

# A correctly rounded binary64 cube root

Robin Leroy (eggrobin)

REMOVE BEFORE FLIGHT 2021-04-36

This document describes the implementation of correctly-rounded binary64 cube root function `principia::numerics::Cbrt` defined in `numerics/cbrt.cpp`.

## Outline

Our computation of the cube root uses root-finding methods that have fallen into disuse, and thus may be unfamiliar to the reader.

So that there can be no confusion on the properties and names of these methods, this document comprises two parts. The first one is about a family of root-finding methods for arbitrary functions, wherein all methods considered for the cube root can be found. The second is about the computation of a binary64 cube root and its error analysis.

The methods considered in the first part originate from a series of works by Thomas Fantet de Lagny<sup>1</sup>. Lagny's methods may be used to find a root of any polynomial; curiously, they are of convergence order  $p$  for polynomials of degree  $p$ . We give a generalization of these methods which applies to arbitrary functions (and in particular to polynomials of degree other than the order).

Throughout the centuries, some of Lagny's methods have been rediscovered, whether from Lagny's works or from first principles, and generalized to arbitrary functions. The reader may be familiar with these special cases, as some of them have remained in use, and are found in more modern literature; most famous of these is perhaps Halley's (rational) method. We thus show that these are indeed special cases of the generalization of Lagny's method, and describe the names we use for specific cases, depending on the order of their discovery.

Our path through more than three centuries of literature to rediscovering—and then generalizing—Lagny's methods led us to many remarkable works; we mention those in a bibliographic note.

The second part starts with the treatment of a faithfully rounded—and very nearly correctly rounded—cube root. For each step in the computation of the cube root, we consider multiple alternatives, each with its error analysis; the one that should be chosen is revealed by the final error analysis.

We then describe how correct rounding is achieved, by determining whether the faithful result might be misrounded, in which case the correct digit needs to be more carefully ascertained.

In order not to interrupt the flow of reading, in both parts, we relegate to the appendices any miscellaneous considerations that are independent of the matter currently being discussed, long-winded proofs, and lists of examples.

---

<sup>1</sup>Thomas Fantet de Lagny (1660–1734), professor of hydrography in Rochefort, 1697–1714, subdirector then director of Law's *Banque générale*. Member of the *Académie Royale des Sciences*. See [Fon34], reprinted in [Fon58, vol. 6, pp. 557 sqq.].

## Part I

# Abridged root-finding methods

We recall and generalize a family of root-finding methods from the late 17th century.

In [Fan91a], Lagny first presents the following iterations for the computation of the cube root  $\sqrt[3]{a^3 + b}$ :

$$a \mapsto \frac{1}{2}a + \sqrt{\frac{1}{4}a^2 + \frac{b}{3a}}, \quad (1)$$

hereafter the (*quadratic*) *irrational method*, and

$$a \mapsto a + \frac{ab}{3a^3 + b}, \quad (2)$$

the *rational method*, mentioning the existence of similar methods for arbitrarily higher powers. In [Fan91b] the above methods are again given, with an outline of the general method for higher powers, and a mention of their applicability to finding roots of polynomials other than  $z^p - r$ .

That general method is given in detail in [Fan92, p. 19]. Modernizing the notation, the general rule is as follows for finding a root of the monic polynomial of degree  $p \geq 2$

$$f(z) := z^p + c_1 z^{p-1} + \dots + c_{p-1} z + c_p =: z^p - R(z)$$

with an initial approximation  $a$ .

Separate the binomial expansion of  $(x + \frac{1}{2}a)^p$  into alternating sums of degree  $p$  and  $p - 1$  in  $x$ ,

$$S_1 := \sum_{\substack{k=0 \\ 2 \nmid k}}^p \binom{p}{k} x^{p-k} \left(\frac{1}{2}a\right)^k \text{ and } S_2 := \sum_{\substack{k=1 \\ 2 \nmid k}}^p \binom{p}{k} x^{p-k} \left(\frac{1}{2}a\right)^k,$$

and consider the following polynomials, of degree  $p$  and  $p - 1$  in  $x$  for almost all  $a$ :

$$E_p := S_1 - \frac{1}{2}R\left(x + \frac{1}{2}a\right) \text{ and } E_{p-1} := S_2 - \frac{1}{2}R\left(x + \frac{1}{2}a\right). \quad (3)$$

Let  $E_{n-1}$  be the remainder of the polynomial division<sup>2</sup> of  $E_{n+1}$  by  $E_n$ ; its degree is  $n - 1$  for almost all  $a$ . The iteration for finding a root of  $f$  is  $a \mapsto x + \frac{1}{2}a$ , where  $x$  is a root of  $E_2$  in the quadratic irrational method, and the root of  $E_1$  in the rational method. Its order is  $p$ .

## Multiplicity of the irrational methods

Lagny does not require that the polynomial division be carried out all the way to  $E_2$ , merely until one gets *une valeur d' $x$  [...] d'un degré commode*, by which he likely means one that is constructible. When  $f$  is a cubic, he uses the term *formule irrationnelle* for  $x + \frac{1}{2}a$  where  $x$  is a root of  $E_2$ , but when it comes to computing the fifth

<sup>2</sup>While the rest of the method is a straightforward translation, this step bears some explanations; its description in [Fan92] is

De ces deux égalitez, ou prifes féparément, ou comparées ensemble felon la methode des problèmes plus que déterminez tirez en une valeur d' $x$  rationelle, ou fimplement d'un degré commode.

It is assumed that the reader is familiar with this “comparison according to the method of more-than-determined problems”. While the application of the root-finding method is described in painstaking detail in [Fan33], which outlines the treatment of overdetermined problems, it is perhaps this remark from [Fan97, p. 494] which lays it out most clearly:

Il n'y a rien de nouveau à remarquer fur les Problemes plus que déterminez du quatrième degré. La Regle générale est d'égalier tout à zero, & de divifer la plus haute équation par la moins élevée, ou l'également élevée l'une par l'autre, continuellement jufques à ce que l'on trouve le refte ou le divifeur le plus fimple.

root, the same term is used to refer to the case where  $x$  is a root of  $E_4$ . In order to avoid confusion, we use the term *quadratic irrational method* when  $x$  is a root of  $E_2$ , and we call the irrational formula for  $\sqrt[5]{a^5 + b}$  from [Fan92, p. 43]<sup>3</sup>

$$a \mapsto \frac{1}{2}a + \sqrt{\sqrt{\frac{1}{4}a^4 + \frac{b}{5a}} - \frac{1}{4}a^2}$$

Lagny's *quartic irrational method* for the fifth root; the quadratic irrational method for the same fifth root would be<sup>4</sup>

$$a \mapsto \frac{a(7b - \sqrt{100a^{10} + 100a^5b - 7b^2})}{4b - 10a^5}.$$

More generally, we call  $a \mapsto x + \frac{1}{2}a$  Lagny's *method of degree  $d$*  when  $x$  is a root of  $E_d$ . Note however that when  $p = 3$ , i.e., when finding a root of a cubic, Lagny's only irrational method is the quadratic one; we can thus unambiguously refer to (1) as *Lagny's irrational method for the cube root*.

## Generalization to arbitrary functions

Lagny's method of degree  $d$  and convergence order  $p$  may be generalized to functions  $f$  other than polynomials of degree  $p$ , by defining  $E_p$  and  $E_{p-1}$  in terms of Taylor polynomials for  $x \mapsto f(x + \frac{1}{2}a)$  around  $x = \frac{1}{2}a$ :

$$E_p := T_p - \frac{1}{2}T_{p-1} \text{ and } E_{p-1} := \frac{1}{2}T_{p-1}, \quad (4)$$

where

$$T_n := \sum_{k=0}^n \frac{f^{(k)}(a)}{k!} \left(x - \frac{1}{2}a\right)^k.$$

The rest of the method remains unchanged; the iteration is given by  $a \mapsto x + \frac{1}{2}a$  for a root  $x$  of  $E_d$ . When  $f$  is a monic polynomial of degree  $p$ , the definitions (4) are equivalent to (3), so that we recover Lagny's method. When  $f$  is not a polynomial of degree  $p$ , we thus call the method defined by (4) the *generalized Lagny method of degree  $d$  and order  $p$* ; we use the terms "rational", "quadratic irrational", etc. for  $d = 1$ ,  $d = 2$ , etc. respectively.

Note that while the  $E_n$  defined in this fashion may not have degree  $n$  if the higher derivatives of  $f$  vanish, e.g., for a polynomial of degree less than  $p$ , the calculation may be carried out formally for an arbitrary  $f$ , and the offending function substituted in the result, taking limits as needed to remove singularities; the generalized methods of high order can thus be applied to polynomials of low degree.

These methods may equivalently be constructed using the Maclaurin series for  $\Delta \mapsto f(a + \Delta)$  in the correction term  $\Delta$ . Let

$$M_n := \sum_{k=0}^n \frac{f^{(k)}(a)}{k!} \Delta^k,$$

and consider the polynomials  $\tilde{E}_p := M_p$  and  $\tilde{E}_{p-1} := M_{p-1}$  of degree  $p$  and  $p - 1$  in  $\Delta$ . Let  $\tilde{E}_{n-1}$  be the remainder of the polynomial division of  $\tilde{E}_{n+1}$  by  $\tilde{E}_n$ . The iteration is then  $a \mapsto a + \Delta$ , where  $\Delta$  is a root of  $\tilde{E}_d$ .

Examples of this method for specific values of the function  $f$ , the order  $p$ , and the degree  $d$  are given in appendix B.

<sup>3</sup>The formula has a misprint in [Fan92, p. 43],  $-\frac{1}{2}a^2$  instead of  $-\frac{1}{4}a^2$  under the radical. Halley remarks on it and gives the corrected formula in [Hal94, pp. 137, 140]. The misprint remains forty years later in [Fan33, p. 440 misnumbered 340]. Bateman writes in [Bat38] "we must not infer that [these expressions] are not correct simply because they differ from Halley's expression", but with Lagny's construction, which was seemingly unknown to Bateman, the error is plain.

<sup>4</sup>Both are of order 5, but the reader who wishes to compute a fifth root should note that leading term of the error of the quartic method is  $\frac{2}{7}$  of that of the quadratic. See appendix B.

## Names

The generalized Lagny methods form a broad family; many of the methods therein are well-known. Since these special cases are better known under the names of their discoverers, we give a list of such occurrences, and use the appropriate names in the remainder of this document.

### Names of the irrational methods

Halley generalized Lagny's (quadratic) irrational method for the cubic to arbitrary<sup>5</sup> polynomials, retaining cubic convergence; when  $f$  is not a polynomial of degree 3, we thus call the generalized Lagny quadratic irrational method of order 3 *Halley's irrational method*. This method was given in terms of derivatives by Bateman in [Bat38, p. 12]:

$$a \mapsto a - \frac{f'(a)}{f''(a)} + \frac{\sqrt{f'^2(a) - 2f(a)f''(a)}}{f''(a)}.$$

### Names of the rational method

Both special cases and generalizations of Lagny's rational method have been discovered multiple times and extensively studied; constructions that take advantage of modern calculus allow us to give a more straightforward expression for the rational method than was available to Lagny. The proof of the following equivalence is given in appendix A.

**Proposition.** *The iteration of the generalized Lagny rational method of order  $p$  for a root of the function  $f$  is*

$$a \mapsto a + (p-1) \frac{(1/f)^{(p-2)}(a)}{(1/f)^{(p-1)}(a)}. \quad (5)$$

□

The iteration (5) is a special case of the *Algorithmen* ( $A_\omega^\lambda$ ) defined by Schröder for an arbitrary polynomial  $f$  in [Sch70, pp. 349 sq.], equation (69); specifically, it is ( $A_{p-1}^0$ ). As seen in the proof of the proposition, it is also a special case of Householder's equation (14) from [Hou70, p. 169], which generalizes it by substituting  $f/g$  for  $f$ . The case  $g \equiv 1$  is mentioned in theorem 4.4.2, and that expression is given explicitly in [SG01].

For  $p = 2$  and  $f$  an arbitrary polynomial, (5) is Newton's method, presented by Wallis in [Wal85, p. 338].

For  $p = 3$  and  $f$  an arbitrary polynomial, it is Halley's rational method, given in [Hal94, pp. 142–143] in an effort to generalize Lagny's (2). It is usually simply known as Halley's method, as his aforementioned irrational method has comparatively fallen into obscurity; see [ST95].

Considering, as remarked by [Sch70, p. 334], that a method can often be generalized from arbitrary polynomials or rational functions to arbitrary analytic functions,

---

<sup>5</sup>Lagny's method is general, in that an iteration is given for any polynomial, albeit one whose order changes with the degree. However, while he refers to its results—and even corrects a misprint therein—, Halley did not have access to a copy of [Fan92].

Has Regulas, cum nondum librum videram, ab amico communicatas habui

and it appears that said friend communicated only the formulæ for the cube and fifth root, as opposed to the general method and its proof, as Halley writes

[...] *D. de Lagny* [...] qui cum totus fere sit in eliciendis Potestatum purarum radicibus, præfertim Cubicâ, pauca tantum eaque perplexa nec fatis demonstrata de affectarum radicum extractione subiungit.

or, about the quartic irrational method for the fifth root, whereon Lagny does not elaborate as it is a direct application of the general method,

Author autem nullibi inveniendi methodum ejufve demonstrationem concedit, etiamfi maxime defiderari videatur [...].

Being unaware of this generality, Halley sets out to generalize (1) and (2) to arbitrary polynomials, and does so by keeping the order constant.

we call the iteration (5)

- Newton’s method when  $p = 2$ , for arbitrary  $f$ ;
- Lagny’s rational method when  $p > 2$  and  $f$  is a polynomial of degree  $p$ ;
- Halley’s rational method when  $p = 3$  and  $f$  is not a polynomial of degree 3;
- the Lagny–Schröder rational method of order  $p$  otherwise.

We do not simply call this last case “Schröder’s method”, as it is only a special case of the methods defined in [Sch70], so that the expression would be ambiguous.

Note that we avoid the name “Householder’s method” which appears in [SG01] and ulterior works (notably *MathWorld* and *Wikipedia*, both citing [SG01]), as it is variably used to refer to either (5) or to a method from a different family, namely  $\varphi_{p+1}$  from [Hou70, p. 168], equation (7), taking  $\gamma_{p+1} \equiv 0$  in the resulting iteration;  $\varphi_3$  is<sup>6</sup> the iteration given in section 3.0.3 of [SG01]. As mentioned by Householder, both of those were described by Schröder a century prior anyway: Householder’s (7) is Schröder’s (18) from [Sch70, p. 327].

## Bibliographic note

Our foray into the history of these methods was prompted by finding the “historical background” section of [ST95] while looking for a reference for Halley’s method: it is mentioned therein that this method, as applied to the cube root, is due to Lagny.

Searching for Lagny’s work led us to the historical note [Can61], wherein a note by the editors Terquem and Gerono reads

Naturellement, en mathématiques, séjour des propositions irréfragables, identiques en toute langue, en tout pays, ces rencontres ne peuvent manquer d’être assez fréquentes; nulle part les plagiats *effectifs* sont si rares, et les plagiats *apparents* si communs que dans la science exacte par excellence; mais les signaler est un devoir, un service rendu à l’histoire scientifique.

The editors then quote a letter by Prouhet, wherein he gives a reference to [Fan92].

Lagny’s work proved far more extensive than we expected: besides the above root finding methods for arbitrary polynomials, it contains an error analysis, and even a discussion of the principles of performance analysis based on a decomposition into elementary operations on—and writing of—decimal digits [Fan92, pp. 5–9], with a remark on applicability to bases other than ten: a 17th century MIX.

Observing that the higher-order examples correspond to the well-known higher order method attributed to Householder in [SG01], we looked for its properties in [Hou70] so as to prove that observation, and found that Householder attributes them to Schröder. As mentioned in the translator’s note by Stewart in [SS93],

A. S. Householder used to claim you could evaluate a paper on root finding by looking for a citation of Schröder’s paper. If it was missing, the author had probably rediscovered something already known to Schröder.

It is possible that the irrational methods could be expressed using Schröder’s results in one way or another, although most of his methods seem to be rational; in any case, such a formulation is unlikely to be something well-known, as irrational methods are far less popular nowadays—unjustifiedly so, as we shall see.

Our generalization of Lagny’s irrational methods to arbitrary  $f$ , which, in the polynomial case, decouples the degree of  $f$  from the convergence order, was inspired by Gander’s rephrasing in [Gan85] of Halley’s construction from [Hal94], wherein the correction term of Halley’s irrational method is defined as a root of  $M_2$ . This change of variables with respect to Lagny’s construction drastically simplifies the proof of the proposition.

Prouhet’s letter in [Can61] ends with these words:

Tout cela est fort abrégé; mais qui nous délivrera des méthodes abrégées, qui n’en finissent pas?

---

<sup>6</sup>We are grateful to Peter Barfuss for this observation.

## Part II

# Computing a real cube root

We now turn to the computation of the cube root of in `numerics/cbrt.cpp`.

## A faithfully rounded cube root

### Overview

Our general approach to computing a faithfully rounded cube root of  $y > 0$  is the one described in [KB01]:

1. integer arithmetic is used to get a an initial quick approximation  $q$  of  $\sqrt[3]{y}$ ;
2. a root finding method is used to improve that that to an approximation  $\xi$  with a third of the precision;
3.  $\xi$  is rounded to a third of the precision, resulting in the rounded approximation  $x$  whose cube  $x^3$  can be computed exactly;
4. a single high order iteration of a root finding method is used to get the faithfully rounded result  $r_0$ .

### Notation

We define the fractional part as  $\text{frac } a := a - \lfloor a \rfloor \in [0, 1[$ , regardless of the sign of  $a$ .

The floating-point format used throughout is binary64; the quantities  $p \in \mathbb{N}$  (precision in bits) and  $\text{bias} \in \mathbb{N}$  are defined as in IEEE 754-2008,  $p = 53$  and  $\text{bias} = 1023$ . Some of the individual methods discussed may be of general use; we thus give all inexact constants used in such methods, as well as the error bounds, rounded to forty decimal places and thirty-three hexadecimal places, which amply suffices for decimal128, binary128, and all smaller formats. A superscript sign after the last digit serves as the sticky bit<sup>7</sup>: the unrounded quantity is in excess of the rounded one if the sign is +, and in default if it is −.

We use capital Latin letters for fixed-point numbers involved in the computation, and  $A > 0$  for the normal floating-point number  $a > 0$  reinterpreted as a binary fixed-point<sup>8</sup> number with  $p - 1$  bits after the binary point,

$$\begin{aligned} A &:= \text{bias} + \lfloor \log_2 a \rfloor + \text{frac}(2^{-\lfloor \log_2 a \rfloor} a) \\ &= \text{bias} + \lfloor \log_2 a \rfloor + 2^{-\lfloor \log_2 a \rfloor} a - 1, \end{aligned} \tag{6}$$

and *vice versa*,

$$a := 2^{\lfloor A \rfloor - \text{bias}} (1 + \text{frac } A). \tag{7}$$

This corresponds to [KB01]’s  $B + K + F$ . For both fixed- and floating-point numbers, given  $\alpha \in \mathbb{R}$ , we write:

- $\llbracket \alpha \rrbracket$  for the nearest representable number, rounding ties to even: IEEE 754-2008 rounding-direction attribute `roundTiesToEven`;
- $\llbracket \alpha \rrbracket_+$  for the nearest representable number no smaller than  $\alpha$ : `roundTowardPositive`;

<sup>7</sup>We learned of this practice from Steve Canon, who found it in a re-edition of [Bru70, p. VIII]; there it is only present on the digit 5, to guard against double-rounding to the nearest decimal place. As mentioned in Hoüel’s foreword to [Sch73, p. II], this practice, originally seen as a way to convey another bit of precision rather than a way to ensure correct rounding, dates back to at least 1827; see [Bab27, p. X], 8th rule. Like Babbage and Schrön, we give this bit regardless of the last digit; this allows for directed rounding.

<sup>8</sup>The implementation uses integers (obtained by multiplying the fixed-point numbers by  $2^{p-1}$ ). For consistency with [KB01] we work with fixed-point numbers here. Since we do not multiply fixed point numbers together, the expressions are unchanged.

- $\llbracket \alpha \rrbracket_-$  for the nearest representable number no larger than  $\alpha$ : roundTowardNegative;
- $\llbracket \alpha \rrbracket_0$  for the nearest representable number no larger in magnitude than  $\alpha$ : roundTowardZero.

We write the unit roundoff  $u := 2^{-p}$  (for rounding to nearest), and, after [Higo2, p. 63],  $\gamma_n := \frac{nu}{1-nu}$ . We discuss other rounding modes in appendix F.

To quote [Tre97], “If rounding errors vanished, 95% of numerical analysis would remain”. While we keep track of rounding errors throughout, they are of very little importance until the last step; when it is convenient to solely study the truncation error, we work with ideal quantities affected with a prime, which correspond to their primeless counterparts by removal of all intervening roundings.

The input  $y$  and all intervening floating-point numbers are taken to be normal; the rescaling performed to avoid overflows also avoids subnormals. We work only with correctly rounded addition, subtraction, multiplication, division, and square root; FMA is treated separately in appendix C.

## 1 Quick approximation

The quick approximation  $q$  is computed using fixed-point arithmetic as

$$Q := C + \left\lfloor \frac{Y}{3} \right\rfloor_0,$$

where the fixed-point constant  $C$  is defined as<sup>9</sup>

$$C := \left\lfloor \frac{2 \text{bias} - \Gamma}{3} \right\rfloor$$

for some  $\Gamma \in \mathbb{R}$ . As we will now show, this step is effectively an argument reduction; we will discuss the choice of the free parameter  $\Gamma$  below.

Let  $\varepsilon_q := \frac{q}{\sqrt[3]{y}} - 1$ , so that  $\sqrt[3]{y}(1 + \varepsilon_q) = q$ ; the relative error of  $q$  as an approximation of  $\sqrt[3]{y}$  is  $|\varepsilon_q|$ . Considering  $Y$ ,  $Q$ ,  $q$ , and  $\varepsilon_q$  as functions of  $y$ , we have

$$\begin{aligned} Y(8y) &= Y(y) + 3, \\ Q(8y) &= Q(y) + 1, \\ q(8y) &= 2q(y), \\ \varepsilon_q(8y) &= \varepsilon_q(y), \end{aligned}$$

so that the properties of  $\varepsilon_q$  need only be studied on some interval of the form  $[\eta, 8\eta[$ .

Pick  $\eta := 2^{\lfloor \Gamma \rfloor}$ , and  $y \in [\eta, 8\eta[ = [2^{\lfloor \Gamma \rfloor}, 2^{\lfloor \Gamma \rfloor+3}[$ , so that  $\log_2 y \in [\lfloor \Gamma \rfloor, \lfloor \Gamma \rfloor + 3[$ . Let  $k := \lfloor \log_2 y \rfloor - \lfloor \Gamma \rfloor$ ; note that  $k \in \{0, 1, 2\}$ . Let  $f := \text{frac}(2^{-\lfloor \log_2 y \rfloor} y) \in [0, 1[$ , so that

$$y = 2^{\lfloor \log_2 y \rfloor} (1 + \text{frac}(2^{-\lfloor \log_2 y \rfloor} y)) = 2^{\lfloor \Gamma \rfloor + k} (1 + f).$$

Up to at most 3 half-units in the last place from rounding (2 from the directed rounding of the division by three and 1 from the definition of  $C$ ), we have, using the definition (6) of  $Y$ ,

$$\begin{aligned} Q &\approx Q' := \text{bias} + \frac{\lfloor \log_2 y \rfloor}{3} + \frac{\text{frac}(2^{-\lfloor \log_2 y \rfloor} y) - \Gamma}{3}, \\ &= \text{bias} + \frac{\lfloor \Gamma \rfloor + k}{3} + \frac{f - \Gamma}{3}, \\ &= \text{bias} + \frac{k + f - \text{frac } \Gamma}{3}. \end{aligned}$$

<sup>9</sup>Note that there is a typo in the corresponding expression  $C := (B - 0.1009678)/3$  in [KBo1]; a factor of 2 is missing on the bias term.

Since  $k \in [0, 2]$ , the numerator  $k + f - \text{frac } \Gamma$  lies in  $] -1, 3[$ . Further, it is negative only if  $k = 0$ , so that

$$\begin{aligned} \lfloor Q' \rfloor &= \begin{cases} \text{bias} - 1 & \text{if } k = 0 \text{ and } \text{frac } \Gamma > f, \text{ and} \\ \text{bias} & \text{otherwise,} \end{cases} \\ \text{frac } Q' &= \begin{cases} 1 + \frac{f - \text{frac } \Gamma}{3} & \text{if } k = 0 \text{ and } \text{frac } \Gamma > f, \\ \frac{k + f - \text{frac } \Gamma}{3} & \text{otherwise.} \end{cases} \end{aligned}$$

Accordingly, for the quick approximation  $q$ , we have, again up to at most 3 half-units in the last place, by the definition (7) of  $q$ ,

$$q \approx q' = \begin{cases} 1 + \frac{f - \text{frac } \Gamma}{3} & \text{if } k = 0 \text{ and } \text{frac } \Gamma > f, \\ 1 + \frac{k + f - \text{frac } \Gamma}{3} & \text{otherwise,} \end{cases}$$

so that  $q$  is a piecewise affine function of  $y$  and may be expressed as an affine function of  $f$  on each piece.

With  $\sqrt[3]{y} = 2^{\frac{\lfloor \Gamma \rfloor + k}{3}} \sqrt[3]{1 + f}$ , we express  $\varepsilon'_q := \frac{q'}{\sqrt[3]{y}} - 1$  piecewise as a function of  $f$  and  $k$ . The maximum of  $\varepsilon'_q$  gives us a bound on the relative error of  $q$ , as

$$|\varepsilon|_q \leq |\varepsilon'_q|(1 + 3u).$$

The values  $\Gamma = 0.1009678$  and  $\varepsilon_q < 3.2\%$  from [KB01] may be recovered by choosing  $\Gamma$  minimizing the maximum of  $|\varepsilon'_q|$  over  $y \in [\eta, 8\eta]$ , or equivalently.

$$\Gamma_{\text{Kahan}} := \operatorname{argmin}_{\Gamma \in \mathbb{R}} \max_{y \in [\eta, 8\eta]} |\varepsilon'_q| = \operatorname{argmin}_{\Gamma \in \mathbb{R}} \max_{(f, k)} |\varepsilon'_q|$$

where the maximum is taken over  $(f, k) \in [0, \text{frac } \Gamma[ \times \{0\} \cup [0, 1[ \times \{1, 2\}$ ,

$$= \operatorname{argmin}_{\Gamma \in \mathbb{R}} \max_{(f, k) \in \mathcal{E} \cup \mathcal{L}} |\varepsilon'_q|,$$

where  $\mathcal{E} := \{(\text{frac } \Gamma, 0)\} \cup \{(0, k) \mid k \in \{0, 1, 2\}\}$  is the set of the endpoints of the intervals whereon  $q'$  is piecewise affine, and  $\mathcal{L} := \left\{ \left( \frac{k - \text{frac } \Gamma}{2}, k \right) \mid k \in \{1, 2\} \right\}$  are the smooth local extrema. We get more precisely<sup>10</sup>

$$\begin{aligned} \Gamma_{\text{Kahan}} &\approx 0.10096\,78121\,55802\,88786\,36993\,42643\,55358\,06490^- \\ &\approx {}_{16}0.19\text{D9}\,06\text{CB}\,2868\,81\text{F4}\,88\text{FD}\,38\text{DF}\,\text{E7F6}\,98\text{DD}\,\text{B}^+ \end{aligned}$$

with

$$\begin{aligned} \max_y |\varepsilon'_q| &\approx 3.15546\,32773\,62480\,60611\,78973\,32817\,13558\,9400^+ \% \\ &\approx {}_{16}1.027\text{E}\,\text{DC79}\,99\text{AB}\,08\text{D3}\,928\text{D}\,83\text{B0}\,17\text{CC}\,\text{E876}^- \cdot 2^{-5} \end{aligned}$$

yielding the constant

$$C_{\text{Kahan}} = {}_{16}2\text{A9F}\,7625\,3119\,\text{D328} \cdot 2^{-52}$$

for IEEE 754-2008 binary64. However, as we will see in the next section, this value does not optimize the final error.

## 2 Getting to a third of the precision

We now consider multiple methods for the refinement of  $q$  to  $\xi$ . The rounding error in this step being both negligible and tedious to bound, its analysis is relegated to appendix E. Here we will study only the truncation error, and thus work only with the primed quantities.

<sup>10</sup>These values may be computed formally, but the expression is unwieldy.



### Lagny's rational method

One way to compute  $\xi'$  is Lagny's rational method,

$$\xi' = q' + \frac{q'(y - q'^3)}{2q'^3 + y},$$

with the error

$$\varepsilon'_\xi := \frac{\xi'}{\sqrt[3]{y}} - 1.$$

With  $q' = \sqrt[3]{y}(1 + \varepsilon'_q)$ , we can express  $\varepsilon'_\xi$  using the transformation of the relative error error by one step of Lagny's rational method on the cube root,

$$\varepsilon'_\xi = \frac{2\varepsilon_q'^3 + \varepsilon_q'^4}{3 + 6\varepsilon_q' + 6\varepsilon_q'^2 + 2\varepsilon_q'^3} = \frac{2}{3}\varepsilon_q'^3 + \mathcal{O}(\varepsilon_q'^4).$$

If  $q'$  is computed using  $\Gamma = \Gamma_{\text{Kahan}}$ , we get  $\max_y |\varepsilon'_\xi| \approx 21.96 \cdot 10^{-6}$ ,  $\log_2 \max_y |\varepsilon'_\xi| \approx -15.47$ . However,  $\Gamma_{\text{Kahan}}$ , which minimizes  $\max_y |\varepsilon_q|$ , does not minimize  $\max_y |\varepsilon'_\xi|$ . This is because while  $\varepsilon'_\xi$  is monotonic as a function of  $\varepsilon_q'$ , it is not odd: positive errors are reduced more than negative errors are, so that the minimum is attained for a different value of  $\Gamma$ . Specifically, we have

$$\begin{aligned} \Gamma_{L^1} &:= \operatorname{argmin}_{\Gamma \in \mathbb{R}} \max_y |\varepsilon'_\xi| \\ &\approx 0.09918\,74615\,29855\,99525\,66149\,20761\,31234\,34720\,2^+ \\ &\approx {}_{16}0.1964\,5977\,71A9\,4DE0\,A8AF\,47A0\,0B1B\,C052\,B^+ \end{aligned}$$

with  $\max_y |\varepsilon_q'| \approx 3.203\%$ , but

$$\begin{aligned} \max_y |\varepsilon'_\xi| &\approx 20.86863\,55363\,95934\,87709\,20083\,98441\,02541\,483^+ \cdot 10^{-6} \\ &\approx {}_{16}1.5E1E\,1B6D\,9718\,42F4\,89C2\,EC7B\,2EC0\,ECC1^- \cdot 2^{-16}, \end{aligned}$$

$\log_2 \max_y |\varepsilon'_\xi| \approx -15.55$ . The corresponding fixed-point constant is

$$C_{L^1} := {}_{16}2A9F\,7893\,782D\,A1CE \cdot 2^{-52}$$

for binary64.

While it is not far from the seventeen bits to which we will round in the next step, this error is still larger, and in any case is not comparatively negligible. As a result, it significantly contributes to misrounding, see (10). Lagny's lesser-known irrational method provides us with a way to improve it.

### Lagny's irrational method

As written in (1), Lagny's irrational method

$$\xi' = \frac{1}{2}q' + \sqrt{\frac{1}{4}q'^2 + \frac{y - q'^3}{3q'}}$$

seems prohibitively computationally expensive in comparison to the rational one: it adds a square root on the critical path, dependent on the result of a division. However, rewriting it as

$$\xi' = \frac{1}{2}q' + \frac{1/\sqrt{12}}{q'} \sqrt{4yq' - q'^4}, \quad (8)$$

one can evaluate it with similar<sup>11</sup> performance to the rational method.

Its error is

$$\varepsilon'_\xi = \frac{-\varepsilon_q'^3}{3\left(\frac{1}{2} + \sqrt{\frac{1}{2} - 2\varepsilon_q'^2 - \frac{4}{3}\varepsilon_q'^3 - \frac{1}{3}\varepsilon_q'^4 - \varepsilon_q'^2}\right)} = -\frac{1}{3}\varepsilon_q'^3 + \mathcal{O}(\varepsilon_q'^4),$$

<sup>11</sup>Other rewritings are preferable on some architectures; we discuss this in appendix D.

whose leading term is half that of the rational method; indeed we find that with  $\Gamma = \Gamma_{\text{Kahan}}$ , we have  $\max_y |\varepsilon'_\xi| \approx 10.48 \cdot 10^{-6}$ ,  $\log_2 \max_y |\varepsilon'_\xi| \approx -16.54$ , gaining one bit with respect to the rational method. Here  $\Gamma = \Gamma_{\text{Kahan}}$  is very close to optimal; with the optimal value

$$\begin{aligned} \Gamma_{L^2} &\approx 0.10096\,82076\,65096\,37285\,40885\,52460\,33434\,63385^- \\ &\approx {}_{16}0.19D9\,0D6D\,DB79\,AE1F\,D556\,591B\,78EF\,F3DD\,B^+, \end{aligned}$$

the error bound

$$\begin{aligned} \max_y |\varepsilon'_\xi| &\approx 10.48337\,57985\,85309\,87229\,03375\,83237\,37064\,369^+ \cdot 10^{-6} \\ &\approx {}_{16}1.5FC3\,832D\,82FF\,67E3\,E4A4\,C2FD\,A877\,7C2E^- \cdot 2^{-17} \end{aligned}$$

improves only in its sixth decimal place. However, we have other ways of improving the error at no cost to performance.

### Canon optimization of Lagny's irrational method

The idea for this optimization comes from [Can18a], reproduced here with the author's permission:

A trick I've used for years and should write up: you can apply optimization to the iteration, not just the starting guess:  $x' = xp(x)$ , select  $p(x)$  to be minimax error on bounded initial error in  $x$ . This yields a nice family of tunable approximations.

Everyone else seems to worry about starting estimate, but use standard iterations, which is appropriate for arbitrary precision, but silly with a fixed precision target.

Note that as  $p$  gets to be high-order, it converges quickly to the Taylor series for the correction, but there's a nice space with cheap initial approximations and order 2–5 or so, because we can evaluate these polynomials with lower latency [than] serially-dependent iterations.

Canon later elaborated on this in [Can18b]:

Quick version: we want to compute  $1/\sqrt{y}$ , we have an approximation  $x_0$ , we want to improve it to  $x_1 = x_0 p(x_0, y)$ . For efficiency, we want  $p$  to be a polynomial correction.

*handwavy motivation for brevity* make  $p$  a polynomial in  $x_0 x_0 y$ , which is approximately 1.

Specifically, if  $x_0$  has relative error  $e$ ,  $x_0 x_0 y$  is bounded by something like  $1 \pm 2e$ . So, we want to find  $p$  that minimizes  $|x/x_0 - p(x_0 x_0 y)|$  on  $[1 - 2e, 1 + 2e]$ . NR<sup>12</sup> uses the  $p = 1$ st order Taylor. We know that we can do better via usual approximation theory techniques.

We can also use higher-order approximations to hit any specific accuracy target in a single step. This isn't always better than iterating, but sometimes it is.

We do not use a polynomial—nor even a rational function—, nor do we express our refinement as a function of a quantity bounded by the error. However, we take advantage of Canon's key idea of “apply[ing] optimization to the iteration, not just the starting guess”; the latter is what we have so far done with  $\Gamma$ .

The constants  $\frac{1}{2}$ ,  $\frac{1}{4}$ , and 3 in Lagny's irrational method may be modified with no effect on performance; altering the first two of these introduces rounding errors, but these need not concern us here. We thus write

$$\xi' = \kappa q' + \sqrt{\lambda q'^2 + \frac{y - q'^3}{\mu q'}}$$

<sup>12</sup>Newton–Raphson, *i.e.*, Newton's method; Raphson described a different formulation of it in [Rap90]. Lagrange notes in [Lag67, vol. 8, p. 161] that Raphson was likely aware of Newton's method, and thus regarded his as entirely different; as Lagrange points out, they are equivalent.

and choose  $\Gamma$ ,  $\kappa$ ,  $\lambda$ , and  $\mu$  minimizing relative error in the Чебышёв norm,

$$(\Gamma_{L^2C}, \kappa_{L^2C}, \lambda_{L^2C}, \mu_{L^2C}) := \operatorname{argmin}_{\Gamma, \kappa, \lambda, \mu} \max_y |\varepsilon'_\xi|.$$

Unfortunately, computing  $\max_y |\varepsilon'_\xi|$  is not as easy as for the standard methods; the introduction of  $\kappa$ ,  $\lambda$ , and  $\mu$  breaks the monotonicity of  $\varepsilon'_\xi(\varepsilon'_q)$ , so that the local extrema of  $\varepsilon'_\xi$  are not found in the same place as those of  $\varepsilon'_q$ . Formally looking for zeros of the derivative of  $\varepsilon'_\xi$  with respect to  $f$  is impractical. Instead we find the local maxima by numerical maximization on the four pieces whereon  $q'$  is a smooth function of  $f$ .

That maximum can be minimized by a straightforward hill-climbing<sup>13</sup> starting from  $\Gamma = \frac{1}{10}$ ,  $\kappa = \frac{1}{2}$ ,  $\lambda = \frac{1}{4}$ , and  $\mu = 3$ . We obtain the values

$$\begin{aligned} \Gamma_{L^2C} &\approx 0.10007\,61614\,69941\,46538\,73178\,74111\,71965\,58348^-, \\ \kappa_{L^2C} &\approx 0.49999\,99381\,08574\,04775\,14291\,72928\,30652\,88838^-, \\ \lambda_{L^2C} &\approx 0.25000\,00000\,00145\,58487\,81104\,01052\,77249\,27607^+, \\ \mu_{L^2C} &\approx 3.00074\,62871\,20756\,72280\,51404\,24030\,90919\,8768^-, \\ \max_y |\varepsilon'_\xi| &\approx 2.61568\,73856\,96087\,03169\,94140\,65268\,27137\,2496^- \cdot 10^{-6}, \end{aligned}$$

in hexadecimal

$$\begin{aligned} \Gamma_{L^2C} &\approx {}_{16}\text{E99.70D0 DEAD BEEF}, \\ \kappa_{L^2C} &\approx {}_{16}\text{E99.70D0 DEAD BEEF}, \\ \lambda_{L^2C} &\approx {}_{16}\text{E99.70D0 DEAD BEEF}, \\ \mu_{L^2C} &\approx {}_{16}\text{E99.70D0 DEAD BEEF}, \\ \max_y |\varepsilon'_\xi| &\approx {}_{16}\text{1.DEAD BEEF} \cdot 2^{-19}, \end{aligned}$$

$\log_2 \max_y |\varepsilon'_\xi| \approx -18.54$ : this optimization gains two bits. The resulting  $\varepsilon'_\xi$  is remarkably equioscillating, as can be seen in figure 1.

With the rewriting (8), the constants  $1/\sqrt{12}$  and 4 should be replaced by

$$\begin{aligned} \sqrt{\frac{1 - \lambda_{L^2C}\mu_{L^2C}}{\mu_{L^2C}}} &\approx 0.28853\,15115\,62316\,71905\,38451\,44194\,38406\,3 \\ &\approx {}_{16}\text{E99.70D0 DEAD BEEF} \end{aligned}$$

and

$$\begin{aligned} \frac{1}{1 - \lambda_{L^2C}\mu_{L^2C}} &\approx 4.00298\,73779\,31697\,18250\,67433\,26901\,80421 \\ &\approx {}_{16}\text{E99.70D0 DEAD BEEF} \end{aligned}$$

respectively.

Note that a similar optimization could be applied to the rational method; however, it would not unconditionally be free: changing the 2 in the denominator turns an addition into a multiplication, and inserting additional constants adds more operations. Whether this hinders performance depends on the architecture. In any case, the optimization can scarcely gain more than two bits; such an optimized rational method would still have double the error of the optimized irrational method.

### 3 Rounded approximation

The number  $x$  is obtained from  $\xi$  by rounding to  $\left\lfloor \frac{p}{3} \right\rfloor$  bits.

<sup>13</sup>It is plausible that some variation of Pemež's algorithm could be used here, much like it can be adapted to rational functions; since the hill-climbing converged satisfactorily, and did so much faster than we were writing this document, we have not investigated this.

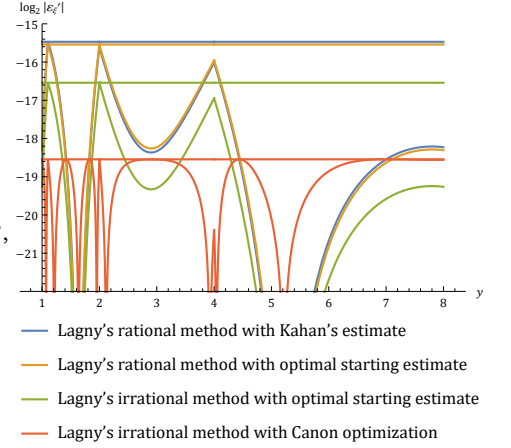


Figure 1. Error of  $\xi'$  for various methods.

### Directed rounding toward zero

An easy solution is to zero all but the most significant  $\lfloor \frac{p}{3} \rfloor$  bits of  $\xi$ .

The resulting relative error  $\left| \frac{x}{\xi} - 1 \right|$  is greatest when the zeroed bits are all 1 and the remaining bits (except for the leading 1) are all 0; this is the case when the significand of  $\xi$  is  $1 + 2^{-\lfloor \frac{p}{3} \rfloor + 1} - 2^{1-p}$ , in which case that of  $x$  is 1, so that

$$\left| \frac{x}{\xi} - 1 \right| \leq 1 - \frac{1}{1 + 2^{-\lfloor \frac{p}{3} \rfloor + 1} - 2^{1-p}} < 2^{-\lfloor \frac{p}{3} \rfloor + 1} = 2^{-16}.$$

For the error of  $x$  as an approximation of the cube root,

$$\varepsilon_x := \frac{x}{\sqrt[3]{y}} - 1,$$

we have the bound  $|\varepsilon_x| < (1 + |\varepsilon_\xi|)(1 + 2^{-16}) - 1$ .

### Rounding to nearest

An alternative which has similar performance on some architectures, but half the error, is to use Dekker's method to round to the nearest number with the desired number of bits; see [Dek71, pp. 235, 241] in *mul12*:

$$\varpi := \llbracket \xi(2^{p-\lfloor \frac{p}{3} \rfloor} + 1) \rrbracket, \quad x := \llbracket \llbracket \xi - \varpi \rrbracket + \varpi \rrbracket.$$

The rounding error is then bounded by  $2^{-\lfloor \frac{p}{3} \rfloor}$ ; we have  $|\varepsilon_x| < (1 + |\varepsilon_\xi|)(1 + 2^{-17}) - 1$ .

## 4 High order iteration

We compute the faithfully rounded result  $r_0$  as the correctly rounded sum of the rounded approximation and a correction term,

$$r := x + \Delta, \quad r_0 := \llbracket r \rrbracket, \tag{9}$$

where the correction term  $\Delta$  is that of a high-order root finding method. As usual, we call the infinite-precision correction term  $\Delta'$ , and  $r' := x + \Delta'$ . The truncation error is

$$\varepsilon_r' := \frac{\llbracket x + \Delta' \rrbracket}{\sqrt[3]{y}} - 1.$$

The rounding error on the correction term is  $\delta := \frac{\Delta}{\Delta'} - 1$ . The error of  $r$  is thus

$$\begin{aligned} \varepsilon_r &:= \frac{x + \Delta}{\sqrt[3]{y}} - 1 \\ &= \frac{x + \Delta'(1 + \delta)}{\sqrt[3]{y}} - 1 \\ &= \frac{x + \Delta' + (x + \Delta' - x)\delta}{\sqrt[3]{y}} - 1 \\ &= \varepsilon_r' + (\varepsilon_r' - \varepsilon_x)\delta, \end{aligned} \tag{10}$$

and that of the faithfully-rounded result  $r_0$  is  $\varepsilon_r(1 + v)$  for some  $|v| \leq u$ .

It is easy to make  $\varepsilon_r'$  negligible by increasing the order of the method; the main contribution to misrounding is then  $\delta\varepsilon_x$ : this is why  $\varepsilon_x$  needed to be kept low in step 2.

We now compute bounds for  $\varepsilon_r'$  and  $\delta$  for two different methods. In the interest of brevity, we have not considered the generalized Lagny irrational methods here; both the rational and quadratic irrational methods of order 4 have overly high truncation error, and the quadratic irrational method of order 5 is costly to evaluate.

TODO(egg): the quadratic irrational method of order 5 *looks* costly, but who knows. It certainly would be tedious to analyse its rounding errors, as it has differences of inexact positive terms—those are far from cancellation though, so the method is likely still viable.

The resulting bound on the numerator is approximately  $4.33u$ , an improvement of  $\frac{2}{3}u$  over the naïve bound of  $5u$ . We may build a similar  $\frac{2}{3}$ -improvement for the denominator.

Expression	Bound on the rounding error
$\llbracket 15x^3 \rrbracket + \llbracket 51y \rrbracket$	$\gamma_1$
$\llbracket \llbracket 15x^3 \rrbracket + \llbracket 51y \rrbracket \rrbracket$	$\gamma_2$
$\llbracket \llbracket \llbracket 15x^3 \rrbracket + \llbracket 51y \rrbracket \rrbracket x^3 \rrbracket$	$\gamma_3$
$\llbracket \llbracket \llbracket 15x^3 \rrbracket + \llbracket 51y \rrbracket \rrbracket x^3 \rrbracket + \llbracket 15 \llbracket y^2 \rrbracket \rrbracket$	$\iota^6 \frac{66\gamma_3 + 15\gamma_2}{81} < \frac{1}{81} \gamma_{228\iota^6}$
$\llbracket \llbracket \llbracket \llbracket 15x^3 \rrbracket + \llbracket 51y \rrbracket \rrbracket x^3 \rrbracket + \llbracket 15 \llbracket y^2 \rrbracket \rrbracket \rrbracket$	$\frac{1}{81} \gamma_{228\iota^6 + 81}$
$\llbracket x^2 \llbracket \llbracket \llbracket \llbracket 15x^3 \rrbracket + \llbracket 51y \rrbracket \rrbracket x^3 \rrbracket + \llbracket 15 \llbracket y^2 \rrbracket \rrbracket \rrbracket \rrbracket$	$\frac{1}{81} \gamma_{228\iota^6 + 2 \cdot 81}$

This is a bound of about  $4.81u$  on the denominator, a more modest improvement of  $\frac{5}{27}u$  over our earlier  $5u$ . Overall, we get the bound

$$\delta < \frac{1 + \frac{1}{27} \gamma_{10\iota^9 + (2 \cdot 26 + 1)\iota^6 + 2 \cdot 27} + \gamma_1}{1 - \frac{1}{81} \gamma_{228\iota^6 + 2 \cdot 81}} - 1 < \frac{1 + \frac{1}{27} \gamma_{10\iota^9 + (2 \cdot 26 + 1)\iota^6 + 3 \cdot 27}}{1 - \frac{1}{81} \gamma_{228\iota^6 + 2 \cdot 81}} - 1,$$

approximately  $10.14u$ .

REMOVE BEFORE FLIGHT: give the resulting bound for  $\varepsilon_r$  using Lagny's irrational method with Canon optimization.

#### Sixth order rational

An alternative is the Lagny–Schröder rational method of order 6:

$$\Delta = \frac{\llbracket \llbracket x(y - x^3) \rrbracket \llbracket \llbracket \llbracket 5x^3 \rrbracket + \llbracket 17y \rrbracket \rrbracket x^3 \rrbracket + \llbracket 5 \llbracket y^2 \rrbracket \rrbracket \rrbracket}{\llbracket \llbracket \llbracket 7x^3 \rrbracket + \llbracket 42y \rrbracket \rrbracket x^6 \rrbracket + \llbracket \llbracket 30x^3 \rrbracket + 2y \rrbracket \llbracket y^2 \rrbracket \rrbracket \rrbracket},$$

where  $x^3$  is exact thanks to the trailing 0s of  $x$ ,  $\llbracket x^6 \rrbracket$  is correctly rounded because it is computed as the square of  $x^3$ , and  $y - x^3$  is exact by Sterbenz's lemma.

Here we have  $|\varepsilon'_r| < 2^{-100}$  if  $|\varepsilon_x| < 2^{-14}$ , so the truncation error is even more negligible. For rounding error, the maximum bound on the error of the sums gives us a naïve bound of  $\gamma_6$  on the numerator,  $\gamma_5$  on the denominator, overall

$$\delta < \frac{1 + \gamma_7}{1 - \gamma_5} - 1 < \frac{\gamma_{12}}{1 - \gamma_5} \approx 12u.$$

The following ~~trivial~~ tighten the bounds on the sums in the numerator and the denominator in the same way as for the fifth order method.

Expression	Bound on the rounding error
$\llbracket 5x^3 \rrbracket + \llbracket 17y \rrbracket$	$\gamma_1$
$\llbracket \llbracket 5x^3 \rrbracket + \llbracket 17y \rrbracket \rrbracket$	$\gamma_2$
$\llbracket \llbracket \llbracket 5x^3 \rrbracket + \llbracket 17y \rrbracket \rrbracket x^3 \rrbracket$	$\gamma_3$
$\llbracket \llbracket \llbracket 5x^3 \rrbracket + \llbracket 17y \rrbracket \rrbracket x^3 \rrbracket + \llbracket 5 \llbracket y^2 \rrbracket \rrbracket$	$\iota^6 \frac{22\gamma_3 + 5\gamma_2}{27} < \frac{1}{27} \gamma_{76\iota^6}$
$\llbracket \llbracket \llbracket \llbracket 5x^3 \rrbracket + \llbracket 17y \rrbracket \rrbracket x^3 \rrbracket + \llbracket 5 \llbracket y^2 \rrbracket \rrbracket \rrbracket$	$\frac{1}{27} \gamma_{76\iota^6 + 27}$
$\llbracket x(y - x^3) \rrbracket \llbracket \llbracket \llbracket \llbracket 5x^3 \rrbracket + \llbracket 17y \rrbracket \rrbracket x^3 \rrbracket + \llbracket 5 \llbracket y^2 \rrbracket \rrbracket \rrbracket$	$\frac{1}{27} \gamma_{76\iota^6 + 2 \cdot 27}$
$\llbracket \llbracket x(y - x^3) \rrbracket \llbracket \llbracket \llbracket \llbracket 5x^3 \rrbracket + \llbracket 17y \rrbracket \rrbracket x^3 \rrbracket + \llbracket 5 \llbracket y^2 \rrbracket \rrbracket \rrbracket \rrbracket$	$\frac{1}{27} \gamma_{76\iota^6 + 3 \cdot 27}$

This bound is approximately  $5.81u$  on the numerator.

Expression	Bound on the rounding error
$\llbracket 30x^3 \rrbracket + 2y$	$t^3 \frac{30u}{32} < \gamma_{\frac{10}{32}} t^3$
$\llbracket \llbracket 30x^3 \rrbracket + 2y \rrbracket$	$\gamma_{\frac{10}{32}} t^3 + 1$
$\llbracket \llbracket 30x^3 \rrbracket + 2y \rrbracket \llbracket y^2 \rrbracket$	$\gamma_{\frac{10}{32}} t^3 + 2$
$\llbracket \llbracket \llbracket 30x^3 \rrbracket + 2y \rrbracket \llbracket y^2 \rrbracket \rrbracket$	$\gamma_{\frac{10}{32}} t^3 + 3$
$\llbracket \llbracket \llbracket 7x^3 \rrbracket + \llbracket 42y \rrbracket \rrbracket \llbracket x^6 \rrbracket \rrbracket + \llbracket \llbracket \llbracket 30x^3 \rrbracket + 2y \rrbracket \llbracket y^2 \rrbracket \rrbracket$	$t^6 \frac{49\gamma_4 + 32\gamma_{\frac{10}{32}} t^3 + 3}{81} < \frac{1}{81} \gamma_{10t^9 + 282t^6}$
$\llbracket \llbracket \llbracket 7x^3 \rrbracket + \llbracket 42y \rrbracket \rrbracket \llbracket x^6 \rrbracket \rrbracket + \llbracket \llbracket \llbracket 30x^3 \rrbracket + 2y \rrbracket \llbracket y^2 \rrbracket \rrbracket$	$\frac{1}{81} \gamma_{10t^9 + 282t^6 + 81}$

This bound is approximately  $4.61u$  on the denominator. We thus have, for the method of order 6,

$$\delta < \frac{1 + \frac{1}{27} \gamma_{76t^6 + 3 \cdot 27} + \gamma_1}{1 - \frac{1}{81} \gamma_{10t^9 + 282t^6 + 81}} < \frac{1 + \frac{1}{27} \gamma_{76t^6 + 4 \cdot 27}}{1 - \frac{1}{81} \gamma_{10t^9 + 282t^6 + 81}}$$

approximately  $11.42u$ .

REMOVE BEFORE FLIGHT: give the resulting bound for  $\varepsilon_r$  using Lagny’s irrational method with Canon optimization.

### The unreasonable effectiveness of the sixth order rational method

The bound on the rounding error in the sixth order rational method is a little more than  $1u$  worse than the one on its fifth order rational method; this is to be expected, as it involves more calculations and does not benefit an exact multiplication by 16. Accordingly, since truncation error is negligible before rounding error for both, one might expect that the fifth order method would be superior. Numerical experiments suggest otherwise: using the method of order 6 leads to 4.33 misroundings per million, whereas the method of order 5 leads to 4.43 per million.

The bounds on  $\delta$  may need to be tightened further; for instance, we have not taken any account of rounding errors systematically compensating each other, the error induced by  $\llbracket y^2 \rrbracket$  is the same in the numerator and the denominator, but we bound it as if it had opposite signs. The explanation may however lie elsewhere; the larger number of operations involved for the overall maximal rounding error may lead to lower average error.

Whatever the reason may be, the method of order 6 appears to be preferable should one wish to implement a faithful cube root—they are about equally fast on modern architectures. However, the poorer bound on the maximal error means that we must instead use the method of order 5 in our correctly-rounded cube root: its better bound means that we go through the “potential misrounding” path less often.

## Correct rounding

We have  $r = \sqrt[3]{y}(1 + \varepsilon_r)$ , thus

$$\sqrt[3]{y} \in \left[ \frac{r}{1 + \bar{\varepsilon}_r}, \frac{r}{1 - \bar{\varepsilon}_r} \right] =: \mathcal{I},$$

where  $\bar{\varepsilon}_r$  is a bound for  $\varepsilon_r$ . Consider the *ties*, i.e., the number halfway between  $r_0$  and its binary64 successor and the number halfway between  $r_0$  and its predecessor. If  $\mathcal{I}$  contains neither of the ties,  $r_0 = \llbracket \sqrt[3]{y} \rrbracket$ : the faithful method returned a correct result.

This criterion, slightly weakened, may be determined as follows. The difference between  $r_0 = \llbracket r \rrbracket$  and the unrounded  $r$  can be computed as described in [Dek71, p. 224],  $r_1 := x - r_0 + \Delta$ , evaluated as written (both operations are exact). The potential other candidate for  $\llbracket \sqrt[3]{y} \rrbracket$  may be computed as  $\tilde{r} := \llbracket r_0 + 2r_1 \rrbracket$ , which is  $r_0$  only if  $r_1$  is below a quarter-unit in the last place—in which case the small size of  $\mathcal{I}$  ensures that it does not contain a tie. If  $\tilde{r} \neq r_0$ , we must ascertain whether the tie lies in  $\mathcal{I}$ ,

$$\frac{\tilde{r} + r_0}{2} \in \left[ \frac{r}{1 + \bar{\varepsilon}_r}, \frac{r}{1 - \bar{\varepsilon}_r} \right],$$

subtracting  $r_0$  to get rid of the unrepresentable  $r$  in the right-hand side,

$$\frac{\tilde{r} - r_0}{2} \in \left[ \frac{r_1 - \bar{\varepsilon}_r r_0}{1 + \bar{\varepsilon}_r}, \frac{r_1 + \bar{\varepsilon}_r r_0}{1 - \bar{\varepsilon}_r} \right],$$

subtracting  $r_1$  to remove the cancellation,

$$\frac{\tilde{r} - r_0}{2} - r_1 \in \left[ -\bar{\varepsilon}_r \frac{r_0 + r_1}{1 + \bar{\varepsilon}_r}, \bar{\varepsilon}_r \frac{r_0 + r_1}{1 - \bar{\varepsilon}_r} \right].$$

The left-hand-side may be computed as written; however the bounds of the interval above are not representable. We must relax them a little,

$$\frac{\tilde{r} - r_0}{2} - r_1 \in \left[ -\frac{\bar{\varepsilon}_r}{1 + \bar{\varepsilon}_r} (1 + u) r_0, \frac{\bar{\varepsilon}_r}{1 - \bar{\varepsilon}_r} (1 + u) r_0 \right].$$

On an architecture where multiplications with directed rounding are not more expensive when the surrounding computation uses the `roundTiesToEven` rounding-direction attribute, these bounds may be used directly, provided that the constants therein have been rounded toward their respective signs, and the multiplication by  $r_0$  be similarly rounded.

If we restrict ourselves to `roundTiesToEven` at runtime, we must relax the bounds some more,

$$\frac{\tilde{r} - r_0}{2} - r_1 \in \left[ \left\lfloor \left\lfloor -\frac{\bar{\varepsilon}_r}{1 + \bar{\varepsilon}_r} \left( 1 + \frac{2u}{1 - u} \right) \right\rfloor r_0 \right\rfloor, \left\lceil \left\lceil \frac{\bar{\varepsilon}_r}{1 - \bar{\varepsilon}_r} \left( 1 + \frac{2u}{1 - u} \right) \right\rceil r_0 \right\rceil \right].$$

We may widen this interval slightly into one that is symmetric about 0, thus requiring only one comparison

$$|\tilde{r} - r_1| \leq \llbracket \tau r_0 \rrbracket,$$

with

$$\begin{aligned} \tau &:= \left\lceil \frac{\bar{\varepsilon}_r}{1 - \bar{\varepsilon}_r} \left( 1 + \frac{2u}{1 - u} \right) \right\rceil_+ \\ &= {}_{16}\text{DEAD BEEF}, \end{aligned}$$

where we have used the constants for Lagny's irrational method with Canon optimization, followed by the method of fifth order.

If this inequality holds, there may be a misrounding. The correctly-rounded result may then readily be computed using the ordinary cube root algorithm, described, *e.g.*, in Lagny's [Fan97, pp. 286 sqq.], used with binary digits.

## Extracting a digit

This method extracts a single digit of  $\llbracket \sqrt[3]{y} \rrbracket_0$  at each step. Since  $\mathcal{I}_\pm$  is small enough that it cannot contain both a floating-point number and a tie; thus if  $\mathcal{I}_\pm$  contains a tie, we know that  $\llbracket \sqrt[3]{y} \rrbracket_0 = \llbracket r \rrbracket_0 =: a$ ; we thus have the first 53 bits already, as

$$a = \min(r_0, \tilde{r}).$$

Correct rounding may be achieved by extracting a single additional bit: the number is in excess of the tie—and thus must be rounded up—if and only if that bit is 1, because, as remarked in [LMoo, p. 15], there are no halfway cases for the cube root.

In the ordinary cube root method, our remainder is  $\rho_{53} := y - a^3$ , and the next bit is 1 if and only if the next remainder would be positive with that bit,

$$\rho_{54|1} := \rho_{53} - 3a^2b - 3ab^2 - b^3 \geq 0,$$

where  $b$  is the power of two corresponding to this bit (the difference between  $a$  and the tie).



Using Veltkamp's<sup>14</sup> algorithm from [Dek71, p. 234] to express  $a^2$  as  $a_0^2 + a_1^2$  with  $a_0^2 = \llbracket a^2 \rrbracket$ , the remainder is

$$\rho_{53} = y - (a_0^2 + a_1^2)a;$$

with two more applications of Veltkamp's algorithm,

$$\rho_{53} = y - a_{00}^3 - a_{01}^3 - a_{10}^3 - a_{11}^3,$$

where the first subtraction is representable by Sterbenz's lemma. The remainder cannot have more digits than two plus twice the number of digits computed for the cube root, *i.e.*,  $2 \cdot 53 + 2$  significant bits, for otherwise the "digit"  $2 = {}_210$  would fit in the place  $b$ . By repeated exact subtractions we may thus express it as

$$\rho_{53} = \rho_{53;0} + \rho_{53;1} + \rho_{53;2},$$

where  $\rho_{53;i} = \llbracket \rho_{53;i} + \rho_{53;i+1} \rrbracket$ , with ample room to spare. The terms  $3a^2b$ ,  $3ab^2$ , and  $b^3$  may then be subtracted exactly while retaining triple precision: they add at most three significant digits. The first two of these terms should be split into representable parts,

$$3a^2b = 2a_0^2b + a_0^2b + 2a_1^2b + a_1^2b,$$

$$3ab^2 = 2ab^2 + ab^2.$$

The sign of  $\rho_{54|1}$  may then be checked.

---

<sup>14</sup>We have not been able to find a copy of *RC-Informatie*, wherein Veltkamp's work was published; we follow Dekker's account of thereof.

# Appendices

## A Proof of the equivalence of Lagny's rational method and Schröder's ( $A_{p-1}^0$ )

We now prove the proposition from part I, which, substituting the definition of the generalized Lagny rational method, is that

$$x + \frac{1}{2}a = a + (p-1) \frac{(1/f)^{(p-2)}(a)}{(1/f)^{(p-1)}(a)} =: \psi(a)$$

if  $x$  is the root of  $E_1$ . This can be expressed equivalently, and more conveniently for the proof, as

$$a + \Delta = a + (p-1) \frac{(1/f)^{(p-2)}(a)}{(1/f)^{(p-1)}(a)} =: \psi(a)$$

if  $\Delta$  is the root of  $\tilde{E}_1$ .

**Proof.** Let  $\tilde{E}_p = d_0\Delta^p + \dots + d_p$ ,  $\tilde{E}_{p-1} = e_0\Delta^{p-1} + \dots + e_{p-1}$ , and As shown in [Hou70, pp. 52–54], the polynomial remainders  $E_k$  are given up to a constant factor by [Hou70, p. 19] equation (23), *i.e.*, for some  $\alpha$ ,

$$\frac{\tilde{E}_n}{\alpha_n} = \det \begin{pmatrix} (\tilde{E}_p)_{p-1-n} \\ (\tilde{E}_{p-1})_{p-n} \end{pmatrix},$$

where the expression on the right-hand side is the *bigradient* defined in [Hou68] (3.4) or [Hou70, p. 19] (20),

$$\frac{\tilde{E}_n}{\alpha_n} = \det \begin{pmatrix} d_0 & d_1 & d_2 & \dots & d_{2(p-n)-3} & \Delta^{p-n-2}E_p \\ & d_0 & d_1 & \dots & d_{2(p-n)-4} & \Delta^{p-n-3}E_p \\ & & \ddots & & \vdots & \\ & & & d_0 & d_1 & \dots & d_{p-n-1} & \Delta^0 E_p \\ & & & & e_0 & \dots & e_{p-n-2} & \Delta^0 E_{p-1} \\ & & & & e_0 & e_1 & \dots & e_{p-n-1} & \Delta^1 E_{p-1} \\ & & & & & \ddots & & \vdots & \\ & & & & & & e_{2(p-n)-4} & \Delta^{p-n-2}E_{p-1} \\ e_0 & e_1 & e_2 & \dots & e_{2(p-n)-3} & \Delta^{p-n-1}E_{p-1} \end{pmatrix} =: \det \tilde{E}_k,$$

with  $p-1-n$  rows of the  $d_k$  and  $p-n$  of the  $e_k$ , where  $d_k := 0$  for  $k > p$ , and  $e_k := 0$  for  $k > p-1$ .

In particular, for  $n = 1$ ,

$$\tilde{E}_1 = \begin{pmatrix} d_0 & d_1 & d_2 & \dots & d_{p-2} & d_{p-1} & d_p & & 0 & \Delta^{p-3}E_p \\ & d_0 & d_1 & \dots & d_{p-3} & d_{p-2} & d_{p-1} & d_p & & \Delta^{p-4}E_p \\ & & \ddots & & & & & & \ddots & \vdots \\ & & & d_0 & d_1 & d_2 & \dots & d_{p-2} & \Delta^0 E_p \\ & & & & e_0 & e_1 & \dots & e_{p-3} & \Delta^0 E_{p-1} \\ & & & & e_0 & e_1 & e_2 & \dots & e_{p-2} & \Delta^1 E_{p-1} \\ & & & & & \ddots & & & \ddots & \vdots \\ & & & & & & e_{p-3} & e_{p-2} & e_{p-1} & \Delta^{p-3}E_{p-1} \\ e_0 & e_1 & e_2 & \dots & e_{p-2} & e_{p-1} & & & 0 & \Delta^{p-2}E_{p-1} \end{pmatrix},$$

with  $p-2$  rows of the  $d_k$  and  $p-1$  of the  $e_k$ . Observe that since the value of  $\Delta$  used in the rational method is the root of  $\tilde{E}_1$ , for that value of  $\Delta$ ,  $\det \tilde{E}_1 = 0$ , *i.e.*,  $\tilde{E}_1$  is singular.

**Lemma.** The matrix  $\tilde{E}_1$  is singular if and only if  $C(a + \Delta)$  is singular, where

$$C(\Psi) := \begin{pmatrix} \Psi - a & f_0 & & \mathbf{0} \\ -1 & f_1 & f_0 & \\ 0 & f_2 & & \ddots \\ \vdots & \vdots & \ddots & f_0 \\ 0 & f_{p-1} & \cdots & f_2 & f_1 \end{pmatrix}$$

and  $f_k := \frac{f^{(k)}(a)}{k!}$ .

**Proof.** Observe that by the definition of  $\tilde{E}_p$  and  $\tilde{E}_{p-1}$ ,  $d_k = f_{p-k}$  and  $e_k = f_{p-1-k}$ , so that

$$\tilde{E}_1 = \begin{pmatrix} f_p & f_{p-1} & f_{p-2} & \cdots & f_2 & f_1 & f_0 & \mathbf{0} & \Delta^{p-3}E_p \\ & f_p & f_{p-1} & \cdots & f_3 & f_2 & f_1 & f_0 & \Delta^{p-4}E_p \\ & & \ddots & & & & & & \vdots \\ \mathbf{0} & & & f_p & f_{p-1} & f_{p-2} & \cdots & f_2 & \Delta^0 E_p \\ & & & f_{p-1} & f_{p-2} & f_{p-3} & \cdots & f_2 & \Delta^0 E_{p-1} \\ & & & f_{p-1} & f_{p-2} & f_{p-3} & \cdots & f_1 & \Delta^1 E_{p-1} \\ & & \ddots & & & & & \vdots \\ f_{p-1} & f_{p-2} & f_{p-1} & \cdots & f_2 & f_1 & f_0 & \mathbf{0} & \Delta^{p-3}E_{p-1} \\ & f_{p-2} & f_{p-1} & \cdots & f_1 & f_0 & & \mathbf{0} & \Delta^{p-2}E_{p-1} \end{pmatrix}.$$

Note that  $\tilde{E}_p - \tilde{E}_{p-1} = f_p \Delta^p$ . Subtracting the penultimate row from the first, the antepenultimate from the second, etc., the determinant is that of

$$\begin{pmatrix} f_p & & & & \mathbf{0} & \Delta^{2p-3}f_p \\ & f_p & & & & \Delta^{2p-4}f_p \\ & & \ddots & & & \vdots \\ \mathbf{0} & & & f_p & & \Delta^p f_p \\ & & & f_{p-1} & f_{p-2} & \cdots & f_2 & \Delta^0 E_{p-1} \\ & & & f_{p-1} & f_{p-2} & f_{p-3} & \cdots & f_1 & \Delta^1 E_{p-1} \\ & & \ddots & & & & \vdots \\ f_{p-1} & f_{p-2} & f_{p-1} & \cdots & f_2 & f_1 & f_0 & \mathbf{0} & \Delta^{p-3}E_{p-1} \\ & f_{p-2} & f_{p-1} & \cdots & f_1 & f_0 & & \mathbf{0} & \Delta^{p-2}E_{p-1} \end{pmatrix}.$$

Since  $\det \tilde{E}_1$  is a polynomial of degree 1, all terms divisible by  $\Delta^2$  must cancel out in the Laplace expansion<sup>15</sup> of the determinant of the above matrix in the last column, so that

$$\det \tilde{E}_1 = \pm(f_0 + f_1 \Delta) \det \mathbf{A} \mp f_0 \Delta \det \mathbf{B},$$

where

$$\mathbf{A} := \left( \begin{array}{c|ccc} f_p \mathbb{1}_{p-2} & & \mathbf{0} & \\ \hline & f_{p-2} & \cdots & f_1 \\ & f_{p-3} & \cdots & f_0 \\ & \vdots & \ddots & \\ \mathbf{X} & f_1 & f_0 & \mathbf{0} \end{array} \right) \text{ and } \mathbf{B} := \left( \begin{array}{c|ccc} f_p \mathbb{1}_{p-2} & & \mathbf{0} & \\ \hline & f_{p-1} & \cdots & f_2 \\ & f_{p-3} & \cdots & f_1 \\ & \vdots & \ddots & \\ \mathbf{Y} & f_1 & f_0 & \mathbf{0} \end{array} \right),$$

where the matrices  $\mathbf{X}$  and  $\mathbf{Y}$  are of no consequence for the determinant:

$$\frac{\det \tilde{E}_1}{f_p^{p-2}} = \pm(f_0 + f_1 \Delta) \det \mathbf{A}' \mp f_0 \Delta \det \mathbf{B}',$$

where the matrices  $\mathbf{A}'$  and  $\mathbf{B}'$  are the bottom right blocks  $\mathbf{A}$  and  $\mathbf{B}$  respectively.

<sup>15</sup>[Lap72, pp. 294–304], reprinted in [Lap78, vol. 8, pp. 395 sqq.].

A few Laplace expansions of

$$\det \mathbf{C}(a + \Delta) = \det \begin{pmatrix} \Delta & f_0 & & & \mathbf{0} \\ -1 & f_1 & f_0 & & \\ 0 & f_2 & & \ddots & \\ \vdots & \vdots & & \ddots & f_0 \\ 0 & f_{p-1} & \cdots & f_2 & f_1 \end{pmatrix}$$

finish the proof; along the first column,

$$\det \mathbf{C}(a + \Delta) = \Delta \det \begin{pmatrix} f_1 & f_0 & & & \mathbf{0} \\ f_2 & & \ddots & & \\ \vdots & & \ddots & & f_0 \\ f_{p-1} & \cdots & f_2 & f_1 & \end{pmatrix} + \det \begin{pmatrix} f_0 & & & & \mathbf{0} \\ f_2 & f_1 & f_0 & & \\ f_3 & f_2 & & \ddots & \\ \vdots & & & \ddots & f_0 \\ f_{p-1} & f_{p-2} & & f_2 & f_1 \end{pmatrix},$$

along the first row of both matrices,

$$\det \mathbf{C}(a + \Delta) = \Delta(f_1 \det \mathbf{\neg A}' - f_0 \det \mathbf{\neg B}') + f_0 \det \mathbf{\neg A}',$$

where  $\mathbf{\neg Z}'$  is obtained by reversing the order of the columns of the transpose of  $\mathbf{Z}$ , so that

$$\det \mathbf{C}(a + \Delta) = \pm \frac{\det \tilde{\mathbf{E}}_1}{f_p^{p-2}}.$$

□

The proposition follows from the lemma and theorem 4.4.2 from [Hou70, p. 169]:  $\psi(a)$  is Householder's (14) with  $g \equiv 1$ ; for that value of  $g$ , theorem 4.4.2 states that (14) is the solution of (12) from the same page, which is  $\det \mathbf{C}(\psi(a)) = 0$ . By the lemma, for the value of  $\Delta$  in the rational method,  $a + \Delta$  solves that equation. □

## B Some applications of the generalized Lagny methods

As previously mentioned, the rational methods are well known; we list them solely so that they may be compared with the irrational ones.

### 1 Formulæ for the cube root

Let  $f(z) := y - z^3$ , and  $a$  be the starting estimate for the root of  $f$ , i.e., for  $\sqrt[3]{y}$ . Let  $b := f(a) = y - a^3$ . Let  $\varepsilon := \frac{a}{\sqrt[3]{y}} - 1$ . We have the following formulæ.

#### 1.1 Rational methods

The rational methods of orders 2 through 5 are given below. The first two of these are Newton's method and Lagny's rational method.

Iteration	Asymptotic relative error
$a + \frac{b}{3a^2}$	$\varepsilon^2 + \mathcal{O}(\varepsilon^3)$
$a + \frac{ab}{3a^3+b}$	$\frac{2}{3}\varepsilon^3 + \mathcal{O}(\varepsilon^4)$
$a + \frac{3ab(3a^3+b)}{26a^6+18a^3b+b^2}$	$\frac{1}{3}\varepsilon^4 + \mathcal{O}(\varepsilon^5)$
$a + \frac{b(27a^6+18a^3b+b^2)}{81a^8+81a^5b+15a^2b^2}$	$\frac{1}{9}\varepsilon^5 + \mathcal{O}(\varepsilon^6)$

#### 1.2 Quadratic irrational methods

The quadratic irrational methods of orders 3 through 5 are given below. The first of these is Lagny's irrational method, for which we give the form from [Fan91a].

Iteration	Asymptotic relative error
$\frac{3a^2+\sqrt{9a^4+12ab}}{6a} = \frac{1}{2}a + \sqrt{\frac{1}{4}a^2 + \frac{b}{3a}}$	$-\frac{1}{3}\varepsilon^3 + \mathcal{O}(\varepsilon^4)$
$\frac{3a^3-b+\sqrt{81a^6+90a^3b+b^2}}{12a^2}$	$-\frac{1}{9}\varepsilon^4 + \mathcal{O}(\varepsilon^5)$
$\frac{-5a^3b+\sqrt{3a^6(108a^6+108a^3b-5b^2)}}{18a^5-2a^2b}$	$-\frac{1}{18}\varepsilon^5 + \mathcal{O}(\varepsilon^6)$

## 2 Formulæ for the fifth root

Let  $f(z) := y - z^5$ , and  $a$  be the starting estimate for the root of  $f$ , i.e., for  $\sqrt[5]{y}$ . Let  $b := f(a) = y - a^5$ . Let  $\varepsilon := \frac{a}{\sqrt[5]{y}} - 1$ . We have the following formulæ.

### 2.1 Rational methods

The rational methods of orders 2 through 5 are given below. The first of these is Newton's method, the second is Halley's rational method, the last is Lagny's rational method.

Iteration	Asymptotic relative error
$a + \frac{b}{5a^4}$	$2\varepsilon^2 + \mathcal{O}(\varepsilon^3)$
$a + \frac{ab}{5a^5+2b}$	$2\varepsilon^3 + \mathcal{O}(\varepsilon^4)$
$a + \frac{ab(5a^5+2b)}{25a^{10}+20a^5b+2b^2}$	$\varepsilon^4 + \mathcal{O}(\varepsilon^5)$
$a + \frac{ab(25a^{10}+20a^5b+2b^2)}{125a^{15}+150a^{10}b+40a^5b^2+b^3}$	$-\frac{1}{5}\varepsilon^5 + \mathcal{O}(\varepsilon^6)$

### 2.2 Quadratic irrational methods

The quadratic irrational methods of orders 2 through 5 are given below. The first of these is Halley's irrational method, for which we give the form from [Hal94, p. 141], the last is Lagny's quadratic irrational method.

Iteration	Asymptotic relative error
$\frac{15a^4+\sqrt{25a^8+40a^3b}}{20a^3} = \frac{3}{4}a + \sqrt{\frac{1}{16}a^2 + \frac{b}{10a^3}}$	$-\frac{1}{3}\varepsilon^3 + \mathcal{O}(\varepsilon^4)$
$\frac{5a^7-a^2b+a^2\sqrt{25a^{10}+30a^5b+b^2}}{10a^6}$	$-\varepsilon^4 + \mathcal{O}(\varepsilon^5)$
$a \frac{-7b+\sqrt{100a^{10}+100a^5b-7b^2}}{10a^5-4b}$	$-\frac{7}{10}\varepsilon^5 + \mathcal{O}(\varepsilon^6)$

### 2.3 Quartic irrational method

Lagny's quartic irrational method is given below, in the form from [Fan92] (with the misprint corrected). The quartic methods involve the solution of a quartic equation, so they are often impractical; in this case however, the equation  $E_4 = 0$  is biquadratic.

Iteration	Asymptotic relative error
$\frac{1}{2}a + \sqrt{\frac{1}{4}a^4 + \frac{b}{5a} - \frac{1}{4}a^2}$	$-\frac{1}{5}\varepsilon^5 + \mathcal{O}(\varepsilon^6)$

## 3 Formulæ for the resolution of Kepler's equation

Let  $f(E) := E - e \sin E - M$ ;  $f(E) = 0$  is Kepler's equation<sup>16</sup>, where  $E$  is the eccentric anomaly,  $M$  the mean anomaly, and  $e$  the eccentricity of an elliptic orbit.

We give the iterations for a starting estimate  $\alpha$  of the root  $E$ . Since  $E$  is an angle and  $M$  a fictitious angle, both in  $[0, 2\pi]$ , we consider absolute errors rather than relative ones. With this two-parameter transcendental function, the two-parameter asymptotics are not particularly useful. Instead we give the approximate maximal error over  $M \in [0, 2\pi]$  for several values of  $e$ , and for the specific choice of starting estimate  $\alpha = M$ , which greatly simplifies the higher-order formulæ.

The errors of the starting estimate  $\alpha = M$  are given in the table below.

<sup>16</sup>From [Kep09, pp. 295–300]. Kepler writes of this equation:

Mihi fufficit credere, folvi a priori non poſſe, propter arcus & ſinus éτερογένηται. Erranti mihi, quicunque viam monſtraverit, is erit mihi magnus Apollonius.

It appears that this honour belongs either to Lagrange, with a power series in  $e$  whose coefficients involve derivatives of powers the sine, given in [Lag71, p. 209] and reprinted in [Lag67, vol. 3, p. 117], or to Bessel, with a sine series in  $M$  whose coefficients involve his eponymous functions, given in an 1818 letter to Olbers [Bes52]; see [Col92].

Approximation of $E$	Max. error for $e = 0.2$	$e = 0.5$	$e = 0.9$	$e = 0.999$
$M$	$10^\circ 37'$	$28^\circ 34'$	$51^\circ 34'$	$57^\circ 14'$

For brevity, we write  $c := \cos \alpha$ ,  $s := \sin \alpha$ .

### 3.1 Rational methods

The rational methods of orders 2 through 4 are given below; the first two of these are Newton's method and Halley's rational method.

$$\begin{aligned}
 &\text{Iteration} \\
 &\alpha + \frac{M - \alpha + es}{1 - ec} \\
 &\alpha + \frac{2(1 - ec)(M - \alpha + es)}{2 - 4ec + e^2c^2 - e\alpha s + eMs + e^2} \\
 &\alpha + \frac{3(M - \alpha + es)(2 - 4ec + e^2c^2 - e\alpha s + eMs + e^2)}{6(1 + e^2 + 2e^2c^2 - e^3e^3 - e\alpha s + eMs) + e((\alpha - M)^2 - 18 + 4e\alpha s - 4eMs - 5e^2s^2)c}
 \end{aligned}$$

For  $\alpha = M$ , we have the following simplified formulæ.

Approximation of $E$	Max. error for $e = 0.2$	$e = 0.5$	$e = 0.9$	$e = 0.999$
$M + \frac{es}{1 - ec}$	$14' 11''$	$4^\circ 27'$	$68^\circ 32'$	$1246^\circ$
$M + \frac{2(1 - ec)es}{2 - 4ec + e^2c^2 + e^2}$	$21'' 43$	$22' 35''$	$13^\circ 07'$	$38^\circ 44'$
$M + \frac{3es(2 - 4ec + e^2c^2 + e^2)}{6(1 + 2e^2c^2 - e^3c^3 + e^2) - e(18 + 5e^2s^2)c}$	$3'' 03$	$7' 56''$	$10^\circ 30'$	$27^\circ 20'$

### 3.2 Quadratic irrational methods

The quadratic irrational methods of orders 3 and 4 are given below. For the method of order 4, the sign  $\pm$  should be positive when  $M < \pi$ , and negative when  $M > \pi$ .

$$\begin{aligned}
 &\text{Iteration} \\
 &\alpha + \frac{ec - 1 + \sqrt{(1 - ec)^2 + 2es(M - \alpha + es)}}{es} \\
 &\alpha + \frac{(\alpha - M + 2es)c - 3s \pm \sqrt{(\alpha - M)^2c^2 - 2(4ec - 3)(\alpha - M)sc + (3 - 4ec)^2s^2 + 18e(M - \alpha)s^3 + 18e^2s^4}}{2c(ec - 1) + 3es^2}
 \end{aligned}$$

For  $\alpha = M$ , we have the following simplified formulæ.

Approximation of $E$	Max. error for $e = 0.2$	$e = 0.5$	$e = 0.9$	$e = 0.999$
$M + \frac{ec - 1 + \sqrt{(1 - ec)^2 + 2e^2s^2}}{es}$	$24'' 37$	$24' 38''$	$10^\circ 54'$	$53^\circ 25'$
$M + \frac{2esc - 3s \pm \sqrt{(3 - 4ec)^2s^2 + 18e^2s^4}}{2c(ec - 1) + 3es^2}$	$3'' 06$	$8' 31''$	$14^\circ 14'$	$27^\circ 23'$

### 3.3 Cubic irrational method

For the extraction of the roots of polynomials, it makes little sense to consider the cubic irrational methods, whose computation involves the extraction of a cube root. However, when it comes to functions whose derivatives involve trigonometric lines, a cube root is no longer prohibitive. The approximation obtained from the method of degree 3 and order 4 with the starting estimate  $\alpha = M$  is given below; its error is remarkably low even at high eccentricities (an order of magnitude better than either the rational or the quadratic method of the same order for  $e \geq 0.9$ ).

Note that this formula is valid for  $M \leq \frac{\pi}{2}$ ; the corresponding expressions for  $\frac{\pi}{2} \leq M \leq \frac{3\pi}{2}$  and for some of the  $M > \frac{3\pi}{2}$  are *casus irreducibiles*<sup>17</sup>, and thus are less practically useful; it is probably best to use one of the formulæ above for the middle, and to get rid of the upper half by argument reduction.

<sup>17</sup>See [Fang7, p. 469] for a definition of the *cas irreductible* of a cubic; see [Wantzel1843] for a proof of its irreducibility.

Approximation of $E$	Max. error for $e = 0.2$	$e = 0.5$	$e = 0.9$	$e = 0.999$
$M + \frac{\frac{2c(ec-1)+es^2}{R} + \frac{R}{e} - s}{c}$	2''71	4'04''	58'33''	1°29'

In the above formula,

$$R := \sqrt[3]{3e^2sc - e^3s^3 + \sqrt{e^3c^2(8(1-ec)^3 - 3es^2(1+4ec(ec-2)) - 6e^3s^4)}}.$$

## C FMA

The overall strategy is different with FMA, since we need only round  $x$  to half the precision (26 bits), rather than a third; this is because the expression  $y - x^3$  may be computed as  $\llbracket y - (x^2)x \rrbracket$ , requiring only an exact square.

This means that  $\varepsilon'_\xi$  should ideally be somewhat less than  $2^{-26}$  or  $2^{-27}$  depending on the manner in which  $x$  is rounded. Even with Canon optimization, Lagny's irrational method cannot achieve that from  $q$ . The error of the Lagny-Schröder rational method of order 4 reaches a below  $2^{-21}$ , and that of the generalized Lagny quadratic irrational method of order 4 below  $2^{-23}$ . It seems that Canon optimization on the latter can readily bring down its error below  $2^{-24}$ , but not much further. With the methods of order 5, the errors are below  $2^{-28}$ ; little stands to be gained from optimization.

Conversely the computation of  $r_0$  from  $x$  may use a lower-order method, since it starts from 26 bits rather than 17, but still cannot add more than 53.

TODO(egg): concrete method.

## D Performance considerations

### Without FMA

In the second step, where rounding errors are not a concern, Lagny's rational method is best evaluated as

$$\xi = \left\llbracket \frac{\llbracket \llbracket q^2 \rrbracket \llbracket q^2 \rrbracket \rrbracket + \llbracket 2yq \rrbracket}{\llbracket 2 \llbracket \llbracket q^2 \rrbracket q \rrbracket + y \rrbracket} \right\rrbracket.$$

On Ivy Bridge<sup>18</sup>, the irrational method evaluated as

$$\xi = \left\llbracket \left\llbracket \left\llbracket \frac{1}{2} q' \right\rrbracket + \left\llbracket \left\llbracket \frac{\llbracket 1/\sqrt{12} \rrbracket}{q} \right\rrbracket \left\llbracket \sqrt{\llbracket \llbracket 4yq \rrbracket - \llbracket \llbracket q^2 \rrbracket \llbracket q^2 \rrbracket \rrbracket} \right\rrbracket \right\rrbracket \right\rrbracket,$$

*mutatis mutandis* with Canon optimization, is a couple cycles faster than the rational method. On Sandy Bridge, it is instead computed fastest as

$$\xi = \left\llbracket \left\llbracket \frac{1}{2} q' \right\rrbracket + \left\llbracket \sqrt{\llbracket \llbracket 4y \rrbracket - \llbracket \llbracket q^2 \rrbracket \llbracket q \rrbracket \rrbracket} \left\llbracket \frac{\llbracket 1/12 \rrbracket}{q} \right\rrbracket \right\rrbracket \right\rrbracket,$$

which is about four cycles slower than the rational method.

Surprisingly, the rounding step  $\xi \mapsto x$  is faster when performed using Veltkamp-Dekker than by directed rounding using `andpd` on Sandy Bridge, and about equally fast on Ivy Bridge. Given the increased precision, we round to the nearest 17-bit value in that step.

TODO(egg): Give some numbers for the rate of passage in the “potential mis-rounding” path, and its cost.

<sup>18</sup>We thank Peter Barfuss for running some benchmarks for us on a machine with an Ivy Bridge processor.

## With FMA

TODO(egg): decide on which methods we use and then report on performance here.

## E Rounding error analysis for the second step

TODO(egg): analysis.

## F Other rounding modes

TODO(egg): short blurb—there is not much to say here.

## G Comparison with other faithful implementations

TODO(egg): <https://twitter.com/eggleroy/status/988114872616484866> (probably only the frequency one, the maximum is noisy), compared with the FMAless faithful method so as to compare like with like; the same with the one by njuffa <https://stackoverflow.com/a/27008559> for our FMA method.

## References

- [Bab27] C. Babbage. *Table of Logarithms of the Natural Numbers, from 1 to 108000*. J. Mawman, 1827.  
eprint: [https://books.google.co.uk/books?id=64o\\_AAAAcAAJ](https://books.google.co.uk/books?id=64o_AAAAcAAJ).
- [Bat38] H. Bateman. “Halley’s methods for solving equations”. In: *The American Mathematical Monthly* 45.1 (Jan. 1938), pp. 11–17.  
DOI: 10.2307/2303467.
- [Bes52] F. W. Bessel. “Bessel an Olbers. Königsberg, 23. April 1818”. In: *Briefwechsel zwischen W. Olbers und F. W. Bessel*. Ed. by A. Erman. Vol. 2. 260. 1852, pp. 84–90.
- [Bru70] C. Bruhns. *Neues logarithmisch-trigonometrisches Handbuch auf sieben Decimalen—A New Manual of Logarithms to Seven Places of Decimals—Nouveau manuel de logarithmes à sept décimales pour les nombres et les fonctions trigonométriques—Nuovo manuale logaritmico-trigonometrico con sette decimali*. German: <https://books.google.de/books?id=Gg1TAAAcAAJ>—English: <https://books.google.de/books?id=v-cGAAAYAAJ>—French: <https://books.google.de/books?id=RXhNAAAAAYAAJ>. Bernhard Tauchnitz, 1870. German, French, English, and Italian versions published simultaneously; we have not been able to find a facsimile of the Italian edition.
- [Can18a] S. Canon. “A trick I’ve used for years and should write up”. Tweets at <https://twitter.com/stephentyrone/status/1016283784067665920> *sqq.*, with a note at <https://twitter.com/stephentyrone/status/1016328842296864770>. 9th July 2018.
- [Can18b] S. Canon. “A trick I’ve used for years and should write up: Quick version”. Tweets at <https://twitter.com/stephentyrone/status/1057788315699687424> *sqq.* 1st Nov. 2018.
- [Can61] M. Cantor. “Note historique sur l’extraction abrégée de la racine carrée”. In: *Nouvelles annales de mathématiques. Journal des candidats aux écoles polytechnique et normale* 1.20 (1861). Ed. by O. Terquem and C.-C. Gerono, pp. 46–47.  
eprint: [http://www.numdam.org/item/?id=NAM\\_1861\\_1\\_20\\_\\_46\\_1](http://www.numdam.org/item/?id=NAM_1861_1_20__46_1).
- [Col92] P. Colwell. “Bessel Functions and Kepler’s Equation”. In: *The American Mathematical Monthly* 99.1 (Jan. 1992), pp. 45–48.  
DOI: 10.2307/2324547.



- [Dek71] T. J. Dekker. “A Floating-Point Technique for Extending the Available Precision”. In: *Numerische Mathematik* 18 (1971), pp. 224–242. DOI: 10.1007/BF01397083.
- [Fan33] T. Fantet de Lagny. “Analyse Générale, ou Méthodes nouvelles pour résoudre les Problèmes de tous les Genres & de tous les Degrez à l’infini”. In: *Recueil des Mémoires de l’Académie Royale des Sciences depuis 1666 jusqu’à 1699*. Ed. by C. Richer. Vol. XI. Par la compagnie des libraires, 1733. eprint: <https://books.google.fr/books?id=KwP-kH6gm1EC>.
- [Fan91a] T. Fantet de Lagny. “Nouvelle methode de Mr. T. F. de Lagny pour l’approximation des Racines cubiques”. In: *Le Journal des sçavans* 1691.17 (14th May 1691), pp. 200–203. eprint: <https://gallica.bnf.fr/ark:/12148/bpt6k56538h/f202.double>.
- [Fan91b] T. Fantet de Lagny. *Méthode nouvelle, infiniment générale et infiniment abrégée, Pour l’Extraction des Racines quarrées, cubiques, &c. & pour l’Approximation des mêmes Racines à l’infini dans toutes sortes d’égalitez. Proposée à examiner aux Mathématiciens de l’Europe*. De l’Imprimerie d’Antoine Lambin, ruë S. Jacques, au Miroir, 1691. eprint: <https://gallica.bnf.fr/ark:/12148/bpt6k1039787>.
- [Fan92] T. Fantet de Lagny. *Methodes nouvelles et abbregees pour l’extraction et l’approximation des racines. Et pour resoudre par le cercle et la ligne droite, plusieurs problèmes solides & sursolides ; comme la duplication du cube, l’invention de deux & de quatre moyennes proportionnelles, &c. dans toute la précision possible, & d’une maniere praticable. Avec une dissertation sur les methodes d’arithmetique & d’analyse ; où l’on établit des principes generaux pour en juger*. De l’Imprimerie de Jean Cusson, ruë saint Jacques, à l’Image de saint Jean Baptiste, 1692. eprint: <https://nubis.univ-paris1.fr/ark:/15733/3415>.
- [Fan97] T. Fantet de Lagny. *Nouveaux elemens d’arithmetique et d’algebre, ou introduction aux mathematiques*. Chez Jean Jombert, près des Augustins, à l’Image Nôtre-Dame, 1697. eprint: [https://books.google.fr/books?id=IbTtzq\\_fixAC](https://books.google.fr/books?id=IbTtzq_fixAC).
- [Fon34] B. de Fontenelle. “Éloge de M. de Lagny”. In: *Histoire de l’Académie royale des sciences avec les mémoires de mathématique et de physique tirés des registres de cette Académie 1734* (1734), pp. 107–114. eprint: <https://gallica.bnf.fr/ark:/12148/bpt6k3531x/f115>.
- [Fon58] B. de Fontelle. *Œuvres de Monsieur de Fontenelle, Des Académies, Françoises, des Sciences, des Belles-Lettres, de Londres, de Nancy, de Berlin & de Rome*. Nouvelle édition. 11 vols. B. Brunet, 1758–1766. DOI: 10.3931/e-rara-25886.
- [Gan85] W. Gander. “On Halley’s Iteration Method”. In: *The American Mathematical Monthly* 92.2 (Feb. 1985), pp. 131–134. DOI: 10.2307/2322644.
- [Hal09] E. Halley. “A new, exact, and easy Method of finding the Roots of any Equations generally, and that without any previous Reduction”. In: *The Philosophical Transactions of the Royal Society of London, from their commencement, in 1665, to the year 1800; Abridged, with notes and biographic illustrations*. Ed. by C. Hutton, G. Shaw and R. Pearson. Vol. III from 1683 to 1694. Translated from the Latin [Hal94]. 1809, pp. 640–649.
- [Hal94] E. Halley. “Methodus Nova Accurata & Facilis Inveniendi Radices Æquationum quarumcumque generaliter, sine prævia Reductione”. In: *Philosophical Transactions of the Royal Society* 18.210 (May 1694), pp. 136–148. DOI: 10.1098/rstl.1694.0029. English translation: [Hal09].
- [Higo2] N. J. Higham. *Accuracy and Stability of Numerical Algorithms*. Society for Industrial and Applied Mathematics, 2002.

- [Hou68] A. S. Householder. “Bigradients and the Problem of Routh and Hurwitz”. In: *SIAM Review* 10.1 (Jan. 1968), pp. 56–66. DOI: 10.1137/1010003.
- [Hou70] A. S. Householder. *The Numerical Treatment of a Single Nonlinear Equation*. International Series in Pure and Applied Mathematics. McGraw-Hill, 1970.
- [KB01] W. Kahan and D. Bindel. “Computing a Real Cube Root”. 2001 retypesetting by Bindel of a purported 1991 version by Kahan, at <https://cscclub.uwaterloo.ca/~pbarfuss/qbrt.pdf>. 21st Apr. 2001.
- [Kep09] J. Kepler. *Astronomia nova ατιολογητος seu physica coelestis tradita commentariis de motibus stellæ martis ex observationibus G. V. Tichonis Brahe. Jussu & sumptibus Rudolphi II romanorum imperatoris &c*, 1609. DOI: 10.3931/e-rara-558.
- [Lag67] J.-L. Lagrange. *Œuvres de Lagrange*. Ed. by J.-A. Serret and G. Darboux. 14 vols. Gauthier-Villars, 1867–1892. eprint: <https://gallica.bnf.fr/ark:/12148/bpt6k2155691>.
- [Lag71] J.-L. Lagrange. “Sur le problème de Kepler”. In: *Mémoires de l’Académie royale des sciences et belles-lettres de Berlin XXV*, année 1769 (1771), pp. 204–233.
- [Lap72] P.-S. de Laplace. “Recherches sur le calcul intégral et sur le système du monde”. In: *Histoire de l’Académie royale des sciences avec les mémoires de mathématique et de physique tirés des registres de cette Académie 1792.2* (1772), pp. 267–376. eprint: <https://gallica.bnf.fr/ark:/12148/bpt6k35711/f433>.
- [Lap78] P.-S. de Laplace. *Œuvres complètes de Laplace*. 14 vols. Gauthier-Villars, 1878–1912. eprint: <https://gallica.bnf.fr/ark:/12148/bpt6k775861>.
- [LM00] T. Lang and J.-M. Muller. *Bound on Run of Zeros and Ones for Images of Floating-Point Numbers by Algebraic Functions*. RR-4045. Institut national de recherche en informatique et en automatique, 2000. eprint: <https://hal.inria.fr/inria-00072593>.
- [Rap90] J. Raphson. *Analysis Æquationum Universalis seu Ad Æquationes Algebraicas Resolvendas Methodus Generalis, et Expedita, Ex nova Infinitarum serierum Doctrina Deducta ac Demonstrata*. Prostant venales apud Abelem Swalle, ad Insigne Monocerotis in Cœmeterio Divi Pauli, 1690. eprint: <https://books.google.co.uk/books?id=A8N1AAAAcAAJ>.
- [Sch70] E. Schröder. “Ueber unendlich viele Algorithmen zur Auflösung der Gleichungen”. In: *Mathematische Annalen* 2 (June 1870), pp. 317–365. DOI: 10.1007/BF01444024. English translation: [SS93].
- [Sch73] L. Schrön. *Tables de logarithmes à sept décimales, pour les nombres depuis 1 jusqu’à 108000 et pour les fonctions trigonométriques de dix en dix secondes*. Gauthier-Villars, 1873. eprint: <https://books.google.fr/books?id=N2a87ajKxIIC>.
- [SG01] P. Sebah and X. Gourdon. “Newton’s method and high order iterations”. In: *Numbers, constants and computation*. 3rd Oct. 2001. eprint: <http://numbers.computation.free.fr/Constants/Algorithms/newton.html>.
- [SS93] E. Schröder and G. W. Stewart. *On Infinitely Many Algorithms for Solving Equations*. Tech. rep. UMIACS-TR-92-121. Translated from the German [Sch70]. Institute for Advanced Computer Studies, University of Maryland, College Park, Jan. 1993. eprint: <http://hdl.handle.net/1903/577>.
- [ST95] T. R. Scavo and J. B. Thoo. “On the Geometry of Halley’s Method”. In: *The American Mathematical Monthly* 102.5 (May 1995), pp. 417–426. DOI: 10.2307/2975033.

- [Tre97] L. N. Trefethen. *Maxims about numerical mathematics, science, computers, and life on Earth*. Maxims. Cornell University, 1997.  
eprint: <https://people.maths.ox.ac.uk/trefethen/maxims.html>.
- [Wal85] J. Wallis. *A Treatise of Algebra, both Historical and Practical. Shewing, the Original, Progress, and Advancement thereof, from time to time; and by what Steps it hath attained to the Heighth at which now it is*. Richard Davis, 1685.  
eprint: <https://books.google.fr/books?id=TXpmAAAACAAJ>.