

Computing predecessor and successor in rounding to nearest

Siegfried Rump, Paul Zimmermann, Sylvie Boldo, Guillaume Melquiond

► **To cite this version:**

Siegfried Rump, Paul Zimmermann, Sylvie Boldo, Guillaume Melquiond. Computing predecessor and successor in rounding to nearest. BIT Numerical Mathematics, Springer Verlag, 2009, 49 (2), pp.419-431. <10.1007/s10543-009-0218-z>. <inria-00337537>

HAL Id: inria-00337537

<https://hal.inria.fr/inria-00337537>

Submitted on 7 Nov 2008

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.

Computing predecessor and successor in rounding to nearest

submitted for publication September 11, 2008

Siegfried M. Rump ^{*} Paul Zimmermann [†] Sylvie Boldo [‡] Guillaume Melquiond [§]

Abstract

We give simple and efficient methods to compute and/or estimate the predecessor and successor of a floating-point number using only floating-point operations in rounding to nearest. This may be used to simulate interval operations, in which case the quality in terms of the diameter of the result is significantly improved compared to existing approaches.

Keywords. Floating-point arithmetic, rounding to nearest, predecessor, successor, directed rounding
AMS subject classification (2000). 68-04, 68N30

1 Introduction and notation

Throughout the paper we assume a floating-point arithmetic according to the IEEE 754 and IEEE 854 arithmetic standards [3, 4] with rounding to nearest. Denote the set of (single or double precision) floating-point numbers by \mathbb{F} , including $-\infty$ and $+\infty$, and let $\text{fl} : \mathbb{R} \rightarrow \mathbb{F}$ denote rounding to nearest according to IEEE 754. This includes especially the default rounding mode, to nearest “ties to even” and the rounding to nearest “ties to away” — away from zero — defined in the revision of the IEEE 754 standard, IEEE 754-2008, that we follow in this paper. Define for single and double precision \mathbf{u} , the relative rounding error unit, and η , the smallest positive subnormal floating-point number:

| | single precision (binary32) | double precision (binary64) |
|--------------|---|---|
| \mathbf{u} | 2^{-24} | 2^{-53} |
| η | 2^{-149} | 2^{-1074} |

Let β be the radix used in this floating-point format. We require β to be even and greater than one. This includes especially 2, 10, and their powers. For a format of precision p “digits” in radix β , we define \mathbf{u} as half the distance between 1 and its successor, i.e.:

$$\mathbf{u} = \frac{\beta^{1-p}}{2}.$$

Then \mathbf{u} and η satisfy:

$$\forall \circ \in \{+, -, \times, \div\} \forall a, b \in \mathbb{F} \setminus \{\pm\infty\} : \text{fl}(a \circ b) = (a \circ b)(1 + \lambda) + \mu \tag{1}$$

^{*}Institute for Reliable Computing, Hamburg University of Technology, Schwarzenbergstraße 95, 21071 Hamburg, Germany, and Visiting Professor at Waseda University, Faculty of Science and Engineering, 3-4-1 Okubo, Shinjuku-ku, Tokyo 169-8555, Japan (rump@tu-harburg.de).

[†]Paul Zimmermann, Centre de Recherche INRIA Nancy - Grand Est, Équipe-projet CACAO, Bâtiment A, 615 rue du jardin botanique, F-54600 Villers-lès-Nancy, France (Paul.Zimmermann@loria.fr).

[‡]Sylvie Boldo, INRIA Saclay - Île-de-France, Parc Orsay Université - ZAC des Vignes, 4 rue Jacques Monod - Bâtiment N, F-91893 Orsay Cedex, France (Sylvie.Boldo@inria.fr).

[§]Guillaume Melquiond, Centre de recherche commun INRIA - Microsoft Research, 28 rue Jean Rostand, F-91893 Orsay Cedex, France (guillaume.melquiond@inria.fr).

with $|\lambda| \leq \mathbf{u}$ and $|\mu| \leq \eta/2$ and at least one of λ, μ is zero, provided $\text{fl}(a \circ b)$ is finite. Note that for addition and subtraction μ is always zero. An important property of the rounding is the monotonicity, that is

$$\forall x, y \in \mathbb{R} : x \leq y \Rightarrow \text{fl}(x) \leq \text{fl}(y). \quad (2)$$

In particular, this implies

$$x \in \mathbb{R}, f \in \mathbb{F}, x \leq f \Rightarrow \text{fl}(x) \leq f, \quad (3)$$

which means that rounding cannot “jump” over a floating-point number. The floating-point predecessor and successor of a *real* number $x \in \mathbb{R}$ are defined by

$$\text{pred}(x) := \max\{f \in \mathbb{F} : f < x\} \quad \text{and} \quad \text{succ}(x) := \min\{f \in \mathbb{F} : x < f\},$$

respectively, where, according to IEEE 754, $\pm\infty$ are considered to be floating-point numbers. For example, $\text{succ}(1) = 1 + 2\mathbf{u}$. Using (1) it is not difficult to see that for finite $c \geq 0 \in \mathbb{F}$, as long as $\text{succ}(c)$ is finite:

$$\min(c(1 - 2\mathbf{u}), c - \eta) \leq \text{pred}(c) \quad \text{and} \quad \text{succ}(c) \leq \max(c(1 + 2\mathbf{u}), c + \eta), \quad (4)$$

and similarly for $c < 0$. For $a, b \in \mathbb{F}$ and finite $c := \text{fl}(a \circ b)$ the monotonicity (2) implies

$$a \circ b \in [c_1, c_2] \quad \text{where} \quad c_1 := \text{fl}(\text{fl}(c - \text{fl}(2\mathbf{u}|c|)) - \eta) \quad \text{and} \quad c_2 := \text{fl}(\text{fl}(c + \text{fl}(2\mathbf{u}|c|)) + \eta). \quad (5)$$

(Note that similarly to (1), the above remains true if $a \circ b$ is replaced by any real y , for example $y = \sin x$, as long as c is the correct rounding of y .) This is the usual basis of interval libraries to emulate directed roundings using only rounding to nearest (see, for example, [5]). It is disadvantageous because for $1.5 \cdot 2^k \leq |a \circ b| < 2^{k+1}$ the interval $[c_1, c_2]$ is twice as wide as it needed to be, i.e., 4 *ulps* (units in the last place) instead of 2 *ulps*.

The IEEE 754-2008 standard [3] requires availability of a function `nextUp` — and similarly `nextDown` — where `nextUp(a)` is the least floating-point number in the format of a that compares greater than a .

The function `nextUp` thus computes the successor of a floating-point number. This just amounts to adding η if directed roundings, as requested by IEEE 754, are available. However, performing a directed computation may not be supported by the programming language in use, or it may depend on a rounding mode change, which usually involves a flush of the processor pipeline. So the function benefits from using the default rounding mode only, which is to nearest with ties to even as requested for binary formats, and recommended for decimal formats by IEEE 754-2008.

In [2] a corresponding algorithm is given by splitting the floating-point number in two parts, treating the second part as an integer and adding/subtracting 1 regarding possible carry. The algorithm does the job, but is slow. In [1] a corresponding routine is given assuming that the exponent range is unlimited and that a fused multiply and accumulate instruction is available, that is $a \cdot b + c$ with only one rounding.

The contributions of this paper are the following. Firstly we describe a simple and efficient routine (Algorithm 1) for any radix to compute an interval $[c_1, c_2]$ — with $c_1, c_2 \in \mathbb{F}$ — containing $a \circ b$ for all $a, b \in \mathbb{F}$ and $\circ \in \{+, -, \times, \div\}$ provided $\text{fl}(a \circ b)$ is finite (Theorem 1). We then focus on binary arithmetic (radix 2), and prove that if $a \circ b$ is not a floating-point number, the result of Algorithm 1 is always best possible¹ except a small range near underflow (Theorem 2). Finally we describe a slightly more complicated variant (Algorithm 2) that always returns the best possible interval in binary arithmetic, and compare its efficiency with the C99 `nextafter` implementation.

2 The results

We use the “unit in the first place” `ufp(x)` defined for $x \in \mathbb{R}$ by

$$\text{ufp}(0) := 0 \quad \text{and} \quad \text{ufp}(x) := \beta^{\lfloor \log_\beta |x| \rfloor} \quad \text{for } x \neq 0.$$

¹Since we don’t know if $c := \text{fl}(a \circ b)$ is smaller or larger than $a \circ b$, the best possible interval is $[\text{pred}(c), \text{succ}(c)]$ of 2 *ulps*, unless c is a power of the radix.

It denotes the weight of the most significant digit in the representation of x . Then

$$\forall 0 \neq x \in \mathbb{R} : \quad \text{ufp}(x) \leq |x| < \beta \text{ufp}(x). \quad (6)$$

This concept introduced by the first author proved to be useful in delicate analyses of the accuracy of floating-point algorithms [6]. The definition is independent of some floating-point format. Define

$$\bar{\mathbb{U}} := \{f \in \mathbb{F} : |f| < \frac{\beta}{2} \mathbf{u}^{-1} \eta\}. \quad (7)$$

For example, in IEEE 754-2008 `binary64` format (double precision), $\bar{\mathbb{U}} = \{f \in \mathbb{F} : |f| < 2^{-1021}\}$. Note that $\frac{1}{2} \mathbf{u}^{-1} \eta$ is the smallest positive *normal* floating-point number. For positive $c \in \mathbb{F}$ such that $\text{succ}(c)$ is finite the following properties are easily verified (see also [6]):

$$\text{if } c \in \bar{\mathbb{U}} : \quad \text{pred}(c) = c - \eta, \quad \text{succ}(c) = c + \eta, \quad (8)$$

$$\text{if } c \notin \bar{\mathbb{U}}, c \neq \beta^k : \quad \text{pred}(c) = c - 2\mathbf{u} \text{ufp}(c), \quad \text{succ}(c) = c + 2\mathbf{u} \text{ufp}(c), \quad (9)$$

$$\text{if } c \notin \bar{\mathbb{U}}, c = \beta^k : \quad \text{pred}(c) = c - \frac{2}{\beta} \mathbf{u} \text{ufp}(c), \quad \text{succ}(c) = c + 2\mathbf{u} \text{ufp}(c). \quad (10)$$

Moreover, define for $c \in \mathbb{F} \setminus \{\pm\infty\}$ with $\text{pred}(c)$ and $\text{succ}(c)$ finite:

$$\mathcal{M}^-(c) := \frac{1}{2}(\text{pred}(c) + c) \quad \text{and} \quad \mathcal{M}^+(c) := \frac{1}{2}(c + \text{succ}(c)). \quad (11)$$

It follows for $c \in \mathbb{F}$, $x \in \mathbb{R}$,

$$\begin{aligned} x < \mathcal{M}^-(c) &\Rightarrow \text{fl}(x) \leq \text{pred}(c) & \text{and} & & \mathcal{M}^+(c) < x &\Rightarrow \text{succ}(c) \leq \text{fl}(x) \\ \mathcal{M}^-(c) < x &\Rightarrow c \leq \text{fl}(x) & \text{and} & & x < \mathcal{M}^+(c) &\Rightarrow \text{fl}(x) \leq c. \end{aligned} \quad (12)$$

2.1 General radix β

For $c \in \mathbb{F}$ consider the following algorithm.

Algorithm 1 *Bounds for predecessor and successor of finite $c \in \mathbb{F}$ in rounding to nearest*

$$\begin{aligned} e &= \text{fl}(\text{fl}(\phi|c|) + \eta) & \% \phi &= \mathbf{u} \left(1 + \frac{4}{\beta} \mathbf{u}\right) = \text{succ}(\mathbf{u}) \\ \text{cinf} &= \text{fl}(c - e) \\ \text{csup} &= \text{fl}(c + e) \end{aligned}$$

Note that we need a reasonable floating-point format to ensure that ϕ is a floating-point number. More precisely, we need $\eta \leq 2\mathbf{u}^2$ (which in turn implies $\eta \leq \frac{4}{\beta} \mathbf{u}^2$, or equivalently $2\mathbf{u} \notin \bar{\mathbb{U}}$, because η is a power of the radix). This does happen on any real-life floating-point format (all IEEE formats included). If not, then the last bit of ϕ is lost due to underflow.

The following results have been formally checked using the Coq automatic proof checker and a previous formalization of floating-point numbers [7, 8].

Theorem 1 *Let finite $c \in \mathbb{F}$ be given, and assume $\mathbf{u} \leq \frac{\beta^{-2}}{2}$. Let `cinf` and `csup` be the quantities computed by Algorithm 1. Then*

$$\text{cinf} \leq \text{pred}(c) \quad \text{and} \quad \text{succ}(c) \leq \text{csup}. \quad (13)$$

REMARK. The technical assumption $\mathbf{u} \leq \frac{\beta^{-2}}{2}$ is satisfied by any practical implementation of floating-point arithmetic: it means that we have at least 3 digits of precision.

PROOF OF THEOREM 1. Since $\mathbb{F} = -\mathbb{F}$ and $\text{fl}(-x) = -\text{fl}(x)$ for $x \in \mathbb{R}$, we may assume without loss of generality $c \geq 0$. If $c \in \overline{\mathbb{U}}$, then $e \geq \eta$ and (8) imply (13).

It remains to prove (13) for $0 \leq c \notin \overline{\mathbb{U}}$. One verifies (13) for c being the largest positive floating-point number, hence we assume without loss of generality that $\text{succ}(c)$ is finite.

We first prove that for e as computed in Algorithm 1 we have

$$e > \mathbf{u} \text{ufp}(c). \quad (14)$$

First note that $0 \leq c \notin \overline{\mathbb{U}}$ means $c \geq \frac{\beta}{2} \mathbf{u}^{-1} \eta$. Also note that $\mathbf{u} \text{ufp}(c) \in \mathbb{F}$. By (2), (6), (9) and (10),

$$c' := \text{fl}(\phi c) = \text{fl}(\text{succ}(\mathbf{u}) c) \geq \text{fl}(\text{succ}(\mathbf{u}) \text{ufp}(c)). \quad (15)$$

If $\text{ufp}(c) \text{succ}(\mathbf{u})$ is normal, i.e., $\text{ufp}(c) \text{succ}(\mathbf{u}) \geq \frac{1}{2} \mathbf{u}^{-1} \eta$, then

$$e = \text{fl}(c' + \eta) \geq c' \geq \text{fl}(\text{ufp}(c) \text{succ}(\mathbf{u})) = \text{ufp}(c) \text{succ}(\mathbf{u}) > \mathbf{u} \text{ufp}(c),$$

which proves (14). If $\text{ufp}(c) \text{succ}(\mathbf{u}) \in \overline{\mathbb{U}}$, i.e., $\text{ufp}(c) \text{succ}(\mathbf{u}) < \frac{\beta}{2} \mathbf{u}^{-1} \eta$, then (15), (2) and (6) imply

$$c' \geq \text{fl}(\text{ufp}(c) \text{succ}(\mathbf{u})) \geq \text{fl}(\mathbf{u} \text{ufp}(c)) = \mathbf{u} \text{ufp}(c). \quad (16)$$

We split the remaining in two cases. First, suppose $c' = \text{fl}(\phi c) < \frac{\beta}{2} \mathbf{u}^{-1} \eta$. Then $c' \in \overline{\mathbb{U}}$, and (8), (16) yield (14) as $e = \text{fl}(c' + \eta) = c' + \eta > \mathbf{u} \text{ufp}(c)$.

The second and last case corresponds to $\text{ufp}(c) \text{succ}(\mathbf{u}) < \frac{\beta}{2} \mathbf{u}^{-1} \eta$ and $c' \geq \frac{\beta}{2} \mathbf{u}^{-1} \eta$. Then

$$c' = \text{fl}(\phi c) \geq \frac{\beta}{2} \mathbf{u}^{-1} \eta > \text{ufp}(c) \text{succ}(\mathbf{u}) > \mathbf{u} \text{ufp}(c),$$

and finally $e = \text{fl}(c' + \eta) \geq c' > \mathbf{u} \text{ufp}(c)$, so that the proof of (14) is finished.

We now prove Theorem 1 for $c \notin \overline{\mathbb{U}}$. By (9) and (10) we know

$$c - \mathbf{u} \text{ufp}(c) \leq \mathcal{M}^-(c) \quad \text{and} \quad \mathcal{M}^+(c) = c + \mathbf{u} \text{ufp}(c),$$

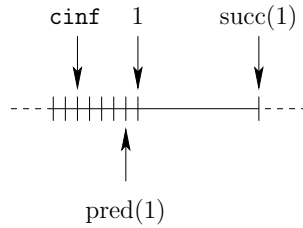
so (14) and (12) yield

$$c - e < \mathcal{M}^-(c) \quad \text{and} \quad \mathcal{M}^+(c) < c + e \quad (17)$$

and prove (13). The theorem is proved. \square

Note that the inequalities are far from sharp in a radix different from 2. The worst case in radix 10 corresponds to computing the predecessor of a power of the radix. For example, let $c = 1$ as in Figure 1, then $e = \phi$ and $\text{cinf} = \text{fl}(1 - \phi) = 1 - \mathbf{u}$ while $\text{pred}(c) = 1 - \frac{\mathbf{u}}{5}$. The computed cinf is equal to $\text{pred}(\text{pred}(\text{pred}(\text{pred}(\text{pred}(c))))))$.

Figure 1: Worst case for computing the predecessor in radix 10.



REMARK. Theorem 1 remains valid if e is computed with a fused multiply-add (fma) in Algorithm 1. Let $e' = \text{fl}(\phi|c| + \eta)$ the value computed with a fma, which might differ from $e = \text{fl}(\text{fl}(\phi|c|) + \eta)$. Either $\mathbf{u} \text{succ}(c)$ is normal, then $e' \geq \text{fl}(\phi|c|) = c' > \mathbf{u} \text{ufp}(c)$ as demonstrated in the proof of Theorem 1; or $\mathbf{u} \text{succ}(c)$ is subnormal, then $e' \geq \text{fl}(\mathbf{u}c + \eta) \geq \text{fl}(\mathbf{u} \text{ufp}(c) + \eta) = \mathbf{u} \text{ufp}(c) + \eta$ since $\mathbf{u} \text{ufp}(c) \in \mathbb{F}$ and is subnormal in that case. Thus in both cases $e' > \mathbf{u} \text{ufp}(c)$ remains valid.

2.2 Binary arithmetic

In this section we assume $\beta = 2$, thus \mathbf{u} is a power of two, and $\phi = \mathbf{u}(1 + 2\mathbf{u}) = \text{succ}(\mathbf{u})$. For a general radix we proved that the inequalities in (13) are satisfied for all finite floating-point numbers. Here we show that equality holds for binary arithmetic except a small range near the underflow threshold.

Theorem 2 *Assume binary arithmetic, i.e., $\beta = 2$, with at least 4 bits of precision, i.e., $\mathbf{u} \leq \frac{1}{16}$. Assume the rounding is “ties to even”, or “ties to away” as defined by IEEE 754-2008. Then for all finite $c \in \mathbb{F}$ with $|c| \notin [\frac{1}{2}, 2]\mathbf{u}^{-1}\eta$ the quantities `cinf` and `csup` computed by Algorithm 1 satisfy*

$$\text{cinf} = \text{pred}(c) \quad \text{and} \quad \text{succ}(c) = \text{csup}. \quad (18)$$

REMARK 1. For IEEE 754 `binary64` format the excluded range $[\frac{1}{2}, 2]\mathbf{u}^{-1}\eta$ is $[2^{-1022}, 2^{-1020}]$, which corresponds to two binades around the subnormal threshold. In this range Algorithm 1 returns `cinf` = `pred(pred(c))` and `csup` = `succ(succ(c))`, i.e., it is one bit off.

REMARK 2. It is easy to see from the proof that when $\mathbf{u} \leq \frac{1}{32}$, the assertions remain true for any tie-breaking rule in rounding to nearest.

PROOF OF THEOREM 2. Without loss of generality we assume $0 \leq c \in \mathbb{F}$. If $c < \frac{1}{2}\mathbf{u}^{-1}\eta$, then $c \leq \frac{1}{2}\mathbf{u}^{-1}\eta - \eta$ by (8), so that

$$\phi c \leq \frac{1}{2}(1 + 2\mathbf{u})\eta - \mathbf{u}(1 + 2\mathbf{u})\eta < \frac{1}{2}\eta$$

implies $\text{fl}(\phi c) = 0$. Hence $e = \eta$, and the theorem is proved for $c < \frac{1}{2}\mathbf{u}^{-1}\eta$.

It remains to prove (18) for $c > 2\mathbf{u}^{-1}\eta$. We first show

$$e \leq \frac{5}{2}\mathbf{u} \text{ufp}(c). \quad (19)$$

First,

$$c \leq \text{pred}(2\text{ufp}(c)) = 2(1 - \mathbf{u})\text{ufp}(c)$$

follows by (6) and (10) and also for $2\text{ufp}(c)$ in the overflow range, so that

$$\phi c = \mathbf{u}(1 + 2\mathbf{u})c < 2\mathbf{u}(1 + \mathbf{u})\text{ufp}(c) =: (1 + \mathbf{u})C. \quad (20)$$

Since $\text{ufp}(c) \geq 2\mathbf{u}^{-1}\eta$ it follows that $C = 2\mathbf{u} \text{ufp}(c) \in \mathbb{F}$. If $C \geq \frac{1}{2}\mathbf{u}^{-1}\eta$, then

$$\phi c < (1 + \mathbf{u})C = \mathcal{M}^+(C),$$

because $C = \text{ufp}(C)$, and (12) yields $c' = \text{fl}(\phi c) \leq C$. If $C < \frac{1}{2}\mathbf{u}^{-1}\eta$, then $C \leq \frac{1}{4}\mathbf{u}^{-1}\eta$ and (20) and $C \in \mathbb{F}$ give

$$c' = \text{fl}(\phi c) \leq \text{fl}(C + \frac{1}{4}\eta) = C,$$

so that $\text{ufp}(c) \geq 2\mathbf{u}^{-1}\eta$ with $c' \leq C$ in both cases — $C \geq \frac{1}{2}\mathbf{u}^{-1}\eta$ and $C < \frac{1}{2}\mathbf{u}^{-1}\eta$ — proves

$$e = \text{fl}(c' + \eta) \leq \text{fl}(C + \frac{1}{2}\mathbf{u} \text{ufp}(c)) = \text{fl}(\frac{5}{2}\mathbf{u} \text{ufp}(c)) = \frac{5}{2}\mathbf{u} \text{ufp}(c)$$

and thus (19).

Now we prove the equalities in (18) for $c > 2\mathbf{u}^{-1}\eta$. From (9) and (19) and (10) we have

$$c + e \leq \text{succ}(c) + \frac{1}{2}\mathbf{u} \text{ufp}(c) < \mathcal{M}^+(\text{succ}(c)),$$

since $\mathcal{M}^+(\text{succ}(c)) \geq \text{succ}(c) + \mathbf{u} \text{ufp}(c)$, so that $\mathbf{c}\text{sup} = \text{fl}(c+e) \leq \text{succ}(c)$ by (11). Together with Theorem 1, this proves the right equality in (18). A similar argument applies when $\text{pred}(\text{pred}(c)) \geq \text{ufp}(c)$ and shows $\mathbf{c}\text{inf} = \text{fl}(c-e) \geq \text{pred}(c)$ in that case. It remains to prove the left equality in (18) for $\text{ufp}(c) \in \{\text{pred}(c), c\}$.

If $\text{pred}(c) = \text{ufp}(c)$, i.e., $c = \text{ufp}(c)(1 + 2\mathbf{u})$, then we prove $e < \frac{5}{2}\mathbf{u} \text{ufp}(c)$ as follows:

$$e = \text{fl}(c' + \eta) \leq \text{fl}(\mathbf{u}(\frac{3}{2} + 4\mathbf{u})\text{ufp}(c)) \leq \text{fl}(2\mathbf{u} \text{ufp}(c)) = 2\mathbf{u} \text{ufp}(c),$$

since we proved above that $C = 2\mathbf{u} \text{ufp}(c)$ is in \mathbb{F} . Hence

$$c - e > c - \frac{5}{2}\mathbf{u} \text{ufp}(c) = \text{pred}(c) - \frac{1}{2}\mathbf{u} \text{ufp}(c),$$

and $\text{fl}(c - e) \geq \text{pred}(c)$ follows.

Finally, if $c = \text{ufp}(c)$, then c is a power of 2 and $c > 2\mathbf{u}^{-1}\eta$ implies $c \geq 4\mathbf{u}^{-1}\eta$. We distinguish three cases. First, assume $c \geq \frac{1}{2}\mathbf{u}^{-2}\eta$. Then $\phi c \in \mathbb{F}$ and $c' = \text{fl}(\phi c) = \phi c = \mathbf{u}c + 2\mathbf{u}^2c \leq \frac{9}{8}\mathbf{u}c$ by $\mathbf{u} \leq \frac{1}{16}$. Furthermore, $\mathbf{u}c \geq \frac{1}{2}\mathbf{u}^{-1}\eta \geq 8\eta$, so that $c' + \eta \leq \frac{9}{8}\mathbf{u}c + \frac{1}{8}\mathbf{u}c$ implies

$$e = \text{fl}(c' + \eta) \leq \text{fl}(\frac{5}{4}\mathbf{u}c) = \frac{5}{4}\mathbf{u}c < \frac{3}{2}\mathbf{u}c.$$

Second, assume $c < \frac{1}{4}\mathbf{u}^{-2}\eta$. Then $\mathbf{u}c \leq \phi c < \mathbf{u}c + \frac{1}{2}\eta$, which shows that $c' = \text{fl}(\phi c) = \mathbf{u}c$. Hence

$$e = \text{fl}(c' + \eta) \leq \text{fl}(\frac{5}{4}\mathbf{u}c) = \frac{5}{4}\mathbf{u}c < \frac{3}{2}\mathbf{u}c.$$

Therefore (10) implies $\text{pred}(c) = (1 - \mathbf{u})c$ and $\mathcal{M}^-(\text{pred}(c)) = (1 - \frac{3}{2}\mathbf{u})c$, and (12) finishes the first and the second case.

Third, assume $c = \frac{1}{4}\mathbf{u}^{-2}\eta$. Then

$$\phi c = \mathbf{u}c + \frac{1}{2}\eta,$$

and this is the only case where the tie-breaking rule is important. If the computation of $\text{fl}(\phi c)$ is rounded to nearest with ties to even, then $c' = \text{fl}(\phi c) = \mathbf{u}c$ and the proof of the second case holds. Let us finally assume rounding to nearest with ‘‘ties to away’’. Then with $\mathbf{u}c = \frac{1}{4}\mathbf{u}^{-1}\eta$:

$$c' = \text{fl}(\phi c) = \text{fl}(\mathbf{u}c + \frac{1}{2}\eta) = \mathbf{u}c + \eta$$

and

$$e = \text{fl}(c' + \eta) = \mathbf{u}c + 2\eta.$$

Hence

$$\mathbf{c}\text{sup} = \text{fl}(c + e) = \text{fl}(\mathbf{u}c(1 + \mathbf{u} + 8\mathbf{u}^2)) = \mathbf{u}c(1 + 2\mathbf{u}) = \text{succ}(c).$$

Finally, $\mathbf{u} \leq \frac{1}{16}$ and rounding ties to away implies

$$\mathbf{c}\text{inf} = \text{fl}(c - e) = \text{fl}(\mathbf{u}c(1 - \mathbf{u} - 8\mathbf{u}^2)) = \mathbf{u}c(1 - \mathbf{u}) = \text{pred}(c).$$

That last subcase ends the proof. □

We mention that the assumption $\mathbf{u} \leq \frac{1}{16}$ is necessary to prove the inequalities in (13) to be equalities. This is seen by $\mathbf{u} = \frac{1}{8}$, $\eta = \frac{1}{32}$ and $c = 1 \notin [\frac{1}{2}, 2]\mathbf{u}^{-1}\eta = [\frac{1}{8}, \frac{1}{2}]$.

Given $a, b \in \mathbb{F}$, rigorous and mostly sharp bounds for $a \circ b, \circ \in \{+, -, \times, \div\}$ can be computed by applying Algorithm 1 to $c := \text{fl}(a \circ b)$. This holds for the square root as well. Although addition and subtraction cause no error if the result is in the underflow range, the extra term η cannot be omitted in the computation of e because it is needed for c slightly outside the underflow range.

Algorithm 1 computes floating-point numbers \mathbf{cinf} and \mathbf{csup} bounding the predecessor and successor of its input c from below and above, respectively. In binary arithmetic, although the inequalities in (13) are proved to be equalities outside a small range near underflow, one may wish to eliminate exceptional cases so that always equality holds in (13). This is achieved by the following algorithm.

Algorithm 2 *Computation of the predecessor and successor of finite $c \in \mathbb{F}$ in binary arithmetic with rounding to nearest*

```

if  $|c| \geq \frac{1}{2}\mathbf{u}^{-2}\eta$  then
   $e = \text{fl}(\phi|c|)$            %  $\phi = \mathbf{u}(1 + 2\mathbf{u}) = \text{succ}(\mathbf{u})$ 
   $\mathbf{cinf} = \text{fl}(c - e)$ 
   $\mathbf{csup} = \text{fl}(c + e)$ 
elseif  $|c| < \mathbf{u}^{-1}\eta$ 
   $\mathbf{cinf} = \text{fl}(c - \eta)$ 
   $\mathbf{csup} = \text{fl}(c + \eta)$ 
else
   $C = \text{fl}(\mathbf{u}^{-1}c)$ 
   $e = \text{fl}(\phi|C|)$ 
   $\mathbf{cinf} = \text{fl}(\text{fl}(C - e) \cdot \mathbf{u})$ 
   $\mathbf{csup} = \text{fl}(\text{fl}(C + e) \cdot \mathbf{u})$ 

```

Theorem 3 *Let finite $c \in \mathbb{F}$ be given in binary radix, and let $\mathbf{cinf}, \mathbf{csup} \in \mathbb{F}$ be the quantities computed by Algorithm 2. Assume $\mathbf{u} \leq \frac{1}{16}$, and that $\frac{1}{2}\mathbf{u}^{-3}\eta$ does not cause overflow. Then*

$$\mathbf{cinf} = \text{pred}(c) \quad \text{and} \quad \text{succ}(c) = \mathbf{csup}. \quad (21)$$

PROOF. By the symmetry of \mathbb{F} and the symmetry of the formulas in Algorithm 2 we may assume without loss of generality that $c > 0$. As for Algorithm 1 one verifies the assertions for c being the largest positive floating-point number, hence we assume without loss of generality that $\text{succ}(c)$ is finite. We first prove (21) for the “if”-clause. Denote

$$\begin{aligned} e' &= \text{fl}(\text{fl}(\phi|c|) + \eta) \\ \mathbf{cinf}' &= \text{fl}(c - e') \\ \mathbf{csup}' &= \text{fl}(c + e') \end{aligned}$$

These are the quantities (with a prime appended) computed by Algorithm 1. The monotonicity (2) implies $e \leq e'$ and therefore

$$\mathbf{cinf}' \leq \mathbf{cinf} \quad \text{and} \quad \mathbf{csup} \leq \mathbf{csup}'.$$

Then $\mathbf{u} \leq \frac{1}{16}$, $|c| \geq 8\mathbf{u}^{-1}\eta$ and Theorem 2 imply

$$\mathbf{cinf}' = \text{pred}(c) \quad \text{and} \quad \text{succ}(c) = \mathbf{csup}'.$$

Hence, (21) is true if we prove

$$c - e < \mathcal{M}^-(c) \quad \text{and} \quad \mathcal{M}^+(c) < c + e$$

because this implies $\mathbf{cinf} = \text{fl}(c - e) \leq \text{pred}(c)$ and $\mathbf{csup} = \text{fl}(c + e) \geq \text{succ}(c)$, respectively. But $\mathbf{u}\text{succ}(c) > \frac{1}{2}\mathbf{u}^{-1}\eta$ and (15) give

$$e = \text{fl}(\phi c) \geq \text{fl}(\mathbf{u}\text{succ}(c)) = \mathbf{u}\text{succ}(c) > \mathbf{u}\text{ufp}(c),$$

and the result follows by (9) and (10).

The validity of (21) in the “elseif”-clause follows by (8). The “else”-clause is applied to $\mathbf{u}^{-1}\eta \leq |c| < \frac{1}{2}\mathbf{u}^{-2}\eta$, with $|C| \geq \mathbf{u}^{-2}\eta$, so that the computation of C cannot cause overflow and $C = \text{fl}(\mathbf{u}^{-1}c) = \mathbf{u}^{-1}c$. Now the code is similar to the one in the “if”-clause applied to a scaled c , where the final multiplication by \mathbf{u} is exact too. \square

If c is not very near to the underflow range, Algorithm 2 needs 2 flops to compute only one neighbor, and 3 flops to compute both. In any case a branch is needed, which we intentionally avoided in Algorithm 1.

Algorithms 1 and 2 may be used to bound the true value of an operation $\circ \in \{+, -, \times, \div\}$: If $c = \text{fl}(a \circ b)$, then

$$\mathbf{cinf} \leq a \circ b \leq \mathbf{csup} \tag{22}$$

for \mathbf{cinf} , \mathbf{csup} computed by Algorithm 1 or 2. However, no rounding error occurs for addition or subtraction of floating-point numbers if the result is less than $\mathbf{u}^{-1}\eta$. Hence one may try to develop an algorithm to compute \mathbf{cinf} , \mathbf{csup} satisfying

$$\begin{aligned} \mathbf{cinf} &= \text{pred}(c) & \text{and} & & \text{succ}(c) &= \mathbf{csup} & \text{for } |c| \geq \mathbf{u}^{-1}\eta, \\ \mathbf{cinf} &= c & \text{and} & & c &= \mathbf{csup} & \text{otherwise.} \end{aligned} \tag{23}$$

Then (22) would still be satisfied. To do this one may try to use

$$\begin{aligned} e &= \text{fl}(\phi'|c|) \\ \mathbf{cinf} &= \text{fl}(c - e) \\ \mathbf{csup} &= \text{fl}(c + e) \end{aligned}$$

with a suitable factor ϕ' other than $\phi = \mathbf{u}(1 + 2\mathbf{u})$. In fact, in binary arithmetic, the factor ϕ in Algorithm 1 or 2 can be varied in a wide range without jeopardizing the quality of the results. With rounding to nearest with ties to even, there is no universal factor ϕ' to satisfy (23). Consider $\phi' = \mathbf{u} \text{pred}(\frac{3}{2})$ and $c = \mathbf{u}^{-1}\eta$, then

$$e = \text{fl}(\text{pred}(\frac{3}{2})\eta) = \eta \quad \text{and} \quad \mathbf{csup} = \text{fl}(\mathbf{u}^{-1}\eta + \eta) = \mathbf{u}^{-1}\eta = c \neq \text{succ}(c),$$

so $\phi' \geq \frac{3}{2}\mathbf{u}$ is necessary. But then we obtain for $c = 1$

$$e = \text{fl}(\frac{3}{2}\mathbf{u}) = \frac{3}{2}\mathbf{u} \quad \text{and} \quad \mathbf{cinf} = \text{fl}(1 - \frac{3}{2}\mathbf{u}) = 1 - 2\mathbf{u} = \text{pred}(\text{pred}(c)) \neq \text{pred}(c).$$

With rounding to nearest ties to away, both cases above give as expected $\text{succ}(c)$ and $\text{pred}(c)$, respectively, and one might hope to find a universal factor for this rounding mode. However, consider $\phi' = \mathbf{u}$ and $c = \text{pred}(\mathbf{u}^{-1}\eta) = \mathbf{u}^{-1}(1 - \mathbf{u})\eta$. Then

$$e = \text{fl}((1 - \mathbf{u})\eta) = \eta \quad \text{and} \quad \mathbf{csup} = \text{fl}(c + e) = \text{succ}(c) \neq c,$$

so $\phi' < \mathbf{u}$ is necessary. But for $c = 1$ then

$$e = \text{fl}(\phi') = \phi' < \mathbf{u} \quad \text{and} \quad \mathbf{csup} = \text{fl}(1 + \mathbf{u}) = 1 \neq \text{succ}(c),$$

and (23) is again not satisfied.

We have implemented Algorithm 2 in the C language, and compared it to the C99 `nextafter(c, +Inf)` function call. We obtained the following timings in seconds for 10 million calls on a 2.667Ghz Core 2 under Linux with gcc 4.3.0 (using inline code):

| input | <code>nextafter</code> | Algorithm 2 |
|--------------------------|------------------------|-------------|
| $[2^{-969}, 2^{1024})$ | 1.020s | 0.464s |
| $[2^{-1021}, 2^{-969})$ | 1.020s | 0.639s |
| $[2^{-1022}, 2^{-1021})$ | 1.020s | 6.220s |
| $[2^{-1074}, 2^{-1022})$ | 12.434s | 30.743s |
| NaN | 0.425s | 0.006s |
| +Inf | 0.596s | 0.208s |
| -Inf | 0.990s | 0.497s |

The range $[2^{-969}, 2^{1024})$ corresponds to the “if”-branch of Algorithm 2, $[2^{-1021}, 2^{-969})$ to the “else”-branch, and $[2^{-1074}, 2^{-1021})$ to the “elseif”-branch. The timings for $[2^{-1022}, 2^{-1021})$ are better than those for $[2^{-1074}, 2^{-1022})$; this is most probably due to the fact that $[2^{-1022}, 2^{-1021})$ is still in the normal range, while $[2^{-1074}, 2^{-1022})$ is in the subnormal range, for which floating-point operations are usually slower in hardware. In conclusion, Algorithm 2 is clearly faster than `nextafter`, except in a tiny range near the subnormal domain.

3 Conclusion

We presented algorithms to compute (bounds for) the predecessor and successor of a floating-point number. The user may choose between the branch-free and fast variant Algorithm 1, the results of which always bound the true result, and are equal to neighbors except a tiny range near underflow. Or the slightly slower Algorithm 2 computing always the neighbors of a floating-point number, without exception, which clearly outperforms `nextafter`.

Our algorithms may be used to simulate interval operations without changing the rounding mode. The advantage of Algorithm 1 over formula (5) — for computing rigorous bounds of $a \circ b$ — is that the width is often halved (see Fig. 2).

| range | formula (5) | Algorithm 1 |
|--------------------------|----------------------|----------------------|
| $[1/2, 1]$ | 2 / 2 / 2.999696 / 4 | 2 / 2 / 2.000000 / 2 |
| $[2^{-1020}, 2^{-1019}]$ | 2 / 4 / 3.250584 / 4 | 2 / 2 / 2.000000 / 2 |
| $[2^{-1021}, 2^{-1020}]$ | 4 / 4 / 4.000000 / 4 | 2 / 2 / 2.999696 / 4 |
| $[2^{-1022}, 2^{-1021}]$ | 4 / 4 / 4.999696 / 6 | 4 / 4 / 4.000000 / 4 |
| $[2^{-1023}, 2^{-1022}]$ | 4 / 4 / 4.000000 / 4 | 2 / 2 / 2.000000 / 2 |
| $[2^{-1024}, 2^{-1023}]$ | 2 / 2 / 2.000000 / 2 | 2 / 2 / 2.000000 / 2 |

Figure 2: Minimal, median, average and maximal difference in ulps between bounds of enclosing interval with formula (5) and Algorithm 1, for IEEE 754 `binary64` format. For each range between consecutive powers of two, 1,000,000 random values of c were used. The excluded range $[\frac{1}{2}, 2]u^{-1}\eta$ from Theorem 2 corresponds to $[2^{-1022}, 2^{-1020}]$.

However, this applies only if $a, b \in \mathbb{F}$. In applications often interval operations $A \circ B$ for thick intervals A, B are executed. The wider A and B are, the less is the gain of Algorithm 1 compared to (5).

References

- [1] S. Boldo and J.-M. Muller. Some Functions Computable with a Fused-mac. In Paolo Montuschi and Eric Schwarz, editors, *Proceedings of the 17th Symposium on Computer Arithmetic*, pages 52–58, Cape Cod, USA, 2005.
- [2] W.J. Cody, Jr. and J.T. Coonen. Algorithm 722: Functions to support the IEEE standard for binary floating-point arithmetic. *ACM Transactions on Mathematical Software*, 19(4):443–451, 1993.

- [3] Institute of Electrical, and Electronic Engineers. IEEE Standard for Floating-Point Arithmetic. *IEEE Standard 754-2008*. Revision of ANSI-IEEE Standard 754-1985. Approved June 12, 2008: IEEE Standards Board.
- [4] American National Standards Institute, Institute of Electrical, and Electronic Engineers. IEEE Standard for Radix-Independent Floating-Point Arithmetic. *ANSI/IEEE Standard, Std. 854-1987*, 1987.
- [5] R.B. Kearfott, M. Dawande, K. Du, and C. Hu. Algorithm 737: INTLIB: A Portable Fortran 77 Interval Standard-Function Library. *ACM Transactions on Mathematical Software*, 20(4):447–459, 1994.
- [6] S.M. Rump, T. Ogita, and S. Oishi. Accurate Floating-point Summation Part I: Faithful Rounding. Accepted for publication in SISC, 2008.
- [7] M. Daumas, L. Rideau, and L. Théry. *A generic library of floating-point numbers and its application to exact computing*. In 14th International Conference on Theorem Proving in Higher Order Logics, Edinburgh, Scotland, pp. 169–184, 2001.
- [8] S. Boldo. Preuves formelles en arithmétiques à virgule flottante. Ph.D. dissertation, École Normale Supérieure de Lyon, 2004.