

# Formal correctness of comparison algorithms between binary64 and decimal64 floating-point numbers

Arthur Blot, Jean-Michel Muller, Laurent Théry

► **To cite this version:**

Arthur Blot, Jean-Michel Muller, Laurent Théry. Formal correctness of comparison algorithms between binary64 and decimal64 floating-point numbers. Numerical Software Verification, Jul 2017, Heidelberg, Germany. hal-01512294

**HAL Id: hal-01512294**

**<https://hal.archives-ouvertes.fr/hal-01512294>**

Submitted on 22 Apr 2017

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Formal correctness of comparison algorithms between binary64 and decimal64 floating-point numbers

Arthur Blot<sup>1</sup>, Jean-Michel Muller<sup>1</sup>, and Laurent Théry<sup>2</sup>

<sup>1</sup> École Normale Supérieure de Lyon, Lyon, France,  
arthur.blot@ens-lyon.org    jean-michel.muller@ens-lyon.fr

<sup>2</sup> INRIA Sophia Antipolis - Méditerranée, France,  
Laurent.Theiry@inria.fr

**Abstract.** We present a full COQ formalisation of the correctness of some comparison algorithms between binary64 and decimal64 floating-point numbers.

## 1 Introduction

Both binary and decimal formats are defined in the IEEE 754-2008 standard. However, no operation is defined to compare them, and the “naïve” approach of converting a decimal floating point number to the closest binary one before comparing it to another binary FP number can lead to inconsistencies. Let us take for example these two FP numbers:

$$\begin{array}{ll} x := 7205759403792794/2^{56} & \text{(in binary64 format) } > 1/10 \\ y := 1/10 & \text{(in decimal64 format)} \end{array}$$

If we convert  $y$  to the binary64 format, we get  $x$ . This means the native approach would give that  $x \leq y$ , but actually  $x > y$ .

The paper [4] by Brisebarre et al. details and proves the correctness of several algorithms which can compare these numbers in a both efficient and entirely reliable way. The initial problem being simple, it would be nice if the algorithms and their correctness proofs were elementary too. This is unfortunately not the case. In particular, the correctness of some of the algorithms requires some non-trivial computation that cannot be performed by hand. The contribution of this paper is to propose a formalisation that carefully checks all the aspects of the correctness proofs.

We use the COQ proof assistant [9] for our formalisation. The paper is organised as follows. In a first section we detail the formal setting in which we have made our formalisation. Then, we present the algorithms and their proof of correctness.

## 2 Formal Setting

COQ is a generic system that lets us define new objects and build formal proofs of their properties. There are two main ways of adding new knowledge into the system. We can work conservatively defining a new notion in term of already defined ones or we can work axiomatically by providing the characteristics properties of the new notion. In our formalisation, the two approaches have been used. We now illustrate this with some elementary examples.

Positive numbers are represented in COQ in a conservative way by a recursive datastructure `positive` that mimics binary numbers. It is composed of 3 constructors: `xI` (a number that starts with a 1), `x0` (a number that starts with a 0), `xH` (the number 1). If we take the little endian convention, the binary number 1101 is represented as `(xI (x0 (xI xH)))`. Thanks to this recursive datastructure, it is then possible to equip positive numbers with usual functions. For example, we give below the definition of the function that computes the successor function. It is defined recursively by a discussion on the structure of its argument `a`.

```
Function Psucc a :=
  match a with
  | xI a1 => x0 (Psucc a1) | x0 a1 => xI a1 | xH => x0 xH
  end.
```

Relative numbers are then defined on top of positive numbers using these constructor `Zero`, `Zpos` for positive numbers and `Zneg` for negative numbers. Usual notations can be associated with these numbers, so that `(Zneg (xI (x0 (xI xH))))` is printed and parsed back as `-13`.

Similarly to the definition of the `Psucc` function, addition (+) is defined for relative numbers. Having this function as an algorithm inside COQ makes it possible to compute with relative numbers like in any programming language. Thus, `-13 + 4` evaluates to `9`. Furthermore, it is possible to reason about this function and prove that the distributive property holds.

```
Lemma add_assoc m n p, m + (n + p) = (m + n) + p.
```

Having a system where one can mix computation and proofs was mandatory for the formalisation of the properties we were aiming at. Coming back to this formalisation, we are using the FLOCC library [2] to reason about floating point numbers. These numbers are parametrised by a radix  $\beta$ . They are defined as a record that contains two fields.

```
Record float (beta : radix) := Float { Fnum : Z; Fexp : Z }.
```

For example, the object `{|Fnum := 7205759403792794; Fexp := -56|}` is a floating-point number whose mantissa is 7205759403792794 and exponent is `-56`.

In our formalisation, we prove the correctness of the comparison algorithms with respect to the comparison over real numbers. The real numbers are axiomatised in COQ. An abstract type `R` is defined. A set of properties is assumed that makes it a complete, totally ordered, Archimedean field. On top of this

basic field, the standard COQ real library and the COQUELICOT[1] library build the usual real analysis. This contains standard notions like series and limits from which usual functions such as exponential, sine and cosine are defined. Thanks to these libraries, it is possible to give an interpretation of a floating point number as a real number.

**Definition** `F2R (f : float beta) := (Fnum f) * beta ^ (Fexp f).`

The floating point number `f` represents the real numbers  $m * \beta^e$  where  $m$  is the mantissa,  $e$  the exponent and  $\beta$  the base. Now the problem is simple, we want to build an algorithm `comp` that takes two floating point numbers `f1` and `f2` in base 2 and in base 10 respectively and returns their comparison. We use the standard COQ type `comparison` for the returned type of the comparison function `comp`. This type contains three elements `Lt`, `Eq`, `Gt` that represent the three possible values of a comparison. Our goal is then to prove that, under some conditions `Cond` that we will explain later, the following theorem holds

**Lemma** `comp_correct (f1 : float 2) (f2 : float 10) :`  
`Cond`  $\rightarrow$  `comp f1 f2 = Lt`  $\leftrightarrow$  `F2R f1 < F2R f2`  
 $\wedge$  `comp f1 f2 = Eq`  $\leftrightarrow$  `F2R f1 = F2R f2`  
 $\wedge$  `comp f1 f2 = Gt`  $\leftrightarrow$  `F2R f1 > F2R f2.`

This lemma relates the result of an evaluation of the algorithm `comp` with a relation between two real numbers. When proving an algorithm `eq` that performs an equality test, the returned type is simply a boolean so the correctness theorem looks like

**Lemma** `eq_correct (f1 : float 2) (f2 : float 10) :`  
`Cond`  $\rightarrow$  `eq f1 f2 = true`  $\leftrightarrow$  `F2R f1 = F2R f2.`

Note that as real numbers are axiomatised in COQ, we can not directly compute with them. Fortunately, there exist some tools that address this issue. One tool that we are using in this formalisation is the `interval` tactic [8]. It tries to automatically solve goals that contain real expressions by computing an approximation using standard techniques of interval arithmetic. For example, the following error bound of an approximation of the exponential function between 0 and 1 can be proved automatically.

**Lemma** `bound_exp_taylor2 x :`  
 $0 \leq x \leq 1 \rightarrow$  `Rabs ((1 + x + x ^ 2 / 2) - exp x)  $\leq$  22 / 100.`

Another technique to approximate real numbers are continued fractions. We are using them in our formalisation. Following closely the material that is presented by Khinchin [7], we have developed a library to compute and reason about continued fractions inside COQ. The main basic result that is needed here is that continued fractions represent the best rational approximations one can get. The formal statement is the following.

**Lemma** `halton_min (n : nat) (p q : Z) (r : R) :`  
 $0 < q < q[r]_{n.+1} \rightarrow$  `t[r]_n  $\leq$  Rabs (q * r - p).`

Approximating a real number  $r$  by continued fractions can be seen as an iterative process. At rank  $n$ , we get an approximation, called convergent,  $s_n = p_n/q_n$ . One way of quantifying the quality of an approximation  $p/q$  is by the value,  $|qr - p|$ . Applied to  $s_n$ , this gives  $\theta_n = |q_n r - p_n|$ , written  $\mathfrak{t}[\mathbf{r}]_n$  in our formalisation. The previous lemma `halton_min` simply states that the convergent at rank  $n$  is the best approximation for all the rational numbers  $p/q$  with  $q$  less than  $q_{n+1}$ .

### 3 Formalisation

We want to compare a binary floating-point number  $x_2$  and a decimal floating-point number  $x_{10}$ . This is reflected in our formalisation by introducing two variables of type `float`.

**Variable** `x2` : `float 2`.

**Variable** `x10` : `float 10`.

Furthermore,  $x_2$  belongs to the binary64 format and  $x_{10}$  to the decimal64 format. This means that we have:

$$\begin{aligned} x_2 &= M_2 \cdot 2^{e_2-52} \\ x_{10} &= M_{10} \cdot 10^{e_{10}-15} \end{aligned}$$

where:

$$\begin{aligned} -1022 \leq e_2 \leq +1023 \quad |M_2| \leq 2^{53} - 1 \\ -383 \leq e_{10} \leq +384 \quad |M_{10}| \leq 10^{16} - 1 \end{aligned}$$

Without loss of generality, we can also suppose that  $x_2 > 0$  and  $x_{10} > 0$ . If they are of opposite sign, the comparison is trivial and if they are both negative, taking the opposite of both numbers brings us back to the case we consider.

In our formalisation, as  $x_2 > 0$ , we shift the exponent in such a way that we have  $2^{52} \leq M_2$ . This corresponds to our `M2_bound` assumption.

**Definition** `M2` := `Fnum x2`.

**Hypothesis** `M2_bound` :  $2^{52} \leq M2 \leq 2^{53} - 1$ .

**Definition** `e2` := `Fexp x2 + 52`.

**Hypothesis** `e2_bound` :  $-1074 \leq e2 \leq 1023$ .

Similarly for  $x_{10}$  but without any shifting, we get

**Definition** `M10` := `Fnum x10`.

**Hypothesis** `M10_bound` :  $1 \leq M10 \leq 10^{16} - 1$ .

**Definition** `e10` := `Fexp x10 + 15`.

**Hypothesis** `e10_bound` :  $-383 \leq e10 \leq 384$ .

The four hypothesis `M2_bound`, `e2_bound`, `M10_bound` and `e10_bound` are implicitly conditions of all the theorems we present in the following.

Let us start deriving new facts from these definitions. First,  $M_{10}$  can also be “normalised”. If we consider  $\nu = 53 - \lfloor \log_2(M_{10}) \rfloor$ , we have that:

$$0 \leq \nu \leq 53$$

$$2^{53} \leq 2^\nu M_{10} \leq 2^{54} - 1$$

These correspond to the theorems `v_bound` and `norm2_bound` of our formalisation.

Now, we express  $x_2$  and  $x_{10}$  with the largest common power of 2 possible. For this, we define  $m = M_2$ ,  $n = M_{10}2^\nu$ ,  $h = \nu + e_2 - e_{10} - 37$  and  $g = e_{10} - 15$ . This gives us the following equalities (theorems `x2_eq` and `x10_eq` respectively).

$$x_2 = m \cdot 2^h \cdot 2^{g-\nu} \qquad x_{10} = n \cdot 5^g \cdot 2^{g-\nu}$$

The initial problem of comparing  $x_2$  and  $x_{10}$  reduces to comparing  $m \cdot 2^h$  with  $n \cdot 5^g$ . We also get bounds over  $m$ ,  $h$ ,  $n$  and  $g$ .

$$2^{52} \leq m \leq 2^{53} - 1$$

$$2^{53} \leq n \leq 2^{54} - 1$$

$$-398 \leq g \leq 369$$

$$-787 \leq h \leq 716$$

These are easily proven in COQ and correspond to the theorems `m_bound`, `n_bound`, `g_bound` and `h_bound`.

## 4 Handling “simple” cases

The first step of the algorithm consists in checking if the result can be determined by looking at the exponents  $g$  and  $h$  only. In order to compare  $5^g$  and  $2^h$ , a function  $\varphi$  is introduced. Its definition is the following.

$$\varphi(h) = \lfloor h \cdot \log_5 2 \rfloor$$

It is relatively easy to check that one can determine the result of the comparison between  $x_2$  and  $x_{10}$  if  $g$  and  $\varphi(h)$  differ. Formally, this means that the following two properties hold:

$$\text{If } g < \varphi(h) \text{ then } x_2 > x_{10}$$

$$\text{If } g > \varphi(h) \text{ then } x_2 < x_{10}$$

They correspond to the theorems `easycomp_lt` and `easycomp_gt` in our formalisation.

In order to get an algorithm, one still needs to provide a way to compute  $\varphi(h)$ . One way to do this is to use the integer value  $\lfloor 2^k \cdot \log_5 2 \rfloor$  for a sufficient large value of  $k$ . In our case,  $k = 19$  is enough, the following property holds.

$$\varphi(h) = \lfloor h \cdot \lfloor 2^{19} \cdot \log_5 2 \rfloor \cdot 2^{-19} \rfloor \qquad \text{for } |h| \leq 1831$$

This corresponds to the theorem `prop2_1`. Unfortunately, such a theorem is outside the reach of the `interval` tactic which does not handle the floor and to the nearest functions automatically.

In order to overcome this problem, we use the characteristic properties of the rounding functions ( $x$  is a real number,  $y$  an integer)

$$\begin{array}{ll} \lfloor x \rfloor = y & \text{iff} & y \leq x < y + 1 \\ \lceil x \rceil = y & \text{if} & |x - y| < 1/2 \end{array}$$

For example, in order to prove that  $\lceil 2^{19} \cdot \log_5 2 \rceil = 225799$ , we reduce the problem to showing that  $|2^{19} \cdot \log_5 2 - 225799| < 1/2$  which is solved automatically by the `interval` tactic. The proof of the `prop2_1` theorem is done in a similar way. We first generate the 3663 subgoals that represents all the possible values for  $h$ . For each of the subgoal, we guess what the value  $z$  for  $\lfloor h \cdot \log_5 2 \rfloor$  is. We then need to prove that  $\varphi(h) = z$  and  $\lceil h \cdot 225799 \cdot 2^{-19} \rceil = z$ . For both of these goals, we use the characteristic property of the floor function. Showing that  $z \leq h \cdot \log_5 2 < z + 1$  and  $z \leq h \cdot 225799 \cdot 2^{-19} < z + 1$  is done automatically by the `interval` tactic. This brute-force method is not the most elegant approach but it works. It takes COQ 40 minutes to check the theorem `prop2_1`.

Once all the properties of the function  $\varphi$  are proved, it is easy to write down an actual algorithm for the simple cases

```
Definition easycomp : comparison :=
  let h := v + e2 - e10 - 37 in
  let g := e10 - 15 in
  let phih := (225799 * h) / (2 ^ 19) in
  if g < phih then Gt
  else if g > phih then Lt
  else Eq.
```

and derive its associated theorem of correctness

```
Theorem easycomp_correct:
  (easycomp = Lt → F2R x2 < F2R x10) ∧
  (easycomp = Gt → F2R x2 > F2R x10) ∧
  (easycomp = Eq ↔ g = phi h).
```

Note that we are using the ideal integer arithmetic of COQ in a more realistic programming language we would have to deal with bounded arithmetic. Nevertheless, we can run our algorithm on the example given in the introduction and get the expected result.

```
Compute easycomp {| Fnum := 7205759403792794; Fexp := -56 |}
              {| Fnum :=
              1; Fexp := -1 |}.
= Eq
: comparison
```

Our correctness theorem tells us that all we can conclude is that  $g = \varphi(h)$ .

## 5 Exact testing

The algorithm in the previous section covers the cases where  $g$  differs from  $\varphi(h)$ . In this section, we assume that  $g = \varphi(h)$ .

### 5.1 Finding the needed precision

We define:

$$f(h) = 5^{\varphi(h)} \cdot 2^{-h} = 2^{\lfloor h \log_5 2 \rfloor - \log_2 5^{-h}}$$

which verifies:

$$f(h) \cdot n > m \Rightarrow x_{10} > x_2$$

$$f(h) \cdot n < m \Rightarrow x_{10} < x_2$$

$$f(h) \cdot n = m \Rightarrow x_{10} = x_2$$

Now, the idea is to get a precise enough approximation of  $f(h) \cdot n$  in order to retrieve the exact comparison. Our goal is therefore to find a lower bound  $\eta$  over the smallest non-zero value of the error  $|m/n - f(h)|$  in the range where

$$\begin{aligned} 2^{52} &\leq m \leq 2^{53} - 1 \\ 2^{53} &\leq n \leq 2^{54} - 1 \\ -787 &\leq h \leq 716 \end{aligned}$$

We can also add an extra constraint that directly comes from the definition of  $n$ .

$$\text{If } 10^{16} \leq n \text{ then even}(n)$$

A good candidate for this lower-bound, as proposed in [4], is  $\eta = 2^{-113.7}$ . In order to formally check that this lower-bound holds, we split in four the interval of  $h$ . For  $0 \leq h \leq 53$ , it is easily to prove that the lower bound is greater than  $2^{-107}$ . This corresponds to the theorem `d_h_easy1`. Similarly, the lower bound is greater than  $2^{-108}$  for  $-53 \leq h \leq 0$  (`d_h_easy2`). For the two remaining intervals  $54 \leq h \leq 716$  and  $-787 \leq h \leq -54$ , we make use of our library for continued fractions. The key idea is that in order to find the lower bound for a given  $h$  it is not necessary to enumerate all the rationals  $m/n$  but only the convergents  $s_n = p_n/q_n$  that approximate  $f(h)$  and whose denominator are bounded by  $2^{53} - 1$ . The property that justifies this computation is the following:

$$\text{If } q_n \leq n < q_{n+1} \text{ and } |f(h) - m/n| < 1/2n^2 \text{ Then } m/n = p_n/q_n$$

This is the `conv_2q2` theorem of our library that directly follows the proof given in [7]. In our case, as  $\eta = 2^{-113.7}$  and  $2^{53} \leq n \leq 2^{54} - 1$ , we are sure that  $\eta < 1/2n^2$ , so our theorem applies.

For the first interval  $54 \leq h \leq 716$ , the lower bound  $\epsilon_p/\epsilon_q$  is reached for  $h = 612$ ,  $m = 3521275406171921$  and  $n = 8870461176410409$  with

$$\begin{aligned} \epsilon_p &= 95837236471514977157488559853369253873503974784033193441759367158 \\ &879170025294997936467046601695051723045003185390596549988774013061 \\ &69390840117819918935869930600601291 \end{aligned}$$



and

```
 $\epsilon_q = 15076604622712586378413085535543158310234113291354064071380852407$   
218620654696227859485444413405009205300052717281615404952311527789  
630428217161580997965373300467263261324839949270053015170403319309  
9264
```

This corresponds to the theorem `pos_correct` of our formalisation. It is proved by a double enumeration. First, we enumerate all the  $h$  from 54 to 716, and for each of them, we enumerate all the convergents of  $f(h)$  whose denominator is smaller than  $2^{53} - 1$  and check that the error is bigger than  $\epsilon_p/\epsilon_q$ . The entire proof takes 80 seconds to be checked by COQ.

For the second interval  $-787 \leq h \leq -54$ , a similar technique is used. This time the lower bound  $\epsilon'_p/\epsilon'_q$  is reached for  $h = 275$ ,  $m = 4988915232824583$  and  $n = 12364820988483254$  with

```
 $\epsilon'_p = 112529423171232400134835569054963071809061903761868569039360251397$ 
```

and

```
 $\epsilon'_q = 186045148427403750191543353736499076917034509405654972460757157204$   
3696050968719646334648132324218750
```

This corresponds to theorem `neg_correct` and takes 100 seconds to be validated by COQ.

The theorem `d_h_min` collects all the previous results and proves the smallest non-zero error is bigger than  $\eta = 2^{-113.7}$ . We thus know that if the error is less than  $\eta$ , the distance is 0 and we must be in the equality case.

## 5.2 Direct method

The first way to use the previous result is to precompute  $f$  in a table indexed by  $h$ . Then, assuming  $\mu$  approximates  $f(h) \cdot n$  with accuracy  $\epsilon < \eta/4$ , we have the following implications:

If $\mu > m + \epsilon \cdot 2^{54}$	Then $x_{10} > x_2$
If $\mu < m - \epsilon \cdot 2^{54}$	Then $x_{10} < x_2$
If $ m - \mu  \leq \epsilon \cdot 2^{54}$	Then $x_{10} = x_2$

This corresponds to the theorem `direct_method_correct` of our formalisation.

Since we have  $\eta = 2^{-113.7}$ , a 128-bit computed value of  $f$  is enough to make the exact comparison. In order to build the algorithm, we first build the table with 1504 entries that contains  $\lceil f(h) \cdot 2^{127} \rceil$ .

```
Definition f_tbl h : Z :=  
  match h with  
  | -787 => 155090275327330814032632713427604519407  
  | -786 => 77545137663665407016316356713802259704
```

```

| -785 => 38772568831832703508158178356901129852
...
...
...
| 714 => 75715339914673581502256102241153698026
| 715 => 37857669957336790751128051120576849013
| 716 => 94644174893341976877820127801442122533
| _    => 0
end.

```

Proving that this table is correct (for each entry  $i$ , the value represents  $\lceil f(h) \cdot 2^{127} \rceil$ ) is done by brute and takes 40 minutes. The algorithm then just needs to look up in the table and compares with  $m * 2^{217}$ .

```

Definition direct_method_alg :=
  let v := n * (f_tbl h) - m * 2 ^ 127 in
  if v < - 2 ^ 57 then Gt
  else if v > 2 ^ 57 then Lt
  else Eq.

```

As we have  $2^{124} \leq f(h) \cdot 2^{127} < 2^{128}$ , our epsilon is equal to  $2^{-124}$ . This explains the 57 in the algorithm,  $57 = -124 + 127 + 54$ . The proof of correctness is a direct consequence of these observations.

```

Theorem direct_method_alg_correct :
  (direct_method_alg = Lt ↔ F2R x2 < F2R x10) ∧
  (direct_method_alg = Eq ↔ F2R x2 = F2R x10) ∧
  (direct_method_alg = Gt ↔ F2R x2 > F2R x10).

```

We have seen in the previous section with the execution of the `easycomp` algorithm that we have  $g = \varphi(h)$  for the example in the introduction. Now, running the `direct_method_alg` algorithm returns the correct value for the comparison.

```

Compute ineq_alg { | Fnum := 7205759403792794; Fexp := -56 | }
              { | Fnum :=                1; Fexp := -1 | }.
= Gt
: comparison

```

The main drawback of this algorithm is that the table may be too large for some realistic implementation. The next section explains how we can alleviate the problem by using a bipartite table.

### 5.3 Bipartite table method

In this method, we use:

$$q = \left\lfloor \frac{\varphi(h)}{16} + 1 \right\rfloor \quad r = 16q - \varphi(h)$$

such that  $f(h) = 5^{16q} \cdot 5^{-r} \cdot 2^{-h}$ . We can easily derive the bounds for  $q$  and  $r$  from the bound of  $\varphi(h)$ .

$$\begin{aligned} -21 &\leq q \leq 20 \\ 1 &\leq r \leq 16 \end{aligned}$$

In this method, we store  $5^{16q}$  and  $5^r$  after a phase of normalisation. For this, we define the function  $\psi$  as follows.

$$\psi(g) = \lfloor g \cdot \log_2 5 \rfloor$$

Similarly to what has been done for  $\varphi(h)$ , it is easy to compute  $\psi(q)$  and  $\psi(16q)$  using integer arithmetic:

$$\begin{aligned} \psi(g) &= \lfloor \lceil 2^{12} \cdot \log_2 5 \rceil \cdot g \cdot 2^{-12} \rfloor && \text{for } |g| \leq 204 \\ \psi(16q) &= \lfloor \lceil 2^{12} \cdot \log_2 5 \rceil \cdot q \cdot 2^{-8} \rfloor && \text{for } |q| \leq 32 \end{aligned}$$

These equalities correspond to the theorems `prop_2_2` and `prop_2_3` in our formalisation. They are proved by brute force in 5 minutes and 40 seconds respectively.

Now, we introduce the two functions  $\theta_1$  and  $\theta_2$  that are defined as follows.

$$\begin{aligned} \theta_1(q) &= 5^{16q} \cdot 2^{-\psi(16q)+127} \\ \theta_2(r) &= 5^r \cdot 2^{-\psi(r)+63} \end{aligned}$$

The following bounds are given by the theorems `theta1_bound` and `theta2_bound`.

$$\begin{aligned} 2^{127} &\leq \theta_1(q) < 2^{128} - 1 \\ 2^{63} &< \theta_2(r) < 2^{64} \end{aligned}$$

Furthermore, we have

$$f(h) = \frac{\theta_1(q)}{\theta_2(r)} 2^{-64-\sigma(h)} \quad \text{with } \sigma(h) = \psi(r) - \psi(16q) + h$$

If we define  $\Delta$  as

$$\Delta = \theta_1(q) \cdot n \cdot 2^{-64+8} - \theta_2(r) \cdot m \cdot 2^{8+\sigma(h)}$$

comparing  $x_2$  and  $x_{10}$  gives the same result as comparing 0 and  $\Delta$  (theorem `delta_ineq`). We also easily get that if  $x_2 \neq x_{10}$ , then  $|\Delta| \geq 2^{124}\eta$ . This allows us to derive an approximated version of  $\Delta$ :

$$\tilde{\Delta} = \lfloor \lceil \theta_1(q) \rceil \cdot n \cdot 2^{8-64} \rfloor - \theta_2(r) \cdot m \cdot 2^{8+\sigma(h)}$$

and proves that comparing 0 with  $\tilde{\Delta}$  gives the same result as comparing  $x_2$  with  $x_{10}$  (theorem `delta'_ineq`).

In order to build the algorithm, we first need to create the two tables for  $\lceil \theta_1(q) \rceil$  and  $\theta_2(r)$ . Our tables are given in Figure 1. The first one contains integers that fit in 128 bits and the second one in 64 bits. The algorithm first computes  $q$  and  $r$ , then reads the values of  $\lceil \theta_1(q) \rceil$  and  $\theta_2(r)$ , and then checks the sign of the computed  $\tilde{\Delta}$ .

```

Definition theta1_tbl q : Z :=
  match q with
  | -21 => 302910693998692996157485768413290076966
  | -20 => 336298426882534191759128470626028036789
  | -19 => 186683128335104582129005107785662008085
  | -18 => 207259907386686073192955235040171322419
  | -17 => 230104721262376436189351064420995165904
  | -16 => 255467559620444135892015707268715336457
  | -15 => 283625966735416996535885333662014114405
  | -14 => 314888078651228693933689466069052580905
  | -13 => 174797997549285651882438866782683340398
  | -12 => 194064761537588616893622436057812819408
  | -11 => 215455166527421378856590945602770070141
  | -10 => 239203286653190548679094257880939433815
  | -9 => 265568996408383549344794103276234313665
  | -8 => 294840814439182918143871451639708507103
  | -7 => 327339060789614187001318969682759915222
  | -6 => 181709681073901722637330951972001133589
  | -5 => 201738271725539733566868685312735302683
  | -4 => 223974474217780421055744228056844427813
  | -3 => 248661618204893321077691124073410420051
  | -2 => 276069853871622551497390234491081018099
  | -1 => 30649910817317771671669405430061836724
  | 0 => 170141183460469231731687303715884105728
  | 1 => 18889465931478580854784000000000000000000
  | 2 => 20971520000000000000000000000000000000000
  | 3 => 23283064365386962890625000000000000000000
  | 4 => 258493941422821148397315216271863391740
  | 5 => 286985925493722536125179818657774823687
  | 6 => 318618382226490455405776079553542361119
  | 7 => 176868732008334225927912486150152183217
  | 8 => 196363738611909062123830878199451025720
  | 9 => 218007543808417316859394750271862213031
  | 10 => 242036994678082392051126914580396990474
  | 11 => 268715044302683550071638623558009085183
  | 12 => 298333629248008269731638612618517353496
  | 13 => 331216864211123806751178713779234900611
  | 14 => 183862294395666818064937594201088633456
  | 15 => 204128152598478183127259193653345185578
  | 16 => 226627774989027955951103258828533066424
  | 17 => 251607373812388019852618613412845800985
  | 18 => 279340299571981831419774226402503504146
  | 19 => 310130032290502989833240994779765547114
  | 20 => 172156751238329846960951049916624692801
  | _ => 0
end.

Definition theta2_tbl r : Z :=
  match r with
  | 1 => 11529215046068469760
  | 2 => 14411518807585587200
  | 3 => 18014398509481984000
  | 4 => 11258999068426240000
  | 5 => 14073748835532800000
  | 6 => 17592186044416000000
  | 7 => 10995116277760000000
  | 8 => 13743895347200000000
  | 9 => 17179869184000000000
  | 10 => 10737418240000000000
  | 11 => 13421772800000000000
  | 12 => 16777216000000000000
  | 13 => 10485760000000000000
  | 14 => 13107200000000000000
  | 15 => 16384000000000000000
  | 16 => 10240000000000000000
  | _ => 0
end.

```

Fig. 1. Tables for  $\lceil \theta_1(q) \rceil$  and  $\theta_2(r)$

```

Definition ineq_alg : comparison :=
  let q := g / 16 + 1 in
  let r := 16 * q - g in
  let psir := (9511 * r) / 2 ^ 12 in
  let psiq := (9511 * q) / 2 ^ 8 in
  let s := psir - psiq + h in
  let a := ((theta1_tbl q) * (n * 2 ^ 8)) / (2 ^ 64) in
  let b := (theta2_tbl r) * (m * 2 ^ (8 + s)) in
  let D := a - b in (0 =?= D)

```

The value 9511 corresponds to  $\lceil 2^{12} \cdot \log_2 5 \rceil$  and  $=?$  to the comparison function for integer. Its associated correctness theorem is the following.

```
Theorem ineq_alg_correct :
  (ineq_alg = Lt  $\leftrightarrow$  F2R x2 < F2R x10)  $\wedge$ 
  (ineq_alg = Eq  $\leftrightarrow$  F2R x2 = F2R x10)  $\wedge$ 
  (ineq_alg = Gt  $\leftrightarrow$  F2R x2 > F2R x10).
```

Running the `ineq_alg` algorithm gives the same value as the one using the direct method.

```
Compute ineq_alg {| Fnum := 7205759403792794; Fexp := -56 |}
              {| Fnum :=
                  1; Fexp := -1 |}.
= Gt
: comparison
```

## 6 Equality case

If we only want to test the equality between  $x_2$  and  $x_{10}$ , an even simpler algorithm can be used. We already know that  $x_2 = x_{10}$  is equivalent to  $m \cdot 2^h = n \cdot 5^g$ . As 2 and 5 are relatively prime, only two situations can occur :

- either  $5^g \mid m$  and  $2^h \mid n$  with  $0 \leq g \leq 22$  and  $0 \leq h \leq 53$ ;
- or  $2^{-h} \mid m$  and  $5^{-g} \mid n$  with  $-22 \leq g \leq 0$  and  $-51 \leq h \leq 0$ .

The algorithm is a direct encoding of this property. It checks that  $5^g \cdot (n2^{-h}) = m$  if we are in the first case or that  $5^{-g} \cdot (m2^h) = n$  if we are in the second case.

```
Definition eq_alg : bool :=
  if (0  $\leq$  h)  $\&\&$  (h  $\leq$  53)  $\&\&$  (0  $\leq$  g)  $\&\&$ 
    (g  $\leq$  22)  $\&\&$  (n mod (2  $\wedge$  h) == 0) then
      let m' := 5  $\wedge$  g * (n / (2  $\wedge$  h)) in m' == m
    else if (h  $\geq$  -51)  $\&\&$  (-22  $\leq$  g)  $\&\&$ 
      (g  $\leq$  0)  $\&\&$  (m mod (2  $\wedge$  (-h)) = 0) then
        let n' := 5  $\wedge$  (-g) * (m / (2  $\wedge$  (-h))) in n' == n
      else
        false.
```

Its correctness theorem proves the equivalence between running the algorithm and testing the two values.

```
Theorem eq_alg_correct : eq_alg = true  $\leftrightarrow$  (F2R x2 = F2R x10).
```

On our favorite example, it returns the expected result.

```
Compute eq_alg {| Fnum := 7205759403792794; Fexp := -56 |}
              {| Fnum :=
                  1; Fexp := -1 |}.
= false
: comparison
```

## 7 Conclusion and Future Works

Our formalisation contains four algorithms for comparing a binary64 floating-point number and a decimal64 floating-point number and their proof of correctness. The code is available at:

<https://gitlab.com/artart78/compbindec>

It is composed of five files: `util.v` and `rfrac.v` for general results (most of which are about real numbers or continued fractions), `frac.v` for the result of the continued fraction problem solution, `compformula.v` for the computation of some formulas requiring tabulation of values, and `compbindec.v` for the main result. Altogether, this amounts to about 5000 lines of code that COQ checks in about 1 hour and 40 minutes. We have been using intensively SSREFLECT [6] for its set of tactics and syntax, the Flocq [2] library for manipulating floating-point numbers, and the Interval [8] tactic for the computation over the real numbers. We had to develop a dedicated library for continued fraction in order to tackle some aspects of the proof. This library is available in the `rfrac.v` file.

The main contribution of this work is a carefully check of the algorithms and the proofs presented in [4]. In particular, we have connected in a formal way the paper proofs with the computation that is required in order to get the accuracy at which  $f(h) \cdot n$  needs to be computed. During the formalisation, we have found some minor mistakes in the original paper and we have also departed at some places from what was presented in the original paper. For example, the original paper states that  $-20 \leq q \leq 21$  instead of  $-21 \leq q \leq 20$ . Fortunately, this mistake has no consequence. More values are tabulated in the actual implementation (the bounds are slightly relaxed). The original paper mentions the bound  $2^{-113.67}$  while our formalisation only makes use of  $2^{-113.7}$ . The statement of the `direct_method_correct` omits the scaling factor of  $f(h)$  in the assumption. It makes its proof easier and simplifies its application. Finally, the statements of the theorems `pos_correct` and `neg_correct` do not mention the extra condition (if  $h \geq 680$ , then  $\nu' = h + \varphi(h) - 971 > 0$  and  $2^{\nu'} \mid n$ ) that is present in the original paper. It does not change the lower-bound and we believe that this omission makes the statements more readable at the cost of a negligible extra computing.

The capability of computing expressions over the real numbers thanks to the `interval` tactic has been a key ingredient on this work. These computations are often hidden in the original paper. It is the case for example for the bounds for  $\theta_1$  and  $\theta_2$ . The bounds that are obtained without taking into account the bound on  $h$  are actually less strict than the ones specified in the paper. For example, only  $\theta_1 < 2^{128}$  could be proved, but it is actually important to get the  $2^{128} - 1$  bound in order to know we can embed the value into a 128-bit value. The  $2^{128}$  bound was obtained by noticing that  $16q \log_2 5 - \lfloor 16q \log_2 5 \rfloor < 1$ , but the  $2^{128} - 1$  bound requires that it is less than  $1 + \log_2(1 - \frac{1}{2^{128}}) \simeq 1 - \frac{1}{2^{128}}$ . We successfully manage to formalise all these proofs but clearly life would have been easier if the tactic `interval` would accommodate the floor and the ceiling functions. It could in particular greatly improve our checking time. 40 minutes

out of the hour and 40 minus of our checking time are spent by our brute-force method in trying to validate the integer equivalent for the function  $\varphi(h)$ .

There are several ways to extend this work. First, we could formalise the counting argument that is present in the original paper that quantifies the percentage of floating point numbers that are comparable using the first partial comparison method (`easycomp`). It would require to develop the notion of cardinal of sets of floating point numbers in the COQ standard library. Second, what we have proven are only algorithms, it would be very interesting to try to prove a realistic software implementation using a tool like WHY3 [5]. Finally a sequel of our reference [4] for this work has been written by the same authors in [3]. The general idea of the algorithm remains the same, but the generalisation of the result also requires a generalisation of all the intermediate results, and computation may become a lot harder since the numbers can then take up to 128 bits. Formalising it would be a real challenge.

## References

1. S. Boldo, C. Lelay, and G. Melquiond. Coquelicot: A user-friendly library of real analysis for Coq. *Mathematics in Computer Science*, 9(1):41–62, 2015.
2. S. Boldo and G. Melquiond. Flocq: A unified library for proving floating-point algorithms in coq. In E. Antelo, D. Hough, and P. Ienne, editors, *20th IEEE Symposium on Computer Arithmetic, ARITH 2011, Tübingen, Germany, 25-27 July 2011*, pages 243–252. IEEE Computer Society, 2011.
3. N. Brisebarre, C. Q. Lauter, M. Mezzarobba, and J. Muller. Comparison between binary and decimal floating-point numbers. *IEEE Trans. Computers*, 65(7):2032–2044, 2016.
4. N. Brisebarre, M. Mezzarobba, J. Muller, and C. Q. Lauter. Comparison between binary64 and decimal64 floating-point numbers. In A. Nannarelli, P. Seidel, and P. T. P. Tang, editors, *21st IEEE Symposium on Computer Arithmetic, ARITH 2013, Austin, TX, USA, April 7-10, 2013*, pages 145–152. IEEE Computer Society, 2013.
5. J.-C. Filliâtre and A. Paskevich. Why3 - where programs meet provers. In M. Felleisen and P. Gardner, editors, *ESOP*, volume 7792 of *Lecture Notes in Computer Science*, pages 125–128. Springer, 2013.
6. G. Gonthier and A. Mahboubi. A small scale reflection extension for the Coq system. Technical Report RR-6455, INRIA, 2008.
7. A. Khinchin and H. Eagle. *Continued Fractions*. Dover books on mathematics. Dover Publications, 1964.
8. E. Martin-Dorel and G. Melquiond. Proving Tight Bounds on Univariate Expressions with Elementary Functions in Coq. *Journal of Automated Reasoning*, 57(3):187–217, octobre 2016.
9. The Coq Development Team. *The Coq Proof Assistant Reference Manual – Version V8.6*, Dec. 2016. <http://coq.inria.fr>.