

# 1 Primitive Floats in Coq

2 **Guillaume Bertholon** 

3 École Normale Supérieure, Paris, [École Normale Supérieure, 45 rue d’Ulm, 75230 Paris Cedex 05],  
4 France

5 guillaume.bertholon@ens.fr

6 **Érik Martin-Dorel** 

7 Lab. IRIT, University of Toulouse, CNRS, [IRIT Université Paul Sabatier, 118 route de Narbonne,  
8 31062 Toulouse Cedex 9], France

9 <https://www.irit.fr/~Erik.Martin-Dorel/>

10 erik.martin-dorel@irit.fr

11 **Pierre Roux** 

12 ONERA, Toulouse, [ONERA, Centre de Toulouse, 2 Avenue Édouard Belin, 31055 Toulouse  
13 Cedex 4], France

14 <https://www.onera.fr/staff/pierre-roux>

15 pierre.roux@onera.fr

## 16 — Abstract —

17 Some mathematical proofs involve intensive computations, for instance: the four-color theorem, Hales’  
18 theorem on sphere packing (formerly known as the Kepler conjecture) or interval arithmetic. For  
19 numerical computations, floating-point arithmetic enjoys widespread usage thanks to its efficiency,  
20 despite the introduction of rounding errors.

21 Formal guarantees can be obtained on floating-point algorithms based on the IEEE 754 standard,  
22 which precisely specifies floating-point arithmetic and its rounding modes, and a proof assistant  
23 such as Coq, that enjoys efficient computation capabilities. Coq offers machine integers, however  
24 floating-point arithmetic still needed to be emulated using these integers.

25 A modified version of Coq is presented that enables using the machine floating-point operators.  
26 The main obstacles to such an implementation and its soundness are discussed. Benchmarks show  
27 potential performance gains of two orders of magnitude.

28 **2012 ACM Subject Classification** Theory of computation → Type theory; Mathematics of computing  
29 → Numerical analysis; General and reference → Performance

30 **Keywords and phrases** Coq formal proofs, floating-point arithmetic, reflexive tactics, Cholesky  
31 decomposition

32 **Digital Object Identifier** 10.4230/LIPIcs.ITP.2019.33

33 **Supplement Material** <https://github.com/coq/coq/pull/9867>

34 **Acknowledgements** The authors would like to thank Maxime Dénès and Guillaume Melquiond for  
35 helpful discussions.

## 36 **1** Motivation

37 The proof of some mathematical facts can involve a numerical computation in such a way  
38 that trusting the proof requires trusting the numerical computation itself. Thus, being able  
39 to efficiently perform this kind of proofs inside a proof assistant eventually means that the  
40 tool must offer efficient numerical computation capabilities.

41 Floating-point arithmetic is widely used in particular for its efficiency thanks to its  
42 hardware implementation. Although it does not generally give exact results, introducing  
43 rounding errors, rigorous proofs can still be obtained by bounding the accumulated errors.



© Guillaume Bertholon and Érik Martin-Dorel and Pierre Roux;  
licensed under Creative Commons License CC-BY

10th International Conference on Interactive Theorem Proving (ITP 2019).

Editors: John Harrison, John O’Leary, and Andrew Tolmach; Article No. 33; pp. 33:1–33:20

Leibniz International Proceedings in Informatics



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

```

R := 0;
for j from 1 to n do
  for i from 1 to j - 1 do
    Ri,j := (Ai,j - Σk=1i-1 Rk,i Rk,j) / Ri,i;
  end for
  Rj,j := √(Mj,j - Σk=1j-1 Rk,j2);
end for

```

■ **Figure 1** Cholesky decomposition: given  $A \in \mathbb{R}^{n \times n}$ , attempts to compute  $R$  such that  $A = R^T R$ .

44 There is thus a clear interest in providing an efficient and sound access to the processor  
 45 floating-point operators inside a proof assistant such as Coq.

## 46 1.1 Proofs Involving Numerical Computations

47 We give below a few examples of proofs involving floating-point computations.

48 As a first example, consider the proof that a given real number  $a \in \mathbb{R}$  is nonnegative.  
 49 One can exhibit another real number  $r$  such that  $a = r^2$  and apply a lemma stating that all  
 50 squares of real numbers are nonnegative. Typically, one could use the square root  $\sqrt{a}$ .

51 A similar method can be applied to prove that a matrix  $A \in \mathbb{R}^{n \times n}$  is positive semidefinite<sup>1</sup>  
 52 as one can exhibit  $R$  such that<sup>2</sup>  $A = R^T R$ . Such a matrix can be computed using an algorithm  
 53 called Cholesky decomposition, given in Figure 1. The algorithm succeeds, taking neither  
 54 square roots of negative numbers nor divisions by zero, whenever  $A$  is positive definite<sup>3</sup>.

55 When executed with floating-point arithmetic, the exact equality  $A = R^T R$  is lost but it  
 56 remains possible to bound the accumulated rounding errors in the Cholesky decomposition  
 57 such that the following theorem holds under mild conditions.

58 ► **Theorem 1** (Corollary 2.4 in [34]). *For  $A \in \mathbb{R}^{n \times n}$ , defining  $c := \frac{(n+1)\epsilon}{1-2(n+1)\epsilon} \text{tr}(A) +$   
 59  $4n(2(n+1) + \max_i A_{i,i})\eta$ , if the floating-point Cholesky decomposition succeeds on  $A - cI$ ,  
 60 then  $A$  is positive definite.  $\epsilon$  and  $\eta$  are tiny constants given by the floating-point format used.*

61 A formal proof in Coq of this theorem can be found in a previous work [33]. Thus,  
 62 an efficient implementation of floating-point arithmetic inside the proof assistant leads to  
 63 efficient proofs of matrix positive definiteness. This can have multiple applications, such as  
 64 proving that polynomials are nonnegative by expressing them as sums of squares [26] which  
 65 can be used in a proof of the Kepler conjecture [24].

66 Interval arithmetic constitutes another example of proofs involving numerical computa-  
 67 tions. Sound enclosing intervals can be easily computed in floating-point arithmetic using  
 68 directed roundings, towards  $\pm\infty$  for lower or upper bounds. The Coq.Interval library [25]  
 69 implements interval arithmetic and could benefit from efficient floating-point arithmetic.

70 More generally, there are many results on rigorous numerical methods [35] that could  
 71 see efficient formal implementations provided efficient floating-point arithmetic is available  
 72 inside proof assistants.

<sup>1</sup> A matrix  $A \in \mathbb{R}^{n \times n}$  is said positive semidefinite when for all  $x \in \mathbb{R}^n$ ,  $x^T A x \geq 0$ .

<sup>2</sup> Since, when  $A = R^T R$ , one gets  $x^T A x = x^T (R^T R) x = (Rx)^T (Rx) = \|Rx\|^2 \geq 0$ .

<sup>3</sup> A matrix  $A \in \mathbb{R}^{n \times n}$  is said positive definite when for all  $x \in \mathbb{R}^n \setminus \{0\}$ ,  $x^T A x > 0$ .

## 1.2 Objectives

The Coq proof assistant has built-in support for computation, which can be used within proofs, and recent progress have been done to provide efficient integer computation (relying on 63-bit machine integers).

The overall goal of this work is to implement efficient floating-point computation in Coq, relying directly on machine `binary64` floats, instead of emulating floats with pairs of integers. Experimentally, that latter emulation in Coq incurs a slowdown of about three orders of magnitude with respect to an equivalent implementation written in OCaml.

## 1.3 Outline

The article is organized as follows: Section 2 provides the background required to position our approach, from proof-by-reflection to the IEEE 754 standard for floating-point arithmetic to interval arithmetic formalized in Coq. Section 3 is devoted to the implementation itself, with a special focus on the interface that it exposes. Section 4 gathers a discussion on several design choices or technicalities that have been important to carry out the implementation and avoid some pitfalls. Section 5 provides benchmarks to evaluate the performance of the implementation. Section 6 finally gives concluding remarks and perspectives for future work.

## 2 Prerequisites and Related Works

In this section, we start by reviewing the two main features that underlie and motivate our work in the Coq proof assistant: Poincaré’s principle and the availability of efficient reduction tactics (in Section 2.1). We then give an overview of all notions of floating-point arithmetic that appear necessary to make this paper self-contained (in Section 2.2). We finally summarize the features of two related Coq libraries that are either a prerequisite for our developments (in Section 2.3), or an important building block for a possible extension of this work (in Section 2.4).

### 2.1 Proof by Reflection and Efficient Numerical Computation

In the family of formal proof assistants, the underlying logic of several systems—including Agda, Coq, Lego, and Nuprl [2]—provides a notion of definitional equality that allows one to automatically prove some equalities by a mere computation. This feature is called *Poincaré’s principle* in reference to Poincaré’s statement that “a reasoning proving that  $2 + 2 = 4$  is not a proof in the strict sense, it is a verification” [32, chap. I]. Based upon this principle, the so-called *proof by reflection* methodology has been developed to take advantage of the computational capabilities of the provers and build efficient (semi)-decision procedures [7]: this approach has been successfully applied to various application domains, such as: graph theory, with the formal verification of the four-color theorem in Coq by Gonthier and Werner [14], discrete geometry, with the formal proof of the Kepler conjecture developed in the Flyspeck project [17], Boolean satisfiability, with the verification of SAT traces in Coq [1], satisfiability modulo theories, with the development of the SMTCoq library [13], or global optimization, with the development of the ValidSDP library [26].

To be able to address the verification of increasingly complex proofs relying on this approach, works have been carried out to increase the computational performance of proof assistants, relying on two complementary approaches: (i) implement alternative evaluation engines, such as evaluators based on compilation to bytecode or native code, and (ii) optimized data structures that might be based on machine values and hardware operators.

116 For example, the Isabelle proof assistant provides (i) several evaluators that can be used  
 117 within proofs, and allows one to generate Standard ML, OCaml, Haskell, or Scala code, then  
 118 (ii) libraries of fast machine words (for fixed size or unspecified size) have been developed  
 119 while ensuring compatibility with all Isabelle’s target languages and evaluators [23].

120 In this work, we specifically focus on the Coq proof assistant which offers in particular  
 121 (i) the reduction tactics `vm_compute`, involving bytecode compilation and evaluation by a  
 122 virtual machine [15] and `native_compute`, involving code generation and native OCaml  
 123 compilation [3], as well as (ii) *machine integers*, upon which the `Bignums` library for multiple-  
 124 precision arithmetic has been developed [16].

125 Regarding machine integers in Coq, the original implementation by Spiwack [1, 39] was  
 126 based on the so-called *retro-knowledge* approach, which consisted in developing a reference  
 127 implementation of 31-bit integer operators in Coq (using lists of bits), then optimizing their  
 128 evaluation in `vm_compute` (and later `native_compute`) by replacing the considered Coq  
 129 operator on-the-fly with the corresponding hardware operator. The implicit assumption here  
 130 is that both implementations match. This implementation has been recently replaced with  
 131 so-called *primitive integers*<sup>4</sup> [12]: this approach required adding a representation of 63-bit  
 132 machine integers in the kernel, and has the two-fold benefit of offering efficient operators for  
 133 all reduction strategies with a compact representation of integers, and making explicit the  
 134 axioms that specify the primitive operators.

135 The overall aim of this work is to provide a similar facility for floating-point arithmetic,  
 136 to be able to compute with *primitive floating-point numbers* in Coq, instead of emulating  
 137 floating-point numbers with pairs of integers.

138 A facility to compute with floating-point numbers for prototyping purposes is available in  
 139 the PVS proof assistant thanks to the PVSio package [31] but to the best of our knowledge,  
 140 no proof assistant currently provides support for machine floating-point computations in the  
 141 scope of proof by reflection.

## 142 2.2 Floating-point Arithmetic

143 This section reviews the main concepts of floating-point arithmetic used in the remainder of  
 144 this paper. The reader interested in more details could find them in reference books [30].

145 Computing in floating-point arithmetic amounts to performing calculations in what is  
 146 often called scientific notation with one digit before the dot, a fixed number of digits following  
 147 it and a power of ten specifying the position of the dot, hence the name *floating-point*  
 148 arithmetic. When results do not fit in the required precision, they have to be rounded, e.g.,  
 149 with a precision of five digits,  $1.234 \cdot 10^2 + 5.678 \cdot 10^{-1} = 1.240 \cdot 10^2$ .

### 150 2.2.1 IEEE 754 Standard

151 Implementations of floating-point arithmetic in hardware nowadays adhere to the IEEE 754  
 152 standard [19]. This standard prescribes sets of floating-point numbers, mostly as subsets  
 153 of the real numbers field  $\mathbb{R}$ , binary representations for them, rounding modes and basic  
 154 arithmetic operators  $+$ ,  $-$ ,  $\times$ ,  $\div$  and  $\sqrt{\cdot}$  defined as functions giving the same result as the  
 155 operator in the real field composed with a rounding.

156 A floating-point format  $\mathbb{F}$  is a subset of  $\mathbb{R}$  such that  $x \in \mathbb{F}$  when

$$157 \quad x = m\beta^e \tag{1}$$

---

<sup>4</sup> See the pull request <https://github.com/coq/coq/pull/6914>.



195 The IEEE 754-2008 standard [19] defines five standard rounding modes:

196 **toward  $-\infty$** :  $\text{RD}(x)$  is the largest floating-point number  $\leq x$ ;

197 **toward  $+\infty$** :  $\text{RU}(x)$  is the smallest floating-point number  $\geq x$ ;

198 **toward zero**:  $\text{RZ}(x)$  is equal to  $\text{RD}(x)$  if  $x \geq 0$ , and to  $\text{RU}(x)$  if  $x \leq 0$ ;

199 **to nearest even**:  $\text{RNE}(x)$  is the floating-point number closest to  $x$ .

200 In case of a tie: the one with an even mantissa;

201 **to nearest away from zero**:  $\text{RNA}(x)$  is the floating-point number closest to  $x$ .

202 In case of a tie: the one with the largest mantissa in absolute value.

203 In this work, we will only rely on the RNE rounding, which is the default rounding mode  
204 in most floating-point programming environments. See Section 4.1 for a more in depth  
205 discussion of this point.

206 Then, all floating-point operators are required to be correctly rounded, that is to say, they  
207 should behave as if they were computed with an infinitely precise mantissa, then rounded  
208 according to the specified rounding mode. To be more precise, for a given floating-point  
209 format  $\mathbb{F}$ , operator  $*$  :  $\mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ , and rounding mode  $\circ$  :  $\mathbb{R} \rightarrow \mathbb{F}$ , a correctly-rounded  
210 implementation  $\otimes$  of  $*$  should verify:

$$211 \quad \forall x, y \in \mathbb{F}, \quad x \otimes y = \circ(x * y).$$

212 The benefits of this definition are two-fold:

- 213 ■ all floating-point operators that are correctly-rounded (the 2008 revision of the standard  
214 requiring this for  $+$ ,  $-$ ,  $\times$ ,  $\div$ ,  $\sqrt{\cdot}$ ) are fully-specified, which straightforwardly ensures the  
215 reproducibility of the results;
- 216 ■ it allows one to devise floating-point algorithms that directly rely upon this specification,  
217 as exemplified in the upcoming Section 2.2.2.

## 218 2.2.2 Error Free Transformations

219 Noticing that the rounding error of a floating-point addition is itself a floating-point number,  
220 algorithms such as Fast2Sum [11] and 2Sum [21, 28] can compute that exact error, taking  
221 advantage of correct rounding.

222 These two “compensated summation algorithms” fall into the larger class of error-free  
223 transformations [22, 37] which constitute an essential building block in the development of  
224 extended precision floating-point algorithms.

## 225 2.2.3 Standard Model

226 Although precise specifications are known for roundings, hence for basic arithmetic operators,  
227 a simpler model is commonly used to prove compound bounds of rounding errors on larger  
228 expressions [18]. Despite being weaker, this model is more amenable to algebraic proofs,  
229 whether pen and paper or mechanized. Called standard model of floating-point arithmetic,  
230 it states the following main properties in the absence of overflow<sup>9</sup>

$$231 \quad \forall x, y \in \mathbb{F}, \quad \exists \delta, \quad |\delta| \leq \epsilon \wedge \circ(x + y) = (1 + \delta)(x + y) \quad (2)$$

$$232 \quad \forall x, y \in \mathbb{F}, \quad \exists \delta, \varphi, \quad |\delta| \leq \epsilon \wedge |\varphi| \leq \eta \wedge \circ(x \times y) = (1 + \delta)(x \times y) + \varphi \quad (3)$$

234 where  $\epsilon$  and  $\eta$  are tiny constants depending on the floating-point format<sup>10</sup>. As a recent  
235 example, the following result is proved in a slightly refined standard model [20].

<sup>9</sup> Overflow can often be handled separately.

<sup>10</sup> For `binary64` and  $\circ$  a rounding to nearest,  $\epsilon = 2^{-53}$  and  $\eta = 2^{-1075}$ .

236 ► **Theorem 2** (Theorem 4.1 in [20]). For  $x \in \mathbb{F}^n$ , denoting  $\hat{s}$  the sum  $\sum_{i=1}^n x_i$  computed with  
 237 floating-point arithmetic in any order<sup>11</sup>, assuming no overflow occurs, it satisfies

$$238 \quad \left| \hat{s} - \sum_{i=1}^n x_i \right| \leq \frac{(n-1)\epsilon}{1+\epsilon} \left( \sum_{i=1}^n |x_i| \right).$$

239 Coq proofs of such results can be performed, and are at the core of the proof of Theorem 1 [33].

## 240 2.3 The Flocq Library

241 Flocq [5, 6] is a Coq library offering a very generic formalization of floating-point arithmetic.  
 242 Radix and precision can be fully parameterized and floating-point values are defined, similarly  
 243 to (1), as a subset of the real numbers  $\mathbb{R}$  provided in the Coq standard library [27, Chapter 1].

244 More specifically, multiple models are available:

- 245 ■ With an unbounded exponent range, i.e., without underflow nor overflow. Although  
 246 unrealistic, this model is attractive for its simplicity and commonly used for error  
 247 bounds [18].
- 248 ■ With an exponent range only lower bounded, i.e., with underflow but without overflow.  
 249 This may still seem unrealistic but overflows can often be studied separately which usually  
 250 proves much harder for underflows [33].
- 251 ■ A binary model of the `binary32` and `binary64` formats defined in the IEEE 754 standard,  
 252 with underflows, overflows to infinities, signed zeros and NaNs with payloads. This model  
 253 is used in the verified C compiler CompCert [4].

254 Along with these models and links between them, the library contains many classical results  
 255 about roundings, about some error-free transformations as presented in Section 2.2.2, and  
 256 basic properties of the standard model described in Section 2.2.3.

257 The library is mainly developed by Sylvie Boldo and Guillaume Melquiond and is available  
 258 at URL <http://flocq.gforge.inria.fr/>.

## 259 2.4 The Coq.Interval Library

260 Another Coq library could benefit from efficient floating-point arithmetic: Coq.Interval [25],  
 261 which offers a modular formalization of interval arithmetic. First, module types (a.k.a. sig-  
 262 natures) are defined for floating-point and interval operators. Then, several implementations  
 263 of the floating-point signature are provided, relying on the Flocq library and specifically its  
 264 model with unbounded exponent range. A generic implementation is provided, as well as  
 265 a specialized implementation assuming radix 2 and representing mantissa and exponent as  
 266 pairs of integers from `Bignums`. Next, a parameterized module implements interval operators  
 267 where intervals are pairs of floating-point numbers, and related computations are performed  
 268 using directed roundings, towards  $-\infty$  or  $+\infty$ . Elementary functions such as `exp`, `ln` or  
 269 `atan` are provided among these interval operators, but correct rounding is not guaranteed  
 270 (namely, the computed intervals can be overestimated, albeit the containment property  
 271 always holds and has been formally proved). Finally, tactics `interval` (decision procedure)  
 272 and `interval_intro` (for forward reasoning) are provided to automatically and formally  
 273 prove inequalities on real-valued expressions.

274 The library is mainly developed by Guillaume Melquiond and is available at URL  
 275 <http://coq-interval.gforge.inria.fr/>.

<sup>11</sup> Floating-point addition is not associative.

### 276 **3 Contributions**

277 In order to provide access to efficient floating-point arithmetic inside proofs, the following  
278 steps have been performed:

- 279 1. Define a minimal working interface for the IEEE 754 `binary64` format. See Section 3.1.
- 280 2. Devise a specification of this interface that enables using `binary64` computations in  
281 proofs. This specification should be compatible with `Flocq`, so that all previously proved  
282 results, both in `Flocq` and based upon it, can be straightforwardly reused, using a simple  
283 compatibility layer. Details are in Section 3.2.
- 284 3. Implement the chosen interface in Coq’s various computation mechanisms, i.e., `compute`,  
285 `vm_compute` and `native_compute` at the OCaml and C levels. A brief summary of the  
286 implementation is given in Section 3.3 and salient points are discussed in Section 4.
- 287 4. Assess the performance by running some benchmarks. Results are given in Section 5.

#### 288 **3.1 Interface**

289 In our modified version of Coq, after typing

```
290 Require Import Floats.
```

293 the user gets access to the following interface<sup>12</sup>:

```
294 Parameter float : Set.
```

297 A type for primitive floating-point values. Inside the kernel, this is mapped to the `float`  
298 type of OCaml<sup>13</sup> that matches `binary64`.

```
299 Parameters add sub mul div : float -> float -> float.  
300 Parameters sqrt opp abs : float -> float.
```

303 The basic arithmetic operators `+`, `-`, `×`, `÷`,  $\sqrt{\cdot}$ , opposite and absolute value.

```
304 Variant float_comparison : Set := FEq | FLt | FGt | FNotComparable.  
305 Parameter compare : float -> float -> float_comparison.
```

308 A comparison function that behaves as specified by the IEEE 754 standard. In particular  
309 `+0` and `-0` are considered equal and NaNs are not comparable to any value, hence the  
310 `FNotComparable` answer.

311 A few functions are then given to examine or craft precise floating-point values by  
312 translating them from or to primitive integers.

```
313 Variant float_class : Set :=  
314 | PNormal | NNormal | PSubn | NSubn | PZero | NZero | PInf | NInf | NaN.  
315 Parameter classify : float -> float_class.
```

318 A function testing whether a given value is a NaN, an infinity (`NInf` and `PInf` for  $-\infty$  and  
319  $+\infty$  respectively), `-0` (`NZero`), `+0` (`PZero`), a denormalized value (`NSubn` and `PSubn`) or a  
320 regular one (`NNormal` and `PNormal`).

<sup>12</sup>Defined in file `theories/Floats/PrimFloat.v` in the implementation.

<sup>13</sup>The implementation language of Coq.

```
321 Definition shift := 2101%int63. (* = 2 × emax + prec *)
```

```
322 Parameter frshiftp : float → float * Int63.int.
```

323 `frshiftp`  $f$  returns a pair  $(m, e)$  such that<sup>14</sup>  $|m| \in [0.5, 1)$  and  $f = m \times 2^{e-\text{shift}}$ . Primitive  
324 integers are unsigned so `shift` is used to ensure that  $e$  is nonnegative.

```
327 Parameter ldshiftp : float → Int63.int → float.
```

328 `ldshiftp`  $f$   $e$  returns  $f \times 2^{e-\text{shift}}$ . This is the reverse of `frshiftp` and it is exact  
329 except when underflow or overflow occurs, in which case the result is rounded using RNE.

```
332 Parameter normfr_mantissa : float → Int63.int.
```

333 When  $f$ , typically obtained from `frshiftp`, satisfies  $|f| \in [0.5, 1)$ , `normfr_mantissa`  $f$   
334 returns the primitive integer  $|f| \times 2^p$ , that is the integer encoding the mantissa of  $f$ .

```
337 Parameter of_int63 : Int63.int → float.
```

338 Converts a primitive integer to a floating-point value. Since primitive integers are unsigned  
339 63-bit integers, they do not all fit into the 53-bit mantissas of the `binary64` format. Values  
340 that do not fit are rounded using RNE.

341 Finally, two functions compute the successor and predecessor of a floating-point value.  
342 They can be used to implement interval arithmetic for instance.

```
345 Parameters next_up, next_down : float → float.
```

346 Equipped with this interface, the Coq user can now perform floating-point computations  
347 using the processor operators and any of the evaluation mechanisms provided by Coq.

```
350 Coq < Require Import Floats. Open Scope float_scope.
```

```
351 Coq < Eval compute in 1 + 0.5.
```

```
352 = 1.5 : float
```

```
353 Coq < Eval vm_compute in 1 / -0.
```

```
354 = neg_infinity : float
```

```
355 Coq < Eval native_compute in 0 / 0.
```

```
356 = nan : float
```

## 359 3.2 Specification

360 Although floating-point computations are possible, they remain entirely useless in proofs at  
361 this point, since there is no specification of their behavior. We thus need a Coq specification  
362 of floating-point arithmetic.

363 First of all, the set of floating-point values itself has to be specified<sup>15</sup>.

```
364 Variant spec_float :=
```

```
365 | S754_zero (sign : bool) (* true for -0, false for +0 *)
```

```
366 | S754_infinity (sign : bool)
```

```
367 | S754_nan
```

```
368 | S754_finite (sign : bool) (mantissa : positive) (exponent : Z).
```

<sup>14</sup>When  $f$  is finite and non zero, otherwise  $(m, e) = (f, 0)$ .

<sup>15</sup>See file `theories/Floats/SpecFloat.v` in the implementation.

### 33:10 Primitive Floats in Coq

371 This is similar to the `full_float` type in the `IEEE754.Binary` module of the `Flocq` library  
 372 except for one point: the sign and payload of NaNs are not modeled here. It is also worth  
 373 noting that this models much more values than the `binary64` format<sup>16</sup> since no bounds on  
 374 mantissas nor exponents are enforced. This makes for a simple specification.

375 Then, each of the above operators must be specified on this `spec_float` type. This  
 376 specification is mostly borrowed<sup>17</sup> from the `IEEE754.Binary` module of the `Flocq` library  
 377 and totals 398 lines in our implementation<sup>18</sup>. We thus only detail the multiplication operator.  
 378 We first need to define a few characteristics of the `binary64` format as seen in Section 2.2.1.1

```
379 Definition prec := 53%Z.
380 Definition emax := 1024%Z.
381 Definition emin := (3 - emax - prec)%Z. (* = -1074 *)
382 Definition fexp e := Z.max (e - prec) emin.
383
384
```

385 When  $|x| \in [2^{e-1}, 2^e)$ , then `fexp e` is the exponent used to encode  $x$  in the `binary64` format.

386 As seen in Section 2.2.1.2, the floating point multiplication is defined by  $x \otimes y = \circ(x \times y)$ .  
 387 When  $x = m_x 2^{e_x}$  and  $y = m_y 2^{e_y}$ , then  $x \times y = (m_x \times m_y) 2^{e_x + e_y}$  and the rounding operator  
 388  $\circ$  has to remove the extra bits in the mantissa to make this value fit in the format. To this  
 389 end, we first abstract the bits to remove as two booleans, the *rounding* bit remembers the  
 390 first forgotten bit whereas the *sticky* bit is `true` when any of the remaining forgotten bits is  
 391 1 and `false` when they are all 0. The function `shr_1` then shifts a mantissa one bit to the  
 392 right, updating the rounding and sticky bits accordingly

```
393 Record shr_record := { shr_m : Z ; shr_r : bool ; shr_s : bool }.
394 Definition shr_1 mrs :=
395   let s := orb (shr_r mrs) (shr_s mrs) in match shr_m mrs with
396     | Z0 (* 0 *) => Build_shr_record Z0 false s
397     | Zpos xH (* 1 *) => Build_shr_record Z0 true s
398     | Zpos (x0 p) (* 2p *) => Build_shr_record (Zpos p) false s
399     | Zpos (xI p) (* 2p+1 *) => Build_shr_record (Zpos p) true s
400     | ... (* same for Zneg _ *) end.
401
402
```

403 Eventually, `shr` can iterate  $n$  shifts and `shr_fexp` removes the required number of bits using  
 404 the above function `fexp` (`Zdigits2 m` is the number of bits of  $m$ )

```
405 Definition shr mrs e n := match n with
406   | Zpos p => (iter_pos shr_1 p mrs, (e + n)%Z) | _ => (mrs, e) end.
407 Definition shr_fexp m e :=
408   shr (Build_shr_record m false false) e (fexp (Zdigits2 m + e) - e).
409
410
```

411 It now remains to round the mantissa according to the values of the rounding and sticky bits

```
412 Definition round_nearest_even mrs := match mrs with
413   | Build_shr_record mx false _ => mx
414   | Build_shr_record mx true false => if Z.even mx then mx else (mx + 1)%Z
415   | Build_shr_record mx true true => (mx + 1)%Z end.
416
417
```

<sup>16</sup> `spec_float` gathers an infinite number of values, whereas `binary64` only contains finitely many values.

<sup>17</sup> Except for the specifications of `frexp`, `ldexp`, `normfr_mantissa`, `succ` and `pred` which were not yet present in `Flocq` and which we took the opportunity to add [https://gitlab.inria.fr/flocq/flocq/merge\\_requests/3](https://gitlab.inria.fr/flocq/flocq/merge_requests/3).

<sup>18</sup> See file `theories/Floats/SpecFloat.v` in the implementation.

418 Finally, the rounding function first shifts the mantissa, rounds it, shifts the result one bit to  
419 the right in case the rounding added an extra bit and handles potential overflows

```
420
421 Definition binary_round_aux sx mx ex :=
422   let '(mrs', e') := shr_fexp mx ex in
423   let '(mrs'', e'') := shr_fexp (round_nearest_even mrs') e'
424   in match shr_m mrs'' with Z0 => S754_zero sx | Zneg _ => S754_nan
425   | Zpos m => if Zle_bool e'' (emax - prec) then S754_finite sx m e''
426   else S754_infinity sx end.
427
```

428 Thus, it remains to the multiplication to handle all particular cases

```
429
430 Definition SFmul x y := match x, y with
431   | S754_nan, _ | _, S754_nan => S754_nan
432   | S754_infinity sx, S754_infinity sy => S754_infinity (xorb sx sy)
433   | S754_infinity sx, S754_finite sy _ => S754_infinity (xorb sx sy)
434   | S754_finite sx _ , S754_infinity sy => S754_infinity (xorb sx sy)
435   | S754_infinity _, S754_zero _ => S754_nan
436   | S754_zero _, S754_infinity _ => S754_nan
437   | S754_finite sx _ , S754_zero sy => S754_zero (xorb sx sy)
438   | S754_zero sx, S754_finite sy _ => S754_zero (xorb sx sy)
439   | S754_zero sx, S754_zero sy => S754_zero (xorb sx sy)
440   | S754_finite sx mx ex, S754_finite sy my ey =>
441   binary_round_aux (xorb sx sy) (Zpos (mx * my)) (ex + ey) end.
442
```

443 In addition to the usual operators, two functions are defined going back and forth from  
444 primitive floats to specification floats.

```
445 Definition Prim2SF : float -> spec_float.
446 Definition SF2Prim : spec_float -> float.
447
448
```

449 Finally, one needs to establish a link between the primitive operators and the specification.  
450 This is done by adding axioms to the system.<sup>19</sup> First, to specify the two functions `Prim2SF`  
451 and `SF2Prim` above, one needs to characterize those values of type `spec_float` that actually  
452 represent a `binary64` floating-point number, i.e., values with appropriately bounded mantissa  
453 and exponent.

```
454 Definition canonical_mantissa m e := Zeq_bool (fexp (Zdigits2 m + e)) e.
455 Definition bounded m e :=
456   andb (canonical_mantissa m e) (Zle_bool e (emax - prec)).
457 Definition valid_binary x := match x with
458   | SF754_finite _ m e => bounded m e | _ => true end.
459
460
```

461 Again, this code comes from the `Flocq` library [5]. So equipped, the following three axioms  
462 can be stated:

```
463 Axiom Prim2SF_valid : forall x, valid_binary (Prim2SF x) = true.
464 Axiom SF2Prim_Prim2SF : forall x, SF2Prim (Prim2SF x) = x.
465 Axiom Prim2SF_SF2Prim :
466   forall x, valid_binary x = true -> Prim2SF (SF2Prim x) = x.
467
468
```

<sup>19</sup>See file `theories/Floats/FloatAxioms.v` in the implementation.

## 33:12 Primitive Floats in Coq

469 These properties allow one to prove that both `Prim2SF` and `SF2Prim` are injective and thereby  
470 form a bijection between primitive floats and the subset of valid specification floats.

```
471 Theorem Prim2SF_inj : forall x y, Prim2SF x = Prim2SF y -> x = y.
```

```
473 Theorem SF2Prim_inj : forall x y, SF2Prim x = SF2Prim y ->
```

```
474   valid_binary x = true -> valid_binary y = true -> x = y.
```

476 Thus, all of the fifteen operators given in Section 3.1 are linked to their specification by an  
477 axiom such as, for the multiplication:

```
478 Axiom mul_spec :
```

```
480   forall x y, Prim2SF (x * y)%float = SFmul (Prim2SF x) (Prim2SF y).
```

482 Since the specification is almost identical to the `IEEE754.Binary` module of `Flocq`, a link  
483 with `Flocq` is straightforwardly built<sup>20</sup>, establishing a bridge towards real numbers and giving  
484 access to all the results already proved in the library. This plays a key role in enabling actual  
485 proofs using primitive floating-point computations. Moreover, this enables to gain additional  
486 confidence in the above non trivial specification, since `Flocq` contains correctness theorems  
487 basically stating for instance<sup>21</sup> that, except when overflow occurs, `SFmul x y` is indeed the  
488 rounding of the real number  $x \times y$ .

### 489 3.3 Implementation

490 The implementation was submitted to be integrated in Coq through the GitHub pull request  
491 <https://github.com/coq/coq/pull/9867>.

492 Below is an overview of the size of the development at the time of writing, summarized  
493 by sub-components (over the  $\approx 3.7$  kLoC added).

494 ■ OCaml and C: 1815 LoC

495 (floats  $\leftrightarrow$  kernel : 1070) (`vm_compute` support: 255) (`native_compute` support: 355)

496 (parsing and pretty-printing: 85) (Coq checker: 50)

497 ■ Coq specifications: 620 LoC [mostly borrowed from `Flocq`]

498 ■ Coq proofs: 340 LoC

499 ■ Tests: 800 LoC

500 ■ Sphinx documentation: 115 LoC

501 This implementation required the addition of some code in the kernel of Coq. Most of it  
502 only consists in wrapping the floating-point operators into the different evaluation mecha-  
503 nisms of Coq and its core, actually dealing with floating-point arithmetic, can be found in  
504 the files `kernel/float64.ml`, `kernel/byterun/coq_interp.c` and `kernel/byterun/coq_`  
505 `float64.h`. Most operators are implemented in C, as required by the `vm_compute` mechanism,  
506 and boil down to calls to the appropriate functions of the C standard library. Thus, no  
507 involved algorithmic happens in this added code itself.

## 508 4 Discussion

### 509 4.1 Rounding Modes

510 We implement only one of the five rounding modes defined in the IEEE 754-2008 standard,  
511 namely rounding to nearest even (RNE). We argue here that implementing other rounding

---

<sup>20</sup> See [https://gitlab.inria.fr/flocq/flocq/merge\\_requests/6](https://gitlab.inria.fr/flocq/flocq/merge_requests/6).

<sup>21</sup> See theorem `Bmult_correct` in module `Flocq.IEEE754.Binary`.

512 modes would not only easily be seriously harmful in terms of performance, notwithstanding  
 513 the potential threat to soundness of the implementation, but also not very useful.

514 Unfortunately on most common processors, operators with different rounding modes  
 515 are not implemented using different opcodes but a status flag. Once the flag is set to a  
 516 particular rounding mode, all subsequent computations are performed with this rounding  
 517 mode. Changing the rounding mode is then costly as it requires flushing pipelines.

518 Interval arithmetic constitutes the main use of rounding modes other than RNE we can  
 519 foresee in a proof assistant. A common solution to the aforementioned performance issue is  
 520 to set the rounding mode once to  $+\infty$  (RU), used to compute upper bounds, and emulate  
 521 rounding toward  $-\infty$  (RD), used to compute lower bounds, by relying on properties like<sup>22</sup>  
 522  $\text{RD}(x + y) = -\text{RU}((-x) + (-y))$ . Although a monadic interface could be a reasonable  
 523 implementation, this remains an imperative programming feature and doesn't integrate well  
 524 within the functional paradigm offered by Coq. Moreover, if no particular care is taken to  
 525 avoid or disable them, wild compiler optimizations—assuming that only RNE is used—could  
 526 easily break the previous property, thus ruining the soundness of the whole approach.

527 However, interval arithmetic doesn't require precise directed roundings but only over-  
 528 and under-approximations thereof. We thus offer the `next_up` and `next_down` functions,  
 529 computing the successor and predecessor of a floating-point value. Together with rounding  
 530 to nearest operators, they satisfy the following property, ensuring soundness of interval  
 531 arithmetic while providing a reasonably precise approximation of directed roundings:

$$532 \quad \forall x \in \mathbb{R}, \quad \text{RU}(x) \leq \text{next\_up}(\text{RNE}(x))$$

$$533 \quad \forall x \in \mathbb{R}, \quad \text{next\_down}(\text{RNE}(x)) \leq \text{RD}(x).$$

## 535 4.2 Parsing and Pretty-Printing

536 Parsing and pretty-printing floating-point values is a non trivial question. We expect the  
 537 following main property: printing a floating-point value and then reparsing the output of  
 538 the printing function should give the initial value, i.e.,  $\text{parse} \circ \text{print}$  should be the identity  
 539 over `binary64`. It is worth noting that this necessarily implies the injectivity of the printing  
 540 function. However, we don't require the parsing function to be injective, i.e., we do accept  
 541 that multiple strings are parsed as the same floating-point value.

542 A simple solution would be to print an exact hexadecimal representation of the floating-  
 543 point values, with a binary exponent, e.g., “0xcp-3”. This fulfills the above requirement.  
 544 Unfortunately, this is not very user-friendly. A decimal output would be much more human  
 545 readable, e.g., “1.5” instead of “0xcp-3”.

546 It is known that printing `binary64` values using at least 17 significant digits and imple-  
 547 menting parsing as a rounding to nearest guarantees the above requirements [30, Table 2.3,  
 548 p. 44]. This is thus the adopted solution. The current version of Coq only offers support for  
 549 parsing and printing integer constants, so we extended this support<sup>23</sup> to decimal constants  
 550 using the ubiquitous format  $\langle \text{integer\_part} \rangle . \langle \text{fractional\_part} \rangle \text{e} \langle \text{decimal\_exponent} \rangle$ , e.g., “1.23e-4”.

## 551 4.3 Soundness

552 During our development, we identified three main potential threats to soundness:

<sup>22</sup>The opposite  $x \mapsto -x$  being exact in floating-point arithmetic (the sign bit is simply flipped).

<sup>23</sup>See the pull request <https://github.com/coq/coq/pull/8764>.

553 **Specification Issues** due to a mismatch w.r.t. the implementation would break the soundness.

554 We hope that taking in extenso our specification from the Flocq library, resulting from  
 555 a few decades of experience in the field and proving links with other models, mitigates  
 556 this risk. Moreover, such an error in the specification can only be harmful when the  
 557 corresponding axiom is used. It is worth noting that all the axioms used in a proved  
 558 theorem explicitly appear in the result of the Coq command `Print Assumptions`.

559 **Incompatible Implementations** in different evaluation mechanisms (`compute`, `vm_compute`  
 560 or `native_compute`) or even on different machines could lead to a proof of `False` by  
 561 evaluating a same term to different results. For instance, the payload of NaNs is not fully  
 562 specified by the IEEE 754 standard and different hardwares can produce different NaNs  
 563 for a same computation. That's why we chose to consider all NaNs as equal and not  
 564 distinguish them. Thus incompatible implementations at the bit level remain compatible  
 565 at the logical level. Double roundings due to the x87 on old 32 bits architectures [29]  
 566 could also be harmful. The OCaml<sup>24</sup> compiler systematically relies on it, forcing us to  
 567 implement all floating-point operators in C and to use the appropriate compiler flags. A  
 568 runtime test<sup>25</sup> is eventually added to prevent Coq from running in case of miscompilation.  
 569 Another extreme example of implementation discrepancy would be a hardware bug such  
 570 as the one encountered in the division of the early Pentium processors.

571 **Incorrect Convertibility Test** that distinguish two values that shouldn't or vice versa is also  
 572 a threat. For instance, implementing this test using the equality test on floating-point  
 573 values (as defined in the IEEE 754 standard) would be wrong as it equates  $-0$  and  $+0$   
 574 which should be distinguished since  $1 \div (-0) = -\infty \neq 1 \div (+0) = +\infty$ . Fortunately  
 575 enough, this keeps a very simple implementation, with the following OCaml code:

```
576 let equal f1 f2 =  

577   let is_nan f = f <> f in  

578   match classify_float f1 with  

579   | FP_normal | FP_subnormal | FP_infinite -> f1 = f2  

580   | FP_nan -> is_nan f2 | FP_zero -> f1 = f2 && 1. /. f1 = 1. /. f2  

581  

582
```

583 A few other, more minor, points appeared during the development. Among them, the fact  
 584 that primitive integers in Coq are unsigned did require some care<sup>26</sup>. Finally, the way OCaml  
 585 optimizes arrays<sup>27</sup> of floating-point values<sup>28</sup> did cause a few nasty bugs, although it is unlikely  
 586 that such bugs could lead to a proof of `False` as they often yield a mere segmentation fault.

## 587 **5** Benchmarks

588 The overall objective of this work is to increase the performance of reflexive tactics involving  
 589 floating-point arithmetic in Coq. Thus we first measure the performance gain on such a tactic,  
 590 then evaluate it on its individual floating-point operators. We first present the reference  
 591 problems under study (Section 5.1), then recap the hardware and software setup for these  
 592 benchmarks (Section 5.2), and finally give the experimental results (Section 5.3).

<sup>24</sup> The implementation language of Coq.

<sup>25</sup> See file `kernel/float64.ml` in the implementation.

<sup>26</sup> We indeed fixed a few soundness bugs in primitive integers, pertaining with unsigned integers, before they were merged in Coq master development branch (<https://github.com/coq/coq/pull/6914>).

<sup>27</sup> Arrays are used to communicate environments between the OCaml implementation of the kernel and the C implementation of the `vm_compute` virtual machine.

<sup>28</sup> This causes other issues in OCaml itself and seems to be a hot topic currently in the OCaml community [9].

## 5.1 Reference Test-suite

We developed a reflexive tactic `posdef_check`, performing some matrix positive definiteness check along the lines of Theorem 1 introduced in Section 1.1. Its implementation was adapted by reusing building blocks from our previous work on the `validsdp` tactic for multivariate polynomial positivity [26].

This tactic is available in four flavors using `vm_compute` or `native_compute` and emulated floats or primitive floats. Emulated floats are a state of the art implementation of floating point arithmetic, based on primitive integers, from the `Coq.Interval` library whereas primitive floats are our new implementation.

Regarding the test-suite, we generated a set of random positive definite matrices (after fixing a given seed to make the random data reproducible) of size  $50 \times 50$  up to  $400 \times 400$ .

We perform two kinds of benchmarks on this test-suite: the overall speedup between the versions of `posdef_check` using emulated vs. primitive floats; and the individual speedup in floating-point operators involved in this tactic.

## 5.2 Hardware/Software Setup

The formalization of the `posdef_check` tactic relies on a large set of dependencies that takes around one hour to compile. For greater convenience, we devised some Docker images containing the benchmark environment, based on Debian Stretch, `opam 2` (the OCaml package manager) and OCaml 4.07.0+`flambda`. The source code of all benchmarks as well as guidelines to install Docker and run the benchmarks are gathered on GitHub at this URL: <https://github.com/validsdp/benchs-primitive-floats/tree/1.0>

The use of Docker (a so-called *OS-level virtualization system*) for these benchmarks yields a number of interesting features, beyond the facility to download and run a pre-built image on different machines: it runs containers in an isolated environment from the host machine, it ensures portability (across OSes such as GNU/Linux, macOS and Windows) and reproducibility, while being more lightweight than traditional virtual machines (VMs).

The experimental results of the upcoming Section 5.3 have been obtained using a Debian GNU/Linux workstation based on a Intel Core i7-7700 CPU clocked at 3.60 GHz, with 16 GB of RAM. All benchmarks have been executed sequentially (namely, without the `-j` option of `make`), with a total elapsed time of about 3h35', using the following image: `"docker pull registry.gitlab.com/erikmd/docker-coq-primitive-floats/master_compiler-edge:9_coq-2ac1f46532264bacf2b1d8f5b6ee3659fe0cde67"`.

## 5.3 Experimental Results

We first measure the execution time of the whole tactic on the test-suite and compare it between emulated floats and primitive floats. The results are displayed in Table 1 for `vm_compute` and `native_compute`. Each timing is measured 5 times. The tables indicate the corresponding average and relative error among the 5 samples.

One can notice that the obtained speedups are far from the three order of magnitudes separating emulated floats from equivalent OCaml implementations. From the above results, it appears that arithmetic operators constitute most of the computation time with emulated floats (at least 95% with `vm_compute`) but nothing tells us this is still the case with primitive floats. In fact, with primitive floats, most of the computation time is dedicated to list

■ **Table 1** Proof time for the reflexive tactic `posdef_check`.

Source	vm_compute			native_compute		
	Emulated	Primitive	Diff.	Emulated	Primitive	Diff.
mat050	0.16s ±2.0%	0.01s ±0.0%	20x	0.05s ±4.0%	0.02s ±5.1%	3x
mat100	1.16s ±1.3%	0.06s ±5.8%	21x	0.28s ±2.5%	0.03s ±2.5%	9x
mat150	3.61s ±1.2%	0.18s ±2.2%	21x	0.75s ±3.0%	0.08s ±3.5%	9x
mat200	8.68s ±0.2%	0.41s ±1.0%	21x	1.71s ±1.0%	0.18s ±3.4%	10x
mat250	17.14s ±1.3%	0.80s ±0.3%	21x	3.34s ±1.4%	0.33s ±2.1%	10x
mat300	30.01s ±1.2%	1.37s ±0.7%	22x	5.77s ±2.4%	0.56s ±1.0%	11x
mat350	48.31s ±1.3%	2.15s ±0.1%	23x	9.09s ±3.0%	0.81s ±1.2%	11x
mat400	70.19s ±1.4%	3.18s ±0.5%	22x	13.56s ±4.0%	1.12s ±0.7%	12x

manipulating functions as our matrices are implemented using lists<sup>29</sup> [8]. Thus, we would like to get an idea of the time actually devoted to floating-point arithmetic in the total proof time of our reflexive tactic. We use the following simple methodology: replace each arithmetic operator with a version, uselessly, performing the computation twice<sup>30</sup>, then subtract the execution time of the original program (“Op” in the tables) to the one of this modified program (“Op×2” in the tables). The obtained time (“Op time” in the tables) corresponds to the time devoted to the considered arithmetic operator. Note that the redundant computations involved in the modified program (“Op×2”) could not be implemented with a mere additional let-in such as `...let m1 := mul a b in let m2 := mul a b in m2` because the virtual machine and the OCaml native compiler would optimize away the unused local definition; but doing so and adding an extra function call `...in select m1 m2` with `Definition select (a b : F.type) := a.` made it possible to use this doubling trick. The results are given in Table 2 for `vm_compute` and Table 3 for `native_compute`, in each case both on addition and multiplication<sup>31</sup>. Again, each timing is measured 5 times. It is worth noting that those last results should be taken more as coarse orders of magnitude than precise results. In particular, due to the overhead stemming from the duplication itself of the operators<sup>32</sup>, the speedups are—maybe seriously—underapproximated. Actual speedups could thus be higher than the ones suggested here.

## 6 Conclusion and Future Work

We developed a theory of floating-point arithmetic for the Coq proof assistant, composed of primitive implementation of basic arithmetic operators ( $+$ ,  $-$ ,  $\times$ ,  $\div$ ,  $\sqrt{\cdot}$ ), using the processor floating-point operators in rounding-to-nearest even, as well as successor and predecessor operators that can be used to approximate directed roundings to  $-\infty$  or  $+\infty$ . This implementation is axiomatized under the assumption that the processor complies with the IEEE 754 standard for floating-point arithmetic. Particular care has been taken to make the implementation compatible across the different reduction engines of Coq, and across different hardware, thereby avoiding soundness issues that could be caused, for example, by the semantics of NaN payloads that is under-specified in the IEEE 754 standard.

<sup>29</sup>This could be improved using primitive “persistent arrays” once they will be integrated in Coq [1].

<sup>30</sup>Or thousand times for primitive floats to avoid getting a result of the same order of magnitude than the variability of computation times.

<sup>31</sup>Additions and multiplications constitute the vast majority of the arithmetic computations performed in a Cholesky decomposition, as seen in Figure 1.

<sup>32</sup>Like expensive function calls.

■ **Table 2** Computation time for individual operators with `vm_compute`.

Op Source	Emulated floats		Op	Primitive floats		Diff.
	CPU times (Op×2–Op)			CPU times (Op×1001–Op)		
add						
mat200	10.78±0.9%	– 8.38±2.8%	2.40s	15.72±0.5%	– 0.45±1.1%	0.02s 157x
mat250	21.46±1.7%	– 16.41±1.5%	5.06s	30.62±0.6%	– 0.82±0.6%	0.03s 170x
mat300	37.43±1.4%	– 28.63±1.4%	8.80s	53.12±2.4%	– 1.40±0.5%	0.05s 170x
mat350	59.42±0.8%	– 45.95±2.9%	13.48s	84.19±0.8%	– 2.19±0.5%	0.08s 164x
mat400	87.78±0.9%	– 66.17±1.7%	21.61s	127.56±8.5%	– 3.21±0.3%	0.12s 174x
mul						
mat200	12.21±1.4%	– 8.38±2.8%	3.83s	16.10±3.0%	– 0.45±1.1%	0.02s 245x
mat250	24.52±1.4%	– 16.41±1.5%	8.11s	31.12±3.7%	– 0.82±0.6%	0.03s 268x
mat300	42.84±1.7%	– 28.63±1.4%	14.21s	53.25±0.8%	– 1.40±0.5%	0.05s 274x
mat350	68.23±1.5%	– 45.95±2.9%	22.28s	84.33±0.7%	– 2.19±0.5%	0.08s 271x
mat400	99.72±1.5%	– 66.17±1.7%	33.55s	125.74±0.8%	– 3.21±0.3%	0.12s 274x

■ **Table 3** Computation time for individual operators with `native_compute`.

Op Source	Emulated floats		Op	Primitive floats		Diff.
	CPU times (Op×2–Op)			CPU times (Op×1001–Op)		
add						
mat200	2.24±1.4%	– 1.78±1.7%	0.46s	17.68±1.4%	– 0.22±0.9%	0.02s 27x
mat250	4.49±4.2%	– 3.41±3.1%	1.08s	34.29±0.7%	– 0.37±1.5%	0.03s 32x
mat300	7.25±1.2%	– 5.83±4.6%	1.42s	59.57±2.5%	– 0.55±0.9%	0.06s 24x
mat350	11.66±3.8%	– 9.28±3.5%	2.39s	93.82±1.1%	– 0.82±0.8%	0.09s 26x
mat400	17.07±2.9%	– 13.14±0.9%	3.93s	141.97±2.6%	– 1.18±0.9%	0.14s 28x
mul						
mat200	2.48±1.5%	– 1.78±1.7%	0.70s	17.81±1.1%	– 0.22±0.9%	0.02s 40x
mat250	4.82±2.4%	– 3.41±3.1%	1.41s	35.14±2.1%	– 0.37±1.5%	0.04s 41x
mat300	8.41±2.4%	– 5.83±4.6%	2.59s	60.66±2.2%	– 0.55±0.9%	0.06s 43x
mat350	13.21±2.4%	– 9.28±3.5%	3.94s	97.25±1.0%	– 0.82±0.8%	0.10s 41x
mat400	19.27±1.5%	– 13.14±0.9%	6.13s	138.61±2.3%	– 1.18±0.9%	0.14s 45x

663 We evaluated the performance on an implementation—carried out in Gallina, the input  
 664 language of Coq—of a Cholesky decomposition that underlies a reflexive tactic for matrix pos-  
 665 itive definiteness, and the experimental results indicate a speedup of two orders of magnitude  
 666 for arithmetic operators using `vm_compute`. This is consistent with the performance factor of  
 667 about three orders of magnitude observed between floating-point arithmetic emulated using  
 668 primitive integers in Coq and equivalent implementations written in OCaml.

669 Now that primitive floats are available in a proof assistant, multiple future works can  
 670 be envisioned. The most obvious one would be to adapt the `Coq.Interval` library to take  
 671 advantage of primitive floats. Still in this direction, it is known that the successor and  
 672 predecessor functions, used to approximate directed roundings, can be efficiently implemented  
 673 using only arithmetic operators [36, 38]. Such an implementation could enable to remove  
 674 these functions from the trusted code base. It would also be interesting to look at more  
 675 elaborate elementary functions such as `exp` or `arctan`, relying for example on the `CR-libm`  
 676 implementation [10]. Finally, in an attempt to improve confidence in the consistency between  
 677 specification and implementation, and while waiting for a fully formally specified hardware  
 678 interface, it is worth noting that this consistency is amenable to some intensive automatic  
 679 testing, although exhaustive testing is out of reach for even unary operators on `binary64`.

## 680 — References

- 681 **1** Michaël Armand, Benjamin Grégoire, Arnaud Spiwack, and Laurent Théry. Extending coq  
682 with imperative features and its application to SAT verification. In Matt Kaufmann and  
683 Lawrence C. Paulson, editors, *Interactive Theorem Proving, First International Conference,*  
684 *ITP 2010, Edinburgh, UK, July 11-14, 2010. Proceedings*, volume 6172 of *Lecture Notes in*  
685 *Computer Science*, pages 83–98. Springer, 2010. doi:10.1007/978-3-642-14052-5\_8.
- 686 **2** Henk Barendregt and Herman Geuvers. Proof-assistants using dependent type systems. In  
687 John Alan Robinson and Andrei Voronkov, editors, *Handbook of Automated Reasoning (in 2*  
688 *volumes)*, pages 1149–1238. Elsevier and MIT Press, 2001.
- 689 **3** Mathieu Boespflug, Maxime Dénès, and Benjamin Grégoire. Full Reduction at Full Throttle.  
690 In Jean-Pierre Jouannaud and Zhong Shao, editors, *Certified Programs and Proofs - First*  
691 *International Conference, CPP 2011, Kenting, Taiwan, December 7-9, 2011. Proceedings*,  
692 volume 7086 of *Lecture Notes in Computer Science*, pages 362–377. Springer, 2011. doi:  
693 10.1007/978-3-642-25379-9\_26.
- 694 **4** Sylvie Boldo, Jacques-Henri Jourdan, Xavier Leroy, and Guillaume Melquiond. A formally-  
695 verified C compiler supporting floating-point arithmetic. In Alberto Nannarelli, Peter-Michael  
696 Seidel, and Ping Tak Peter Tang, editors, *21st IEEE Symposium on Computer Arithmetic,*  
697 *ARITH 2013, Austin, TX, USA, April 7-10, 2013*, pages 107–115. IEEE Computer Society,  
698 2013. URL: <http://ieeexplore.ieee.org/xpl/mostRecentIssue.jsp?punumber=6545866>,  
699 doi:10.1109/ARITH.2013.30.
- 700 **5** Sylvie Boldo and Guillaume Melquiond. Flocq: A unified library for proving floating-point  
701 algorithms in coq. In Elisardo Antelo, David Hough, and Paolo Ienne, editors, *20th IEEE*  
702 *Symposium on Computer Arithmetic, ARITH 2011, Tübingen, Germany, 25-27 July 2011*,  
703 pages 243–252. IEEE Computer Society, 2011. URL: [http://ieeexplore.ieee.org/xpl/](http://ieeexplore.ieee.org/xpl/mostRecentIssue.jsp?punumber=5991607)  
704 [mostRecentIssue.jsp?punumber=5991607](http://ieeexplore.ieee.org/xpl/mostRecentIssue.jsp?punumber=5991607), doi:10.1109/ARITH.2011.40.
- 705 **6** Sylvie Boldo and Guillaume Melquiond. *Computer Arithmetic and Formal Proofs: Verifying*  
706 *Floating-point Algorithms with the Coq System*. Elsevier, 2017.
- 707 **7** Samuel Boutin. Using reflection to build efficient and certified decision procedures. In Martín  
708 Abadi and Takayasu Ito, editors, *Theoretical Aspects of Computer Software, Third International*  
709 *Symposium, TACS '97, Sendai, Japan, September 23-26, 1997, Proceedings*, volume 1281 of  
710 *Lecture Notes in Computer Science*, pages 515–529. Springer, 1997. doi:10.1007/BFb0014565.
- 711 **8** Cyril Cohen, Maxime Dénès, and Anders Mörtberg. Refinements for free! In Georges Gonthier  
712 and Michael Norrish, editors, *Certified Programs and Proofs*, volume 8307 of *Lecture Notes*  
713 *in Computer Science*, pages 147–162. Springer, 2013. URL: [http://dx.doi.org/10.1007/](http://dx.doi.org/10.1007/978-3-319-03545-1_10)  
714 [978-3-319-03545-1\\_10](http://dx.doi.org/10.1007/978-3-319-03545-1_10).
- 715 **9** Simon Colin, Rodolphe Lepigre, and Gabriel Scherer. Unboxing mutually recursive type  
716 definitions in ocaml. In *Proceedings of JFLA, Les Rousses, France, 30th January to 2nd*  
717 *February 2019.*, 2019.
- 718 **10** Catherine Daramy-Loirat, David Defour, Florent De Dinechin, Matthieu Gallet, Nicolas  
719 Gast, Christoph Lauter, and Jean-Michel Muller. Cr-libm. a library of correctly rounded  
720 elementary functions in double-precision. Technical report, LIP, ENS Lyon, 2006. URL:  
721 <https://hal-ens-lyon.archives-ouvertes.fr/ensl-01529804/>.
- 722 **11** T. J. Dekker. A floating-point technique for extending the available precision. *Numerische*  
723 *Mathematik*, 18(3):224–242, 1971.
- 724 **12** Maxime Dénès. Towards primitive data types for COQ 63-bits integers and persistent arrays.  
725 In *Coq-5, the Coq Workshop 2013*, Rennes, France, July 2013. Extended abstract. URL:  
726 [https://coq.inria.fr/files/coq5\\_submission\\_2.pdf](https://coq.inria.fr/files/coq5_submission_2.pdf).
- 727 **13** Burak Ekici, Alain Mebsout, Cesare Tinelli, Chantal Keller, Guy Katz, Andrew Reynolds,  
728 and Clark W. Barrett. Smtcoq: A plug-in for integrating SMT solvers into Coq. In Rupak  
729 Majumdar and Viktor Kuncak, editors, *Computer Aided Verification - 29th International*  
730 *Conference, CAV 2017, Heidelberg, Germany, July 24-28, 2017, Proceedings, Part II*, volume

- 731 10427 of *Lecture Notes in Computer Science*, pages 126–133. Springer, 2017. doi:10.1007/  
732 978-3-319-63390-9\_7.
- 733 **14** Georges Gonthier. Formal Proof—The Four-Color Theorem. *Notices of the American*  
734 *Mathematical Society*, 55(11):1382–1393, 2008. URL: [http://www.ams.org/notices/200811/  
735 tx081101382p.pdf](http://www.ams.org/notices/200811/tx081101382p.pdf).
- 736 **15** Benjamin Grégoire and Xavier Leroy. A compiled implementation of strong reduction. In  
737 Mitchell Wand and Simon L. Peyton Jones, editors, *Proceedings of the Seventh ACM SIGPLAN*  
738 *International Conference on Functional Programming (ICFP '02), Pittsburgh, Pennsylvania,*  
739 *USA, October 4-6, 2002.*, pages 235–246. ACM, 2002. doi:10.1145/581478.581501.
- 740 **16** Benjamin Grégoire and Laurent Théry. A Purely Functional Library for Modular Arithmetic  
741 and Its Application to Certifying Large Prime Numbers. In Ulrich Furbach and Natarajan  
742 Shankar, editors, *IJCAR*, volume 4130 of *Lecture Notes in Computer Science*, pages 423–437.  
743 Springer, 2006. doi:10.1007/11814771\_36.
- 744 **17** Thomas Hales, Mark Adams, Gertrud Bauer, Tat Dat Dang, John Harrison, Le Truong  
745 Hoang, Cezary Kaliszyk, Victor Magron, Sean Mclaughlin, Tat Thang Nguyen, and et al.  
746 a formal proof of the Kepler conjecture. *Forum of Mathematics, Pi*, 5:e2, 2017. 29 pages.  
747 doi:10.1017/fmp.2017.1.
- 748 **18** Nicholas Higham. *Accuracy and Stability of Numerical Algorithms*. Society for Industrial and  
749 Applied Mathematics, Philadelphia, PA, USA, 1996.
- 750 **19** IEEE Computer Society. IEEE Standard for Floating-Point Arithmetic. *IEEE Standard*  
751 *754-2008*, 2008.
- 752 **20** Claude-Pierre Jeannerod and Siegfried M. Rump. On relative errors of floating-point operations:  
753 Optimal bounds and applications. *Math. Comput.*, 87(310):803–819, 2018. doi:10.1090/  
754 mcom/3234.
- 755 **21** D. Knuth. *The Art of Computer Programming*, volume 2. Addison Wesley, Reading, MA,  
756 1973.
- 757 **22** Peter Kornerup, Vincent Lefèvre, Nicolas Louvet, and Jean-Michel Muller. On the computation  
758 of correctly rounded sums. *IEEE Trans. Computers*, 61(3):289–298, 2012. doi:10.1109/TC.  
759 2011.27.
- 760 **23** Andreas Lochbihler. Fast Machine Words in Isabelle/HOL. In Jeremy Avigad and Assia  
761 Mahboubi, editors, *Interactive Theorem Proving - 9th International Conference, ITP 2018,*  
762 *Held as Part of the Federated Logic Conference, FloC 2018, Oxford, UK, July 9-12, 2018,*  
763 *Proceedings*, volume 10895 of *Lecture Notes in Computer Science*, pages 388–410. Springer,  
764 2018. doi:10.1007/978-3-319-94821-8\_23.
- 765 **24** Victor Magron, Xavier Allamigeon, Stéphane Gaubert, and Benjamin Werner. Formal  
766 proofs for nonlinear optimization. *Journal of Formalized Reasoning*, 8(1):1–24, 2015. URL:  
767 <http://jfr.unibo.it/article/view/4319>.
- 768 **25** Érik Martin-Dorel and Guillaume Melquiond. Proving tight bounds on univariate expressions  
769 with elementary functions in Coq. *Journal of Automated Reasoning*, 57(3):187–217, October  
770 2016. doi:10.1007/s10817-015-9350-4.
- 771 **26** Érik Martin-Dorel and Pierre Roux. A reflexive tactic for polynomial positivity using numerical  
772 solvers and floating-point computations. In Yves Bertot and Viktor Vafeiadis, editors, *Proceed-*  
773 *ings of the 6th ACM SIGPLAN Conference on Certified Programs and Proofs, CPP 2017, Paris,*  
774 *France, January 16-17, 2017*, pages 90–99. ACM, 2017. doi:10.1145/3018610.3018622.
- 775 **27** Micaela Mayero. *Formalisation et automatisé de preuves en analyses réelle et numérique*.  
776 PhD thesis, Université Paris VI, décembre 2001.
- 777 **28** O. Møller. Quasi double-precision in floating-point addition. *BIT*, 5:37–50, 1965.
- 778 **29** David Monniaux. The pitfalls of verifying floating-point computations. *ACM Trans. Program.*  
779 *Lang. Syst.*, 30(3):12:1–12:41, 2008. URL: <https://doi.org/10.1145/1353445.1353446>, doi:  
780 10.1145/1353445.1353446.

- 781 **30** Jean-Michel Muller, Nicolas Brisebarre, Florent de Dinechin, Claude-Pierre Jeannerod, Vincent  
782 Lefèvre, Guillaume Melquiond, Nathalie Revol, Damien Stehlé, and Serge Torres. *Handbook*  
783 *of Floating-Point Arithmetic*. Birkhäuser, 2010. doi:10.1007/978-0-8176-4705-6.
- 784 **31** César Muñoz. Rapid prototyping in PVS. Technical Report CR-2003-212418, NASA, 2003.
- 785 **32** Henri Poincaré. *La science et l'hypothèse*. Flammarion, Paris, 1902.
- 786 **33** Pierre Roux. Formal Proofs of Rounding Error Bounds - With Application to an Automatic  
787 Positive Definiteness Check. *J. Autom. Reasoning*, 57(2):135–156, 2016. doi:10.1007/  
788 s10817-015-9339-z.
- 789 **34** Siegfried M. Rump. Verification of positive definiteness. *BIT Numerical Mathematics*, 46:433–  
790 452, 2006.
- 791 **35** Siegfried M. Rump. Verification methods: Rigorous results using floating-point arithmetic.  
792 *Acta Numerica*, 19:287–449, 2010.
- 793 **36** Siegfried M Rump, Takeshi Ogita, Yusuke Morikura, and Shin'ichi Oishi. Interval arithmetic  
794 with fixed rounding mode. *Nonlinear Theory and Its Applications, IEICE*, 7(3):362–373, 2016.
- 795 **37** Siegfried M. Rump, Takeshi Ogita, and Shin'ichi Oishi. Accurate floating-point summation  
796 part I: faithful rounding. *SIAM J. Scientific Computing*, 31(1):189–224, 2008. doi:10.1137/  
797 050645671.
- 798 **38** Siegfried M. Rump, Paul Zimmermann, Sylvie Boldo, and Guillaume Melquiond. Computing  
799 predecessor and successor in rounding to nearest. *BIT Numerical Mathematics*, 49(2):419–431,  
800 Jun 2009. doi:10.1007/s10543-009-0218-z.
- 801 **39** Arnaud Spiwack. *Verified Computing in Homological Algebra. (Calculs vérifiés en algèbre*  
802 *homologique)*. PhD thesis, École Polytechnique, Palaiseau, France, 2011. URL: <https://tel.archives-ouvertes.fr/pastel-00605836>.  
803