayAcc` to the type of expressions to represent accesses to tables of constants. The constructor takes as argument an array  $a$  of hardware floating-point numbers and an integer expression  $i$  and is interpreted as follows:

$$\llbracket \text{ArrayAcc } a \ i \rrbracket_{\text{prim}} := a.\llbracket i \rrbracket_{\text{prim}}.$$

For an access to be well-behaved, we need the index to be well-behaved and smaller than the length of the array, and all the entries of the array to be finite floating-point numbers.

### 4.3 Macro-operations

As we can see in Equation (4), the argument reduction also requires the `nearbyint` function. This poses a problem as the latter is not provided by the `PrimFloat` module. Since we only need to compute the exponential on inputs in the  $[-746; 710]$  range, we can use the following trick<sup>3</sup> to compute the integer part:

$$\lceil f \rceil = \circ(\circ(f + 1.5 \cdot 2^{52}) - 1.5 \cdot 2^{52}). \quad (5)$$

Using the language of abstract expressions, we could simply represent this sequence of operations as `Op SUB (Op ADD  $f$  (BinF1 0x1.8p52)) (BinF1 0x1.8p52)`. However, that would not be very helpful in proofs because it leaves us the tedious work of showing that those operations behave as `nearbyint`. Instead, we want to treat those operations as if they were one single `nearbyint` operation. For this, we define a new constructor `FastNearbyint` in the language whose interpretation as a floating-point expression is the sequence of operations above, but whose interpretation as a rounded expression is the integer part:

$$\begin{aligned} \llbracket \text{FastNearbyint } e \rrbracket_{\text{flt/prim}} &:= \llbracket e \rrbracket_{\text{flt/prim}} \oplus 0\text{x}1.8\text{p}52 \ominus 0\text{x}1.8\text{p}52, \\ \llbracket \text{FastNearbyint } e \rrbracket_{\text{rnd}} &:= \lceil \llbracket e \rrbracket_{\text{rnd}} \rceil. \end{aligned}$$

Since these interpretations are no longer in one-to-one correspondence, proving Theorems 1 and 2 for these constructors required significantly more work on our part. This, however, saves the user from having to do the work themselves. Note that Equation (5) is only meaningful for inputs  $|f| \leq 2^{51}$ , so  $\text{WB}(\text{FastNearbyint } e)$  contains a conjunct  $|\llbracket e \rrbracket_{\text{rnd}}| \leq 2^{51}$ .

The macro-operation we have just defined computes the integer part as a floating-point number, but the algorithm in Listing 1 also needs it as an integer. Hence, we define another constructor `FastNearbyintToInt` which extracts the mantissa<sup>4</sup> after adding  $1.5 \cdot 2^{52}$ :

$$\lceil f \rceil_{\mathbb{Z}} = \text{mantissa}(\circ(f + 0\text{x}1.8\text{p}52)) - 3 \cdot 2^{51}.$$

## 5 Application: a state-of-the-art exponential

We now have all the ingredients to state and prove the correctness of the algorithm shown in Listing 1. It is stated as follows, with  $x$  the floating-point input, and with  $flb$  and  $fub$  the components of the pair computed by the algorithm:

► **Theorem 3.** *If  $x$  is finite, then  $flb \leq \exp x \leq fub$ .*

Proof of this theorem for large positive and negative values of  $x$  is straightforward, as those are the cases where the exponential either overflows or degenerates to 0. This section presents the methodology we have followed for proving the correctness theorem for  $x \in [-745.13; 709.78]$ .

For a given floating-point input  $x$ , to find an enclosure of  $\exp x$ , our algorithm performs only one approximation  $y$ , but then subtracts (resp. adds) an error term  $d$  to find the lower (resp. upper) bound of the enclosure. Correctness of the algorithm therefore relies on

<sup>3</sup> If  $|f| \leq 2^{51}$  then  $f + 1.5 \cdot 2^{52}$  is between  $2^{52}$  and  $2^{53}$  with  $\text{ulp}(2^{52}) = 1$ , which means the result of the addition is rounded to the nearest integer.

<sup>4</sup> For hardware numbers, we have implemented it as `normfr_mantissa (fst (frshiftexp f))`.

whether  $d$  is big enough to cancel out the inaccuracy of the approximation. For this reason, an essential step in our proof is to find some  $\varepsilon$  such that any choice of  $d$  with  $|d| > \varepsilon$  makes  $\circ(\circ(p_0 \cdot y) + d)$  an upper bound of  $\exp(x - k \ln 2/64) - 2^{k_r/64}$  (and similarly for the lower bound). The value we have experimentally found for  $\varepsilon$  is characterized by the following lemma:

► **Lemma 4.** *For any finite floating-point number  $d$  such that  $|d| \leq 2^{-52}$ , we have*

$$\left| 2^{k_r/64} + \circ(\circ(p_0 \cdot y) + d) - (\exp(x - k_q \ln 2) + d) \right| < \varepsilon \simeq 1.14 \cdot 2^{-57}.$$

The proof of this lemma is too intricate to be comprehensively explained. Instead, we will illustrate our methodology on the parts that verify the argument reduction and the polynomial evaluation (Section 5.1). We will also show part of the proof for the reconstruction as it involves some unusual facts about floating-point arithmetic (Section 5.2).

## 5.1 Illustration of the methodology

Among other facts about the argument reduction, we need to prove that the computation of  $\mathbf{ki}$  causes no exceptional behaviors and that it is indeed an integer part equal to  $k$  despite its convoluted code. To do so in the Coq proof, we have defined an abstract expression  $\mathbf{ki}'$  whose interpretation in the hardware numbers – namely,  $\llbracket \mathbf{ki}' \rrbracket_{\text{prim}}$  – is the value stored in  $\mathbf{ki}$ , and whose interpretation in the rounded real numbers – namely,  $\llbracket \mathbf{ki}' \rrbracket_{\text{rnd}}$  – is  $k$ . Then we can use Theorem 2 to transform a goal about  $\mathbf{ki}$  into a goal about  $k$ .

In practice, we not only want to transform the goal, but we also want to assert some property on the transformed subexpression. Hence, we have implemented a strategy `assert_float` which takes as argument a predicate  $Q$ , looks for an expression of the form  $\llbracket e \rrbracket_{\text{prim}}$ , and applies the following corollary of Theorem 2:

$$\text{WB}(e) \implies Q(\llbracket e \rrbracket_{\text{rnd}}) \implies (\forall x, x = \llbracket e \rrbracket_{\text{rnd}} \wedge Q(x) \implies G(x)) \implies G(\llbracket e \rrbracket_{\text{prim}}).$$

The strategy also invokes `simplify_wb` to discharge as many conjuncts of  $\text{WB}(e)$  as possible. When using this strategy, the Coq proof usually looks as follows:

```
set (ki' := FastNearbyIntToInt (Op MUL (Var 0) InvLog2_64)).
change (normfr_mantissa _ - _)
  with (evalPrim ki' [:x:]). (* Var 0 is mapped to x *)
assert_float (fun ki => -68736 <= ki <= 65536).
{ ... proof of the assertion ... }
```

We have used the `assert_float` strategy 8 times in the proof. Here is another example of its usage to state the main property about the reduced argument  $t$ , which contains exact operations just like in the original Cody-Waite algorithm (see Section 2.1):

```
set (t' := Op SUB (OpExact SUB (Var 1) ...) ...).
change (x - _ - _) with (evalPrim t' [:k, x:]).
assert_float (fun t => abs t <= 355 / 65536
  /\ abs (t - (x - k * ln 2)) <= 65537 * pow2 (-77)).
```

Contrary to the other uses of `assert_float` in the proof, `simplify_wb` is not able to completely discharge the subgoal  $\text{WB}(t)$ . So we have to manually prove the remaining conjuncts, which are the proof obligations of exact operations:

$$\circ(k \cdot c_1) = k \cdot c_1 \quad \wedge \quad \circ(x - k \cdot c_1) = x - k \cdot c_1.$$

The proof of both equalities relies on bit-counting reasoning, which Gappa is specifically designed for [3, §4.3.4]. But to avoid introducing a dependency over Gappa just to prove these two conjuncts, we have performed this reasoning by hand.

As a last illustration, let us consider Equation (3), which bounds the error caused by both the polynomial approximation and its floating-point evaluation. The corresponding proof script looks as follows. For the sake of readability, we have removed a few administrative steps (*e.g.*, unfolding of definitions) from the script.

```
change (Papprox t') with (evalPrim g0 [:t':]).
assert_float (fun y => abs y <= 0.0055
              /\ abs (1 + y - exp t) <= 11 * pow2 (-62)).
{ split.
- interval.
- interval with (i_taylor t, i_bisect t, i_prec 80). }
```

Thanks to the automation provided by `assert_float` and `interval`, in just a few lines, we have proved that none of the floating-point operations had any exceptional behavior, that the image of the floating-point function was bounded, and more importantly, that its error was bounded too. More generally, if a user needed to prove the correctness of a simple floating-point implementation with no intricate argument reduction (*e.g.*, a piece-wise polynomial approximation), that would be the whole of the script.

## 5.2 Correctness of reconstruction

At this point in the proof, thanks to Lemma 4, we know  $2^{k_r/64} + y_\ell \leq \exp x \cdot 2^{-k_q} \leq 2^{k_r/64} + y_u$ , with  $y_\ell$  and  $y_u$  some intermediate floating-point results. To complete the proof, we need to deduce the following enclosure:

$$flb = \text{pred}(\circ(\circ(p_0 + y_\ell) \cdot 2^{k_q})) \leq \exp x \leq \text{succ}(\circ(\circ(p_0 + y_u) \cdot 2^{k_q})) = fub.$$

To do so, we want to factor out the multiplication by  $2^{k_q}$ , but the possibility that the result might fall into the subnormal range makes this factorization impossible. So we have proved the following lemma:

► **Lemma 5.** *Let  $y$  be a binary64 floating-point number greater than  $2^{-1021}$ . Then, for any integer  $k$ ,  $\text{pred}(\circ(y \cdot 2^k)) \leq \text{pred}(y) \cdot 2^k$  and  $\text{succ}(\circ(y \cdot 2^k)) \geq \text{succ}(y) \cdot 2^k$ .*

By transitivity, we are thus left to prove the following inequalities:

$$\text{pred}(\circ(p_0 + y_\ell)) \leq 2^{k_r/64} + y_\ell \quad \wedge \quad 2^{k_r/64} + y_u \leq \text{succ}(\circ(p_0 + y_u)).$$

These inequalities hold because the predecessor and successor functions are enough to compensate both the error between  $p_0$  and  $2^{k_r/64}$  and the rounding error of the final addition. Indeed, there are two cases, depending on the value of  $k_r$ :

- If  $k_r = 0$ , then  $p_0 = 1$ . So, the first inequality reduces to  $\text{pred}(\circ(1 + y_\ell)) \leq 1 + y_\ell$ , which is a general property of the predecessor function. Proof of the upper bound is similar.
- If  $k_r \neq 0$ , then we have  $1.01 < p_0 < 1.99$ . Therefore,  $1.001 < \circ(p_0 + y_\ell) < 1.999$  and hence  $\text{pred}(\circ(p_0 + y_\ell)) = \circ(p_0 + y_\ell) - 2^{-52}$ . Moreover, we have  $|p_0 - 2^{k_r/64}| \leq 2^{-53}$ . Similarly,  $|\circ(p_0 + y_\ell) - (p_0 + y_\ell)| \leq 2^{-53}$ . As a consequence,

$$\begin{aligned} \text{pred}(\circ(p_0 + y_\ell)) &= 2^{k_r/64} + y_\ell + (p_0 - 2^{k_r/64}) + (\circ(p_0 + y_\ell) - (p_0 + y_\ell)) - 2^{-52} \\ &\leq 2^{k_r/64} + y_\ell + 2^{-53} + 2^{-53} - 2^{-52} = 2^{k_r/64} + y_\ell \end{aligned}$$

which completes the proof for the lower bound. Proof of the upper bound is similar.

## 6 Related work

While there had been some earlier works to formalize hardware arithmetic operators [18], formal verification of mathematical libraries really started with the impressive work by John Harrison. Among other things, he used the HOL Light system to prove the correctness of a *binary32* approximation of the exponential function which presents many similarities with our own algorithm [8]. But being a *binary32* function, it could nowadays be validated by sheer exhaustive testing. So, perhaps more interesting is Harrison’s subsequent work on the formal verification of the implementation of sin and cos for IA-64 architectures, as it sets the bar even higher [9]. Indeed, these approximations perform 80-bit floating-point computations and are accurate to 0.574 ulp for inputs smaller than  $2^{63}$ . They use an intricate argument reduction: first a pre-reduction, followed by a 3-term Cody-Waite reduction, resulting in a double-*binary80* reduced argument. Then a degree-17 polynomial is evaluated, followed by a simple reconstruction of the result so as to take the lower part of the reduced argument into account. During both the argument reduction and the reconstruction, several floating-point operations are actually exact and need to be considered as such, in order to be able to prove anything interesting about the result. Our methodology could be used to automate various parts of this proof, but the representation of the reduced argument as a non-evaluated sum of two floating-point numbers would presumably warrant adding a few more macro-operations to our expression language, *e.g.*, a FastTwoSum operator [15, §1.3].

A more recent work is the large verification using the Coq proof assistant of the power function of the CORE-MATH library by Laurence Rideau and Laurent Théry [10]. This includes the correctness of an exponential function whose implementation shares some similarities with ours, but it is a lot more subtle, since both input and output are double-*binary64* numbers. Their formalization, however, ignores the issue of exceptional behaviors and just assumes that numbers can be arbitrarily large, as is traditionally the case in pen-and-paper proofs. Again, our methodology could help transition to a complete proof, especially since they are already making heavy use of CoqInterval.

Regarding the use of hardware floating-point numbers in the Coq proof assistant, Érik Martin-Dorel and Pierre Roux have implemented and verified a checker for semi-definite positive matrices [12]. The algorithm performs a Choleski decomposition using floating-point arithmetic on a slightly perturbed input matrix. The correctness theorem states that, if this decomposition succeeds, then a Choleski decomposition using exact arithmetic would have succeeded on the original input matrix, which guarantees that it was indeed semi-definite positive. The perturbation, and hence the correctness proof, depends on the ability to compute a bound on the rounding error of the floating-point decomposition [16]. Our approach would have been of little help for that use case, as the algorithm and the error bound highly depend on the dimension of the matrix.

Regarding the automation of proofs of bounds on rounding errors in a proof assistant, one can cite the FPTaylor tool, which can generate proofs for HOL Light [19]. Given a floating-point expression, it computes an affine form that encloses it, using elementary rounding errors as variables of the affine form. The strength of that tool is that the coefficients of the affine form are kept as symbolic expressions rather than intervals. This approach separates the concerns between the global optimizer used for computing enclosures and the formalization of affine forms for floating-point arithmetic. Indeed, enclosures of the symbolic expressions are only needed when they occur in terms of order 2 or more, as these terms cannot be represented as part of the affine form. Therefore, the verification of these enclosures can be done in a rather naive way, since they are only used for higher-order error term and thus do

not have to be tight. It should be noted that FPTaylor supports both the standard model of floating-point arithmetic and a tighter model (see Section 3.2), but it can only generate proofs for the former. Moreover, the global optimizer used to compute the enclosures in FPTaylor is not the same as the one used to verify them in HOL Light, which might cause difficulties if the latter procedure is not strong enough or too slow.

A similar tool is PRECiSA, which targets the PVS proof assistant [13]. As with FPTaylor, errors are kept as symbolic expressions, and a global optimizer is used to compute their enclosures. Higher-order error terms, however, are not eliminated as computations progress, which might cause some performance issues compared to FPTaylor. PRECiSA, however, uses a tight formal model of floating-point errors, and the tool can detect exact subtractions (Sterbenz' lemma). Moreover, it supports conditional expressions, including the cases where rounding causes a different branch to be taken.

Finally, one should mention the VCFloater tool, which targets the Coq proof assistant [1]. As with the previous two tools, the error is kept as a symbolic expression. Before being fed to a global optimizer (namely CoqInterval), this expression is first simplified by expanding it through distributivity and discarding the sub-expressions that cancel. This expansion might cause some performance issues, due to combinatorial explosion. A user-provided threshold is used to further discard negligible terms, at the expense of a potentially worse error bound. It can also use a technique similar to Gappa to reduce the correlation between sub-expressions, and thus improve the tightness of the computed enclosures. The tool uses the standard model of floating-point errors, but the user can annotate operations that are supposed to be exact and the tool will verify that the conditions hold (Sterbenz' lemma). Moreover, the tool supports user-defined operations, which means that it can easily be extended with double-word arithmetic, as long as the user has formalized it beforehand.

## 7 Conclusion

In this article, we have presented a floating-point approximation of the exponential function, its mechanized proof of correctness, and the tools we have developed to ease the verification work. One peculiarity of this work is that the verified code is not just modeled using the Coq proof assistant, it can actually run in the logic of the system and therefore be used to perform proofs by computations. Indeed, the correctness theorem tells how the result of the approximation can be used as a lower/upper bound of the mathematical exponential. The specification of the code of Listing 1 and its correctness proof take about 600 lines of Coq script;<sup>5</sup> Lemma 5 is about 130 lines; extending the proof of Theorem 1 to support macro-operations and arrays takes about 500 lines; the tighter bounds on rounding errors take about 200 lines. This work was integrated in release 4.10.0 of CoqInterval.

### 7.1 Integration to CoqInterval and performances

As explained in the introduction, the CoqInterval library provides an interval extension of the exponential function that can use the floating-point unit of the processor to speed up proof checking [12]. Its implementation, however, is based on a truncated power series, which is effective but rather naive, compared with the implementations that can be found in mathematical libraries targeting hardware floating-point formats. We have thus plugged our verified implementation in place of the original one. Consider the following Coq script.

---

<sup>5</sup> [https://gitlab.inria.fr/coqinterval/interval/-/blob/interval-4.11.0/src/Interval/Float\\_full\\_primfloat.v](https://gitlab.inria.fr/coqinterval/interval/-/blob/interval-4.11.0/src/Interval/Float_full_primfloat.v)

```
Goal forall x, 10 <= x <= 11 -> Rabs (exp x - exp x) <= 0x1p-6.
Proof. intros x Hx. interval with (i_bisect x, i_depth 30). Qed.
```

It states that, for any real number  $x$  between 10 and 11, the difference between  $\exp x$  (mathematical exponential) and itself is less than  $2^{-6}$ . From a mathematical point of view, this statement is useless, since it could be trivially proved by rewriting  $\exp x - \exp x$  to zero, but it is a good way to exercise the computations performed by CoqInterval. Indeed, the way the `interval` strategy is invoked, it will not try to use anything fancier than naive interval arithmetic. As a consequence, because  $\exp x$  is strongly correlated with itself, the formal proof generated by the tactic ends up considering around 6.7 million sub-intervals of the input enclosure  $x \in [10; 11]$  (and as many interval evaluations of the exponential function).

Using the original implementation, this computationally intensive proof takes about 160 seconds to be checked by the Coq proof assistant on an Intel 13th-generation 4GHz processor. With the implementation verified in this work, the proof is checked in less than 8 seconds. Taking the average of three runs, the speedup is  $20.5\times$ . Since the argument reduction of the original implementation is more costly the further away from zero the input is, the speedup can grow even larger, up to  $24\times$ .

As for the accuracy of the new implementation, one can get an intuitive feel of it by considering the distance between the bounds of the output interval. Ideally, it should be one ulp (except for the input 0), meaning that the bounds of the interval should be consecutive floating-point numbers. This property, called *correct rounding* [14, §12.3], is still an open research question for floating-point formats larger than *binary32* and completely out of reach of a formal proof, as of today. So, the best we can hope to achieve is a distance of up to two ulps, that is, one component is optimal, while the other is off-by-one. In the code shown in Listing 1, if the constant  $d$  was zero, this would be the case. As it is not quite zero here, when  $\exp x$  is close to the midpoint between two consecutive floating-point numbers, the distance might end up being three ulps. The proportion of inputs that cause a 3-ulp interval output is roughly  $d \cdot 2^{51} \simeq 1/60$ .

## 7.2 Real-life performances

Being able to perform about one million faithful interval evaluations of exponential per second inside the logic of Coq is impressive, but it is nowhere near the actual throughput of the floating-point unit of the processor. Indeed, disregarding any concern about the guaranteed accuracy of a mathematical library, one should expect a state-of-the-art implementation to take 25–50 cycles to compute two floating-point approximations of exponential<sup>6</sup> (and thus one interval enclosure), so about  $100\times$  faster than what we currently achieve in the logic of Coq. There are several reasons for the remaining gap. First of all, the code of our implementation is not directly run by the processor, but interpreted by a virtual machine. Second, this bytecode interpreter boxes floating-point numbers, and thus performs a large amount of memory allocations. Third, while our code only performs computations on values, the interpreter still needs to account for the possibility of open terms (*e.g.*, free variables) appearing as operands to the floating-point computations.

The first issue can be worked around by using the `native_compute` machinery of the Coq system, which compiles the code using the OCaml compiler and then executes it directly [2]. This machinery also partially avoids the second issue, since the compiler can optimize away

<sup>6</sup> <https://core-math.gitlabpages.inria.fr/>



■ **Listing 2** Floating-point exponential in OCaml.

```
let fexp x =
  if x < -0x1.74385446d71c4p9 then 0. else
  if x > 0x1.62e42fefa39efp9 then infinity else
  if x <> x then nan else
  let k' = x *. invLog2_64 +. 0x1.8p52 in
  ...
  let p0 = cst.(ki land 63) in
  ldexp (p0 +. p0 *. y) (ki asr 6)
```

the boxing of some intermediate floating-point results. But the third issue is still present and makes it hard to avoid pessimization in the generated code. As a consequence, this only improves proof checking by a factor  $3\times$  to  $4\times$  for the longer proofs.

To get a better feel of the actual performances of our implementation, we can instead implement the function directly in OCaml, as shown in Listing 2. This is roughly the same code as Listing 1, except that the original last three lines, which were computing an enclosure of  $\exp x$ , have been replaced by a single floating-point value: `ldexp (p0 +. p0 *. y) (ki asr 6)`. Accordingly, the first few lines return a single value for the exceptional cases. The code is run on about  $1.5 \cdot 10^9$  inputs uniformly distributed among those that lead to a finite output. Compiling the code with OCaml 5.1.1, we get that the floating-point exponential from the GNU C Library is about  $1.45\times$  faster than our implementation.

Even if the GNU C Library has been heavily tuned, this is still a rather large gap. Part of the reason is its use of the FMA operation. This ternary operation computes  $\circ(x \cdot y + z)$  at once, which halves the number of operations performed during the argument reduction and the polynomial evaluation. Modifying our code accordingly, this reduces its slowdown to  $1.32\times$ . When translating the code to C and compiling it with GCC, the slowdown is brought down to  $1.24\times$ . Obviously, using FMAs in place of multiplications and additions invalidates the correctness proof, since they do not compute the same values (notice the lack of rounding operator around the product). Fortunately, the proof can be easily adapted. Indeed, exact operations during the argument reduction are still exact when performed with an FMA, and having a more accurate polynomial evaluation only makes the proof simpler. Note that, while our framework supports reasoning about the FMA operation, it is not one of the native floating-point operations provided by the Coq system, so it cannot be used to speed up the implementation of CoqInterval. One would instead have to use larger tables, as does the GNU C Library, so as to reduce the degree of the polynomial approximation.

### 7.3 Future works

First, it should be noted that, while the GNU C Library does not implement correct rounding either, it is nonetheless slightly more accurate than our implementation. In about 20% of cases, the code of Listing 2 returns a floating-point result that is off by one, while for the GNU C Library, the probability is  $10^{-5}$ . In the context of CoqInterval, this hardly matters, since we want to compute an enclosure of the mathematical result rather than the nearest floating-point number. But for a mathematical library, people might prefer a code that is experimentally a bit more accurate to a code whose correctness has been formally verified. Most of the inaccuracy comes from the factor `p0`. There are two ways to improve it, both of which require adding a new table along `cst`. In the first approach, the new table contains the error on `p0`, which can then be reintroduced in the computation. In the second approach, the new table tells how to shift the input, such that the error on `p0` becomes negligible.

A natural extension of this work is to convert all the other mathematical functions of CoqInterval to use some state-of-the-art implementation when hardware floating-point numbers are used as interval bounds. For functions such as log and arctan, our approach should work without difficulty, as they are quite similar to exp. For trigonometric functions such as sin and cos, the situation is slightly different. First of all, they are not monotone, so considering the lower and upper bounds of the input interval separately might be counter-productive; it might be better to perform a simultaneous argument reduction on both bounds. Second, the Cody-Waite approach to argument reduction does not scale well to extremely large inputs, while some other algorithms for argument reduction take advantage of the periodicity of the trigonometric functions [14, §11.4].

---

## References

---

- 1 Andrew Appel and Ariel Kellison. VCFloat2: Floating-point error analysis in Coq. In *13th ACM SIGPLAN International Conference on Certified Programs and Proofs*, pages 14–29, London, United Kingdom, 2024. doi:10.1145/3636501.3636953.
- 2 Mathieu Boespflug, Maxime Dénès, and Benjamin Grégoire. Full reduction at full throttle. In *1st International Conference on Certified Programs and Proofs*, pages 362–377, Kenting, Taiwan, December 2011. doi:10.1007/978-3-642-25379-9\_26.
- 3 Sylvie Boldo and Guillaume Melquiond. *Computer Arithmetic and Formal Proofs*. ISTE Press – Elsevier, 2017.
- 4 William J. Cody, Jr. and William Waite. *Software Manual for the Elementary Functions*. Prentice-Hall, Englewood Cliffs, NJ, 1980.
- 5 Maxime Dénès. Towards primitive data types for Coq 63-bits integers and persistent arrays. In *5th Coq Workshop*, Rennes, France, July 2013. URL: [https://coq.inria.fr/files/coq5\\_submission\\_2.pdf](https://coq.inria.fr/files/coq5_submission_2.pdf).
- 6 Paul Geneau de Lamarlière and Guillaume Melquiond. CoqInterval. Software, version 4.10.0., swhId: swh:1:dir:78da3e6e98b7ef018180119255ce1e10a048cc88 (visited on 2024-08-22). URL: <https://gitlab.inria.fr/coqinterval/interval.git>.
- 7 Paul Geneau de Lamarlière, Guillaume Melquiond, and Florian Faissole. Slimmer formal proofs for mathematical libraries. In Theo Drane and Anastasia Volkova, editors, *30th IEEE International Symposium on Computer Arithmetic*, Portland, OR, USA, September 2023. doi:10.1109/ARITH58626.2023.00026.
- 8 John Harrison. Floating-point verification in HOL Light: The exponential function. Technical Report 428, University of Cambridge Computer Laboratory, 1997.
- 9 John Harrison. Formal verification of floating point trigonometric functions. In Warren A. Hunt and Steven D. Johnson, editors, *3rd International Conference on Formal Methods in Computer-Aided Design*, volume 1954 of *Lecture Notes in Computer Science*, pages 217–233, 2000.
- 10 Tom Hübner, Claude-Pierre Jeannerod, Paul Zimmermann, Laurence Rideau, and Laurent Théry. Towards a correctly-rounded and fast power function in binary64 arithmetic, 2024. URL: <https://inria.hal.science/hal-04159652>.
- 11 Érik Martin-Dorel and Guillaume Melquiond. Proving tight bounds on univariate expressions with elementary functions in Coq. *Journal of Automated Reasoning*, 57(3):187–217, 2016. doi:10.1007/s10817-015-9350-4.
- 12 Érik Martin-Dorel, Guillaume Melquiond, and Pierre Roux. Enabling floating-point arithmetic in the Coq proof assistant. *Journal of Automated Reasoning*, 67, 2023. doi:10.1007/s10817-023-09679-x.
- 13 Mariano Moscato, Laura Titolo, Aaron Dutle, and César A. Muñoz. Automatic estimation of verified floating-point round-off errors via static analysis. In Stefano Tonetta, Erwin Schoitsch, and Friedemann Bitsch, editors, *36th International Conference on Computer Safety, Reliability*,

- and Security*, volume 10488 of *Lecture Notes in Computer Science*, pages 213–229, Trento, Italy, 2017. doi:10.1007/978-3-319-66266-4\_14.
- 14 Jean-Michel Muller. *Elementary Functions, Algorithms and Implementation*. Birkhäuser, Boston, MA, 3rd edition, 2016. doi:10.1007/978-1-4899-7983-4.
  - 15 Jean-Michel Muller, Nicolas Brunie, Florent de Dinechin, Claude-Pierre Jeannerod, Mioara Joldes, Vincent Lefèvre, Guillaume Melquiond, Nathalie Revol, and Serge Torres. *Handbook of Floating-Point Arithmetic*. Birkhäuser Basel, 2nd edition, 2018. doi:10.1007/978-3-319-76526-6.
  - 16 Pierre Roux. Formal proofs of rounding error bounds – with application to an automatic positive definiteness check. *Journal of Automated Reasoning*, 57(2):135–156, 2016. doi:10.1007/s10817-015-9339-z.
  - 17 Siegfried M. Rump. Verification methods: Rigorous results using floating-point arithmetic. *Acta Numerica*, 19:287–449, May 2010. doi:10.1017/S096249291000005X.
  - 18 David M. Russinoff. *Formal Verification of Floating-Point Hardware Design*. Springer Cham, 2nd edition, 2022. doi:10.1007/978-3-030-87181-9.
  - 19 Alexey Solovyev, Marek S. Baranowski, Ian Briggs, Charles Jacobsen, Zvonimir Rakamarić, and Ganesh Gopalakrishnan. Rigorous estimation of floating-point round-off errors with symbolic Taylor expansions. *ACM Transactions on Programming Languages and Systems*, 41(1), December 2018. doi:10.1145/3230733.