

xjb: Fast Float to Decimal Algorithm

Xiang JunBo¹, Wang TieJun^{2,*}

¹ Chengdu University of Information Technology; 1302035400@qq.com

² Chengdu University of Information Technology; tjw@cuit.edu.cn

* Correspondence: tjw@cuit.edu.cn; Tel.: +86-138-8090-0111

Abstract

This paper introduces an algorithm for converting floating-point numbers to decimal based on schubafch[2], which performs better than other existing algorithms.

Keywords: floating-point; decimal; performance

1. Introduction

In 1990, Steele and White[1] published the paper *how to print floating-point numbers Accurately* and proposed the optimal principle of floating-point number printing algorithms (hereinafter referred to as the SW principle) :

- **Information preservation:** The print result can be parsed back to the original floating-point number.
- **Minimum length:** The print result should be as short as possible.
- **Correct rounding:** On the basis of satisfying 1 and 2, if there are two candidate values, they should be correctly rounded (i.e., the even value should be selected).
- **Generate from left to right:** The print result is generated from the left.

Floating-point number printing algorithms that satisfy the SW principle convert floating-point numbers into real values with unique and definite results. The algorithm in this paper is based on the schubfach[2] algorithm. The algorithm in this paper has the same output result as algorithms such as ryu[3,4], schubfach[2], and dragonbox[5]. This article only introduces two floating-point number types, IEEE754-binary32 and IEEE754-binary64. To simplify the content, in this article, float represents IEEE754-binary32 and double represents IEEE754-binary64. This article involves the python code and algorithm implementation code at <https://github.com/xjb714/xjb>.

Received:

Revised:

Accepted:

Published:

Citation: JunBo, X.; TieJun, W.; xjb:

Fast Float to Decimal Algorithm.

Journal Not Specified **2025**, *1*, 0.

<https://doi.org/>

Copyright: © 2025 by the authors.

Submitted to *Journal Not Specified* for

possible open access publication under
the terms and conditions of the

Creative Commons Attribution (CC
BY) license (<https://creativecommons.org/licenses/by/4.0/>).

2. IEEE754 Floating Point Number Representation

Since the print result of a negative floating-point number only has one more negative sign than the print result of its absolute value, this article only discusses positive floating-point numbers and does not include special values such as 0, NaN, and Inf.

The IEEE754 double-precision floating-point number consists of 64 bits, including 1 sign bit (*sign*), 11 exponent bits (*exp*), and 52 fraction bits (*frac*). *sign*'s range is 0 or 1, *exp*'s range is [0, 2047], and *frac*'s range is $[0, 2^{52} - 1]$.

The IEEE754 single-precision floating-point number consists of 32 bits, including 1 sign bit (*sign*), 8 exponent bits (*exp*), and 23 fraction bits (*frac*). *sign*'s range is 0 or 1, *exp*'s range is [0, 255], and *frac*'s range is $[0, 2^{23} - 1]$.

When *frac* = 0, it is an irregular floating-point number.

The real value of the positive floating-point number v can be expressed as the following expression:

$$\begin{aligned} \text{double:}v &= \left(\text{frac} + \left(\text{exp} \neq 0?2^{52}:0 \right) \right) \cdot 2^{\max(\text{exp},1)-1075} = c \cdot 2^q \\ \text{float:}v &= \left(\text{frac} + \left(\text{exp} \neq 0?2^{23}:0 \right) \right) \cdot 2^{\max(\text{exp},1)-150} = c \cdot 2^q \end{aligned} \quad (1)$$

There are two cases in total. When exp equals 0 (referred to as subnormal floating-point numbers), there are:

$$\begin{aligned} \text{double:}v &= \text{frac} \cdot 2^{-1074} \\ \text{float:}v &= \text{frac} \cdot 2^{-149} \end{aligned} \quad (2)$$

When exp is not equal to 0 (referred to as a normal floating-point number), there is:

$$\begin{aligned} \text{double:}v &= \left(\text{frac} + 2^{52} \right) \cdot 2^{\text{exp}-1075} \\ \text{float:}v &= \left(\text{frac} + 2^{23} \right) \cdot 2^{\text{exp}-150} \end{aligned} \quad (3)$$

In the rounding interval R_v of floating-point numbers, all real numbers will be rounded to this floating-point number when parsed. R_v is:

$$\begin{aligned} v_l &= \begin{cases} \left(c - \frac{1}{2} \right) \cdot 2^q, & \text{if } \text{frac} \neq 0 \text{ or } \text{exp} \leqslant 1 \\ \left(c - \frac{1}{4} \right) \cdot 2^q, & \text{if } \text{frac} = 0 \end{cases} \\ v_r &= \left(c + \frac{1}{2} \right) \cdot 2^q \\ R_v &= \begin{cases} [v_l, v_r], & \text{if } \text{frac}\%2 = 0 \\ (v_l, v_r), & \text{if } \text{frac}\%2 = 1 \end{cases} \end{aligned} \quad (4)$$

When the floating-point number is a regular floating-point number, 2^{q-1} is the rounded radius.

3. Principle of Algorithm

At present, other algorithms use a large number of branches, which can easily lead to branch prediction failure penalties and excessive high multiplication overhead. The algorithm in this paper will minimize the overhead of branch prediction failures and reduce the number of multiplication operations to improve performance. Moreover, the core code for the algorithm implementation in this paper is only about twenty lines and it also supports parallel computing. The process of printing floating-point numbers is usually divided into two parts: the first part is to convert the floating-point number to a decimal number, and the second part is to convert the decimal number to a string. And this article will only introduce the first part. All double-precision floating-point numbers are classified into two types: irregular values and regular values. An irregular value is one where all the lower 52 bits are 0, meaning the frac value is 0. There are a total of 2046 valid irregular values (i.e., exp values range from 1 to 2046). Dividing by the irregular values yields the regular value. Similarly, there are a total of 254 irregular values in a single-precision floating-point number. When exp is 0, it is called a subnormal floating-point number.

The valid range for c and q in regular floating-point numbers is:

$$\begin{aligned} \text{float} : & \begin{cases} 1 \leq c \leq 2^{24} - 1, c \neq 2^{23}; q = -149 \\ 2^{23} + 1 \leq c \leq 2^{24} - 1; -148 \leq q \leq 104 \end{cases} \\ \text{double} : & \begin{cases} 1 \leq c \leq 2^{53} - 1, c \neq 2^{52}; q = -1074 \\ 2^{52} + 1 \leq c \leq 2^{53} - 1; -1073 \leq q \leq 971 \end{cases} \end{aligned} \quad (5)$$

The valid range for c and q in irregular floating-point numbers is:

$$\begin{aligned} \text{float} : & \begin{cases} c = 2^{23}; -149 \leq q \leq 104 \end{cases} \\ \text{double} : & \begin{cases} c = 2^{52}; -1074 \leq q \leq 971 \end{cases} \end{aligned} \quad (6)$$

The valid range for c and q in subnormal floating-point numbers is:

$$\begin{aligned} \text{float} : & \begin{cases} c \leq 2^{23} - 1; q = -149 \end{cases} \\ \text{double} : & \begin{cases} c \leq 2^{52} - 1; q = -1074 \end{cases} \end{aligned} \quad (7)$$

Floating-point numbers that do not fall within the subnormal range are called normal floating-point numbers.

regular floating-point numbers account for the vast majority of all possible values of floating-point numbers and are the most worthy of discussion part. Therefore, unless otherwise specified, only regular floating-point numbers will be discussed below. Suppose the floating-point number v is converted to the optimal solution that satisfies the SW principle as opt , d is a positive integer and k is an integer, which is expressed as:

$$\begin{aligned} v = c \cdot 2^q \rightarrow opt = d \cdot 10^k \\ opt \in R_v; d \in N^+; k \in Z \end{aligned} \quad (8)$$

For example: IEEE754-binary64 floating-point number "1.3", the real value of the floating-point number is 1.300000000000000444089209850062616169452667236328125, hexadeciml representation of floating-point Numbers is 3ff4ccccccccc, Then the opt value that meets the SW principle is 1.3. The IEEE754-binary32 floating-point number "1.3" has an actual value of 1.2999999523162841796875, and its hexadeciml representation is 3FA66666. Therefore, the opt value that satisfies the SW principle is 1.3.

3.1. Review the schubfach algorithm and the derivation of the algorithm in this paper

According to the schubfach algorithm, the possible values of d can be one of the following four situations:

$$10 \cdot \lfloor v \cdot 10^{-k-1} \rfloor, \lfloor 10 \cdot (v \cdot 10^{-k-1}) \rfloor, \lfloor 10 \cdot (v \cdot 10^{-k-1}) \rfloor + 1, 10 \cdot \lfloor v \cdot 10^{-k-1} \rfloor + 10 \quad (9)$$

The calculation method of k in equation (9) is as follows:

$$k = \lfloor q \cdot \lg(2) \rfloor \text{ if } v \in \text{regular} \text{ else } \lfloor q \cdot \lg(2) - \lg\left(\frac{4}{3}\right) \rfloor \quad (10)$$

In the range of float and double, equation (10) can be equivalent to:

$$k = (q \cdot 315653 - (v \in \text{regular}?0:131237)) \gg 20 \quad (11)$$

Suppose the integer part of $v \cdot 10^{-k-1}$ is m and the decimal part is n , then we have:

$$\begin{aligned} \lfloor v \cdot 10^{-k-1} \rfloor &= m \\ v \cdot 10^{-k-1} &= m + n \\ 0 \leq n = v \cdot 10^{-k-1} - \lfloor v \cdot 10^{-k-1} \rfloor &< 1 \end{aligned} \tag{12}$$

Then the decimal part of $v \cdot 10^{-k}$ is expressed as:

$$v \cdot 10^{-k} - \lfloor v \cdot 10^{-k} \rfloor = 10m + 10n - \lfloor 10m + 10n \rfloor = 10n - \lfloor 10n \rfloor \tag{13}$$

The possible values of d obtained from equation (9) are:

$$10m, \lfloor 10(m+n) \rfloor, \lfloor 10(m+n) \rfloor + 1, 10m + 10 \tag{14}$$

The possible values of d in equation (14) can be simplified to:

$$10m, 10m + \lfloor 10n \rfloor, 10m + \lfloor 10n \rfloor + 1, 10m + 10 \tag{15}$$

Among them, $10m$ represents the minimum possible value and $10m + 10$ represents the maximum possible value. Suppose ten is used to represent $10m$. There are four possible values for one , with $d = ten + one$, denoted as:

$$\begin{aligned} ten &= 10m \\ one &\in \{0, \lfloor 10n \rfloor, \lfloor 10n \rfloor + 1, 10\} \\ d &= ten + one \end{aligned} \tag{16}$$

Calculating d will be converted to calculating ten and one .

The final possible values of d are as follows:

- $10m$

When the following conditions are met, the result is $10m$ (or equivalent to $one = 0$).

That is, the floating-point number v minus the minimum possible value of $10m$ is less than the rounded radius 2^{q-1} .

$$\begin{aligned} c \cdot 2^q - 10m \cdot 10^k &< 2^{q-1} \\ c \cdot 2^q - \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor \cdot 10^{k+1} &< 2^{q-1} \\ c \cdot 2^q \cdot 10^{-k-1} - \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor &< 2^{-1} \cdot 2^q \cdot 10^{-k-1} \\ n &< 2^{-1} \cdot 2^q \cdot 10^{-k-1} \\ 2^{-1} \cdot 2^q \cdot 10^{-k-1} &> n \end{aligned} \tag{17}$$

Or when $2^{-1} \cdot 2^q \cdot 10^{-k-1} = n$, $c \% 2 = 0$ must also be satisfied. Therefore, the following conditions are valid:

$$\text{if } 2^{-1} \cdot 2^q \cdot 10^{-k-1} > n \text{ or } (2^{-1} \cdot 2^q \cdot 10^{-k-1} = n \ \&\& \ c \% 2 = 0) : one = 0 \tag{18}$$

- $10m + 10$

When the following conditions are met, the result is $10m + 10$ (or equivalent to $one = 10$). The maximum possible value of $10m + 10$ minus the floating-point number v is less than the rounded radius 2^{q-1} .

$$\begin{aligned}
 (10m + 10) \cdot 10^k - c \cdot 2^q &< 2^{q-1} \\
 \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor \cdot 10^{k+1} + 10^{k+1} - c \cdot 2^q &< 2^{q-1} \\
 \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor - c \cdot 2^q \cdot 10^{-k-1} + 1 &< 2^{-1} \cdot 2^q \cdot 10^{-k-1} \\
 1 - n &< 2^{-1} \cdot 2^q \cdot 10^{-k-1} \\
 2^{-1} \cdot 2^q \cdot 10^{-k-1} &> 1 - n
 \end{aligned} \tag{19}$$

Or when $2^{-1} \cdot 2^q \cdot 10^{-k-1} = 1 - n$, $c \% 2 = 0$ must also be satisfied. Therefore, the following conditions are valid:

$$\text{if } 2^{-1} \cdot 2^q \cdot 10^{-k-1} > 1 - n \text{ or } \left(2^{-1} \cdot 2^q \cdot 10^{-k-1} = 1 - n \text{ \&\& } c \% 2 = 0 \right) : one = 10 \tag{20}$$

- $10m + \lfloor 10n \rfloor$ **or** $10m + \lfloor 10n \rfloor + 1$

When none of the conditions are met as $d = 10m$ or $d = 10m + 10$, d is either $10m + \lfloor 10n \rfloor$ or $10m + \lfloor 10n \rfloor + 1$. The final value is determined based on the decimal part of $10n$. If the decimal part is 0.5, it is rounded to the nearest even value; if it is not 0.5, it is rounded to the nearest value. For irregular floating-point numbers, it is also necessary to determine whether $10m + \lfloor 10n \rfloor$ is within the rounding interval R_v . If it is not, then $10m + \lfloor 10n \rfloor + 1$.

In summary, the steps of the schubfach algorithm variants are as follows, that is, the algorithms proposed in this paper (xjb32(for float), xjb64(for double)):

```

input : c, q
output : d, k
convert  $c \cdot 2^q$  to  $d \cdot 10^k$ 
( 1 )  $v = c \cdot 2^q$ 
( 2 )  $k = \lfloor q \cdot \lg(2) \rfloor$  if  $v \in regular$  else  $\lfloor q \cdot \lg(2) - \lg(\frac{4}{3}) \rfloor$ 
( 3 )  $m = \lfloor v \cdot 10^{-k-1} \rfloor$ ,  $n = v \cdot 10^{-k-1} - m$ 
( 4 )  $ten = 10m$ 
( 5 ) if  $10n - \lfloor 10n \rfloor = 0.5$  :  $one = \lfloor 10n \rfloor$  if  $(\lfloor 10n \rfloor \% 2 = 0)$  else  $\lfloor 10n \rfloor + 1$ 
( 6 ) if  $10n - \lfloor 10n \rfloor < 0.5$  :  $one = \lfloor 10n \rfloor$ 
( 7 ) if  $10n - \lfloor 10n \rfloor > 0.5$  :  $one = \lfloor 10n \rfloor + 1$ 
( 8 ) if  $v \in irregular$  :
( 9 ) if  $10n - \lfloor 10n \rfloor > 2^{q-2} \cdot 10^{-k}$  :  $one = \lfloor 10n \rfloor + 1$ 
( 10 ) if  $2^{q-2} \cdot 10^{-k-1} \geq n$  :  $one = 0$ 
( 11 ) else :
( 12 ) if  $2^{q-1} \cdot 10^{-k-1} > n$  or  $\left( 2^{q-1} \cdot 10^{-k-1} = n \text{ \&\& } c \% 2 = 0 \right)$  :  $one = 0$ 
( 13 ) endif
( 14 ) if  $2^{q-1} \cdot 10^{-k-1} > 1 - n$  or  $\left( 2^{q-1} \cdot 10^{-k-1} = 1 - n \text{ \&\& } c \% 2 = 0 \right)$  :  $one = 10$ 
( 15 )  $d = ten + one$ 

```

This algorithm process (21) is applicable to float and double floating-point numbers. Taking a floating-point number v as input, c and q are extracted, and the calculation results d (line 15) and k (line 2) are returned. The real value represented by the returned results is $d \cdot 10^k$, which conforms to the SW principle. The calculation process of k is relatively simple and can be obtained from (11). Therefore, the following only focuses on introducing the rapid calculation process of d .

The following will be divided into five parts to introduce the algorithm process (21) :

1. Introduce the pre-computation process of the algorithm's lookup table.
2. Quickly calculate m .
3. Quickly determine whether $one = 0$ or $one = 10$.
4. Quickly calculate $\lfloor 10n \rfloor$ and determine whether $one = \lfloor 10n \rfloor$ or $one = \lfloor 10n \rfloor + 1$ based on the decimal part of $10n$.
5. Processing of irregular floating-point numbers.

3.2. Pre-computation of Lookup Table

The algorithm in this paper uses a lookup table to store the values of 10^{-k-1} for q in the range of $[-149, 104]$ for float and $[-1074, 971]$ for double. In the algorithm of this paper, float uses 64-bit precision and double uses 128-bit precision lookup tables. The code implementation in this section is gen.py. Suppose the bit length of a single value data in the lookup table is B . For float, it has $B = 64$, and for double, it has $B = 128$. Suppose there are integers e_{10} and real numbers e_2 , where $1 \leq f < 2$. There are:

$$f \cdot 2^{\lfloor e_2 \rfloor} = 2^{e_2} = 10^{e_{10}} \quad (22)$$

Then:

$$\lfloor e_2 \rfloor = \lfloor e_{10} \cdot \lg(2) \rfloor \quad (23)$$

The calculation leads to f , and the following conclusions are drawn:

$$f = \frac{10^{e_{10}}}{2^{\lfloor e_{10} \cdot \lg(2) \rfloor}} \quad (24)$$

The way to calculate the lookup table is as follows (using the upward rounding method) :

$$lookup[e_{10}] = \lceil f \cdot 2^{B-1} \rceil = \lceil \frac{10^{e_{10}}}{2^{\lfloor e_{10} \cdot \lg(2) \rfloor}} \cdot 2^{B-1} \rceil = \lceil 10^{e_{10}} \cdot 2^{B-1-\lfloor e_{10} \cdot \lg(2) \rfloor} \rceil \quad (25)$$

For float, when $0 \leq e_{10} \leq 27$, $f \cdot 2^{B-1}$ is an integer in equation (25). For double, when $0 \leq e_{10} \leq 55$, $f \cdot 2^{B-1}$ is an integer in equation (25). The detailed calculation process is as follows:

- Float

The range of $-k - 1$ is calculated to be $[-32, 44]$ through the q value range in equation (5), so the lookup table contains representation values from 10 to the power of -32 to 10 to the power of 44. The calculation process is as follows:

$$-32 \leq e_{10} \leq 44$$

$$e_2 = \lfloor e_{10} \cdot \log_2(10) \rfloor - 63$$

$$pow10t = \begin{cases} 2^{e_2} / 10^{\lfloor e_{10} \rfloor}; & \text{if } e_{10} < 0 \\ 10^{\lfloor e_{10} \rfloor} / 2^{e_2}; & \text{if } e_{10} \geq 20 \\ 10^{\lfloor e_{10} \rfloor} \cdot 2^{e_2}; & \text{if } 1 \leq e_{10} \leq 19 \end{cases} \quad (26)$$

$$f_{1,e_{10}} = pow10 = pow10t + (e_{10} \geq 0 \& e_{10} \leq 27? 0 : 1)$$

When $0 \leq e_{10} \leq 27$, the lookup table variable indicates that the values $f_{1,e_{10}} \cdot 2^{\lfloor e_{10} \cdot \log_2(10) \rfloor - 63}$ and $10^{e_{10}}$ are equal. In other cases, the relative error is less than 2^{-63} . Expressed as:

$$r_{1,e_{10}} = \frac{f_{1,e_{10}} \cdot 2^{\lfloor e_{10} \cdot \log_2(10) \rfloor - 63}}{10^{e_{10}}} \in \begin{cases} 1; & \text{if } 0 \leq e_{10} \leq 27 \\ (1, 1 + 2^{-63}); & \text{if } e_{10} < 0 \text{ or } e_{10} > 27 \end{cases} \quad (27)$$

- Double

The range of $-k - 1$ is calculated to be $[-293, 323]$ through the q value range in equation (5), so the lookup table contains representation values from 10 to the power of -293 to 10 to the power of 323. The calculation process is as follows:

$$\begin{aligned} -293 &\leq e_{10} \leq 323 \\ e_2 &= |\lfloor e_{10} \cdot \log_2(10) \rfloor - 127| \\ pow10t &= \begin{cases} 2^{e_2} / 10^{\lfloor e_{10} \rfloor}; & \text{if } e_{10} < 0 \\ 10^{\lfloor e_{10} \rfloor} / 2^{e_2}; & \text{if } e_{10} \geq 39 \\ 10^{\lfloor e_{10} \rfloor} \cdot 2^{e_2}; & \text{if } 1 \leq e_{10} \leq 38 \end{cases} \\ f_{1,e_{10}} &= pow10 = pow10t + (e_{10} \geq 0 \& e_{10} \leq 55? 0 : 1) \end{aligned} \quad (28)$$

When $0 \leq e_{10} \leq 55$, the lookup table variable indicates that the values $f_{1,e_{10}} \cdot 2^{\lfloor e_{10} \cdot \log_2(10) \rfloor - 127}$ and $10^{e_{10}}$ are equal. In other cases, the relative error is less than 2^{-127} . Expressed as:

$$r_{1,e_{10}} = \frac{f_{1,e_{10}} \cdot 2^{\lfloor e_{10} \cdot \log_2(10) \rfloor - 127}}{10^{e_{10}}} \in \begin{cases} 1; & \text{if } 0 \leq e_{10} \leq 55 \\ (1, 1 + 2^{-127}); & \text{if } e_{10} < 0 \text{ or } e_{10} > 55 \end{cases} \quad (29)$$

The following uses r_1 to represent all possible errors of the lookup table values within the float range, r_2 to represent all possible errors of the lookup table values within the double range, and r to represent all possible errors of the lookup table values within either the float or double range. In algorithm process (21), an approximate representation value of 10 to the power of $-k - 1$ needs to be obtained through a lookup table. From equation (27) and equation (29), the lookup table representation value is error-free when q is within the following range:

$$\begin{aligned} float: 0 &\leq -k - 1 \leq 27 \Rightarrow -93 \leq q \leq -1 \\ double: 0 &\leq -k - 1 \leq 55 \Rightarrow -186 \leq q \leq -1 \end{aligned} \quad (30)$$

When q is not within the range of equation (30), the error range of the value represented by the lookup table can be concluded as follows:

$$\begin{aligned} float: 0 &< r_1 - 1 < 2^{-63} \\ double: 0 &< r_2 - 1 < 2^{-127} \end{aligned} \quad (31)$$

The introduction of the lookup table calculation process is complete. The storage space required for a float range lookup table is 616 bytes, and that for a double range lookup table is 9872 bytes.

3.3. Quickly Calculate m

Relevant theorems (partially from the dragonbox[5] algorithm paper) : Suppose there are positive integers n, P , and Q , where P and Q are coprime, $P < Q$, $1 \leq n \leq n_{max}$, $Q > n_{max}$, P^*/Q^* is the best rational approximation result greater than or equal to P/Q , P_*/Q_* is the best rational approximation result less than or equal to P/Q , and it satisfies $Q^* \leq n_{max}$, $Q_* \leq n_{max}$. And if $n \cdot P$ does not divide Q evenly, it is expressed as:

$$\lfloor n \cdot \frac{P}{Q} \rfloor + 1 = \lceil n \cdot \frac{P}{Q} \rceil \quad (32)$$

Suppose the following holds true:

$$\lfloor n \cdot \frac{P}{Q} \rfloor = \lfloor n \cdot \xi \rfloor \quad (33)$$

Then there are:

$$\frac{P_*}{Q_*} = \max_{1 \leq n \leq n_{max}} \frac{\lfloor n \cdot \frac{P}{Q} \rfloor}{n} \leq \xi < \min_{1 \leq n \leq n_{max}} \frac{\lfloor n \cdot \frac{P}{Q} \rfloor + 1}{n} = \min_{1 \leq n \leq n_{max}} \frac{\lceil n \cdot \frac{P}{Q} \rceil}{n} = \frac{P^*}{Q^*} \quad (34)$$

Therefore, the range of values for ξ is:

$$\frac{P_*}{Q_*} \leq \xi < \frac{P^*}{Q^*} \quad (35)$$

And the range of the decimal part with $n \cdot \frac{P}{Q}$ is:

$$\left[\frac{(Q_* P) \% Q}{Q}, \frac{(Q^* P) \% Q}{Q} \right] \quad (36)$$

That is, when $n = Q_*$, the decimal part is the smallest; when $n = Q^*$, the decimal part is the largest.

The definition of the best rational approximation function is as follows (this function is implemented on line 15 of the test1.py file):

$$(DN, UP) = f(C, P, Q) \quad (37)$$

The function (37) Calculate the best rational approximation result with a denominator not exceeding C based on the mean term theorem of the Farey sequence. DN and UP are two adjacent terms in the C -order Farey sequence F_C .

In algorithm process (21), m is calculated as $\lfloor v \cdot 10^{-k-1} \rfloor$ (line 3). Just prove that the following equation holds:

$$m = \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor c \cdot 2^q \cdot r \cdot 10^{-k-1} \rfloor \quad (38)$$

Where r is the error of the lookup table value, as defined in equation (27) and equation (29). When the condition (30) is met, r is 1, and the equation (38) clearly holds. When r is not 1, there is:

$$\begin{aligned} float &: 1 < r < 1 + 2^{-63} \\ double &: 1 < r < 1 + 2^{-127} \end{aligned} \quad (39)$$

Calculate the range of $2^q \cdot 10^{-k-1}$ and we get:

$$2^q \cdot 10^{-k-1} = 10^{-1} \cdot \left(10^{q \cdot \lg(2) - \lfloor q \cdot \lg(2) \rfloor} \right) \quad (40)$$

When q is not 0, equation (40) exists:

$$\begin{aligned} q \cdot \lg(2) &\neq \lfloor q \cdot \lg(2) \rfloor \\ 0 < q \cdot \lg(2) - \lfloor q \cdot \lg(2) \rfloor &< 1 \end{aligned} \quad (41)$$

When q is 0, $q \cdot \lg(2) - \lfloor q \cdot \lg(2) \rfloor = 0$, so the final conclusion is:

$$10^{-1} \leq 2^q \cdot 10^{-k-1} < 1 \quad (42)$$

Because there is:

$$c \cdot 2^q \cdot 10^{-k-1} = c \cdot \frac{2^{q-k-1}}{5^{k+1}} \in [0.1c, c) \quad (43)$$

Therefore:

$$c \cdot 2^q \cdot 10^{-k-1} = \begin{cases} \frac{c \cdot 2^{q-k-1}}{5^{k+1}}; q \geq 1 \\ \frac{c}{2^{1+k-q} \cdot 5^{k+1}} = \frac{c}{10}; q = 0 \\ \frac{c \cdot 5^{-k-1}}{2^{1+k-q}}; q < 0 \end{cases} \quad (44)$$

Suppose:

$$c \cdot 2^q \cdot 10^{-k-1} = c \cdot \frac{x}{y} < c \quad (45)$$

Then there are:

$$(x, y) = \begin{cases} (2^{q-k-1}, 5^{k+1}); q \geq 1 \\ (1, 10); q = 0 \\ (5^{-k-1}, 2^{1+k-q}); q < 0 \end{cases} \quad (46)$$

Suppose:

$$\begin{aligned} float : c &\leq c_{\max} = C_1 = 2^{24} - 1 \\ double : c &\leq c_{\max} = C_2 = 2^{53} - 1 \end{aligned} \quad (47)$$

The following is represented by C as C_1 or C_2 . C within the float range is C_1 , and C within the double range is C_2 .

When $y > C$, calculate the P^* and Q^* corresponding to each q by calling $f(C, x, y)$ according to function (37). And calculate the minimum BIT value when the following conditions are met:

$$\frac{x}{y} \left(1 + 2^{-BIT} \right) < \frac{P^*}{Q^*} \quad (48)$$

When $y \leq C$, there is:

$$c \cdot \frac{x}{y} \left(1 + \frac{1}{Cy} \right) = \frac{cx + \frac{c}{C} \cdot \frac{x}{y}}{y} < \frac{cx + 1}{y} \quad (49)$$

Therefore:

$$\lfloor c \cdot \frac{x}{y} \rfloor = \lfloor c \cdot \frac{x}{y} \left(1 + \frac{1}{Cy} \right) \rfloor \quad (50)$$

Similarly, calculate the minimum BIT value:

$$\frac{x}{y} \left(1 + 2^{-BIT} \right) < \frac{x}{y} \left(1 + \frac{1}{Cy} \right) \quad (51)$$

In summary, the calculation results of the maximum value among the minimum *BIT* values corresponding to different q are as follows (the running result is in the test1.py file, and the running time of this code is only about 1 to 2 seconds) :

$$\begin{aligned} \text{float : } & \text{BIT}_{\max} = 52 \\ \text{double : } & \text{BIT}_{\max} = 113 \end{aligned} \quad (52)$$

Therefore, the following conclusions exist:

$$\begin{aligned} \text{float : } & \lfloor c \cdot \frac{x}{y} \rfloor = \lfloor c \cdot \frac{x}{y} \cdot (1 + 2^{-52}) \rfloor = \lfloor c \cdot \frac{x}{y} \cdot r_1 \rfloor \\ \text{double : } & \lfloor c \cdot \frac{x}{y} \rfloor = \lfloor c \cdot \frac{x}{y} \cdot (1 + 2^{-113}) \rfloor = \lfloor c \cdot \frac{x}{y} \cdot r_2 \rfloor \end{aligned} \quad (53)$$

This section has been verified. After quickly calculating m , the value of $ten = 10m$ can be obtained very quickly.

3.4. Quickly Determine Whether $one = 0$ or $one = 10$

In algorithm process (21), the conditions for determining $one = 0$ and $one = 10$ are on lines 12, and 14. This section will introduce how to quickly determine whether $one = 0$ or $one = 10$ holds by using equivalent conditions.

When discussing the case of $2^{-1} \cdot 2^q \cdot 10^{-k-1} = n$ (line 12, one might be 0), it is equivalent to:

$$\begin{aligned} c \cdot 2^q \cdot 10^{-k-1} - \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor &= 2^{-1} \cdot 2^q \cdot 10^{-k-1} \\ (2c - 1) \cdot 2^{q-1} \cdot 10^{-k-1} &= \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor \end{aligned} \quad (54)$$

When discussing the case of $2^{-1} \cdot 2^q \cdot 10^{-k-1} = 1 - n$ (line 14, one might be 10), it is equivalent to:

$$\begin{aligned} \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor - c \cdot 2^q \cdot 10^{-k-1} + 1 &= 2^{-1} \cdot 2^q \cdot 10^{-k-1} \\ (2c + 1) \cdot 2^{q-1} \cdot 10^{-k-1} &= \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor + 1 \end{aligned} \quad (55)$$

Since equation (42), we have:

$$2^{q-1} \cdot 10^{-k-1} \in [0.05, 0.5] \quad (56)$$

Therefore, there is:

$$\begin{aligned} \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor - 1 &< c \cdot 2^q \cdot 10^{-k-1} - 0.5 \\ &< (2c - 1) \cdot 2^{q-1} \cdot 10^{-k-1} \\ &\leqslant c \cdot 2^q \cdot 10^{-k-1} - 0.05 < \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor + 1 \end{aligned} \quad (57)$$

Therefore, for equation (54), when $(2c - 1) \cdot 2^{q-1} \cdot 10^{-k-1}$ is an integer, it must be equal to $\lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor$. Similarly, for equation (55), there is:

$$\begin{aligned} \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor &< c \cdot 2^q \cdot 10^{-k-1} + 0.05 \\ &\leqslant (2c + 1) \cdot 2^{q-1} \cdot 10^{-k-1} \\ &< c \cdot 2^q \cdot 10^{-k-1} + 0.5 < \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor + 2 \end{aligned} \quad (58)$$

Therefore, for equation (55), when $(2c + 1) \cdot 2^{q-1} \cdot 10^{-k-1}$ is an integer, it must be equal to $\lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor + 1$.

In conclusion, it is equivalent to discussing whether $(2c \pm 1) \cdot 2^{q-1} \cdot 10^{-k-1}$ is an integer. There are:

$$(2c \pm 1) \cdot 2^{q-1} \cdot 10^{-k-1} = (2c \pm 1) \cdot 2^{q-k-2} \cdot 5^{-k-1} \quad (59)$$

According to the range of q , there are:

$$\begin{cases} q - k - 2 \geq 0, -k - 1 < 0; q \geq 2 \\ q - k - 2 < 0, -k - 1 < 0; 1 \geq q \geq 0 \\ q - k - 2 < 0, -k - 1 \geq 0; q < 0 \end{cases} \quad (60)$$

Therefore, equation (59) is equivalent to:

$$(2c \pm 1) \cdot 2^{q-1} \cdot 10^{-k-1} = \begin{cases} \frac{(2c \pm 1) \cdot 2^{q-k-2}}{5^{k+1}}; q \geq 2 \\ \frac{(2c \pm 1)}{2^{2+k-q} \cdot 5^{k+1}}; 1 \geq q \geq 0 \\ \frac{(2c \pm 1) \cdot 5^{-k-1}}{2^{2+k-q}}; q < 0 \end{cases} \quad (61)$$

According to the different ranges of q , the following situations are discussed:

- $q \geq 2$

From $q \geq 2$, we get $k \geq 0$. When $q \geq 2$, it is equivalent to discussing whether $(2c \pm 1) \cdot 2^{q-k-2}$ is divisible by 5^{k+1} . Since 2 and 5 are coprime, it is equivalent to discussing whether $(2c \pm 1)$ is divisible by 5^{k+1} .

$$(2c \pm 1) \% 5^{k+1} = 0 \quad (62)$$

Suppose t is a positive integer:

$$2c \pm 1 = t \cdot 5^{k+1}; t \geq 1 \quad (63)$$

Since $2c \pm 1$ is odd, t is also odd. Because the following conditions exist:

$$\begin{aligned} float: 2c - 1 &\in [2^{24} + 1, 2^{25} - 3]; 2c + 1 \in [2^{24} + 3, 2^{25} - 1]; \\ double: 2c - 1 &\in [2^{53} + 1, 2^{54} - 3]; 2c + 1 \in [2^{53} + 3, 2^{54} - 1]; \end{aligned} \quad (64)$$

Therefore, the following satisfies:

$$\begin{aligned} float: 2^{24} + 1 &\leq t \cdot 5^{k+1} \leq 2^{25} - 1 \\ double: 2^{53} + 1 &\leq t \cdot 5^{k+1} \leq 2^{54} - 1 \end{aligned} \quad (65)$$

Therefore, the following conclusions are drawn:

$$\begin{aligned} float: \frac{2^{24} + 1}{5^{k+1}} &\leq t \leq \frac{2^{25} - 1}{5^{k+1}}; \\ double: \frac{2^{53} + 1}{5^{k+1}} &\leq t \leq \frac{2^{54} - 1}{5^{k+1}}; \end{aligned} \quad (66)$$

For the above equation (66), the maximum value of k when t can obtain at least one odd number is:

$$\begin{aligned} float: k_{\max} &= 9 \Rightarrow q_{\max} = 33, t = 3 \\ double: k_{\max} &= 22 \Rightarrow q_{\max} = 76, t = 1 \end{aligned} \quad (67)$$

Therefore, the maximum value of k is 9 within the float range and 22 within the double range. Therefore, when k exceeds the above range, $(2c \pm 1)$ is not divisible by 5^{k+1} .

- $1 \geq q \geq 0$

Because the denominator $2^{2+k-q} \cdot 5^{k+1}$ is even and the numerator $(2c \pm 1)$ is odd, the condition is not met.

- $q < 0$

Because the denominator 2^{2+k-q} is even and the numerator $(2c \pm 1) \cdot 5^{-k-1}$ is odd, the condition is not met.

In summary, the situations where $(2c \pm 1) \cdot 2^{q-1} \cdot 10^{-k-1}$ may be an integer are as follows:

$$\begin{aligned} \text{float : } & 2 \leq q \leq 33 \text{ && } (2c \pm 1) \% 5^{k+1} = 0; \\ \text{double : } & 2 \leq q \leq 76 \text{ && } (2c \pm 1) \% 5^{k+1} = 0; \end{aligned} \quad (68)$$

And, the range of $-k - 1$ is:

$$\begin{aligned} \text{float : } & -10 \leq -k - 1 \leq -1 \\ \text{double : } & -23 \leq -k - 1 \leq -1 \end{aligned} \quad (69)$$

When $2^{-1} \cdot 2^q \cdot 10^{-k-1} = n$, the following conclusions can be drawn:

$$\begin{aligned} \text{float : } & \left\{ 2^{35} \cdot 2^q \cdot 10^{-k-1} = n \cdot 2^{36} \Rightarrow \lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor n \cdot 2^{36} \rfloor \right. \\ \text{double : } & \left\{ 2^{63} \cdot 2^q \cdot 10^{-k-1} = n \cdot 2^{64} \Rightarrow \lfloor 2^{63} \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor n \cdot 2^{64} \rfloor \right. \end{aligned} \quad (70)$$

When $2^{-1} \cdot 2^q \cdot 10^{-k-1} = 1 - n$, the following conclusions can be drawn:

$$\begin{aligned} \text{float : } & \left\{ \begin{array}{l} 2^{35} \cdot 2^q \cdot 10^{-k-1} = 2^{36} - n \cdot 2^{36} \Rightarrow \\ \lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor 2^{36} - n \cdot 2^{36} \rfloor = 2^{36} - 1 - \lfloor n \cdot 2^{36} \rfloor \end{array} \right. \\ \text{double : } & \left\{ \begin{array}{l} 2^{63} \cdot 2^q \cdot 10^{-k-1} = 2^{64} - n \cdot 2^{64} \Rightarrow \\ \lfloor 2^{63} \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor 2^{64} - n \cdot 2^{64} \rfloor = 2^{64} - 1 - \lfloor n \cdot 2^{64} \rfloor \end{array} \right. \end{aligned} \quad (71)$$

The discussion on whether $\lfloor 2^{36} - n \cdot 2^{36} \rfloor = 2^{36} - 1 - \lfloor n \cdot 2^{36} \rfloor$ in equation (71) holds true, that is, whether $2^{36} \cdot n$ in equation (71) is an integer, or equivalent to discussing whether the following values are integers when equation (68) holds true (the same applies to double) :

$$\begin{aligned} \text{float : } & c \cdot 2^{q+36} \cdot 10^{-k-1} = c \cdot 2^{q-k+35} \cdot 5^{-k-1} = c \cdot \frac{2^{q-k+35}}{5^{k+1}} \\ \text{double : } & c \cdot 2^{q+64} \cdot 10^{-k-1} = c \cdot 2^{q-k+63} \cdot 5^{-k-1} = c \cdot \frac{2^{q-k+63}}{5^{k+1}} \end{aligned} \quad (72)$$

Suppose c can divide 5^{k+1} evenly (where t is a temporary integer variable):

$$c = t \cdot 5^{k+1}; t \geq 1 \quad (73)$$

Therefore, when equation (73) was established, there were:

$$2c \pm 1 = 2 \cdot t \cdot 5^{k+1} \pm 1 \quad (74)$$

Expression (74) cannot divide 5^{k+1} evenly, which contradicts equation (68), so c cannot divide 5^{k+1} evenly. Therefore, for float, $c \cdot 2^{q+36} \cdot 10^{-k-1}$ is not an integer; For double, $c \cdot 2^{q+64} \cdot 10^{-k-1}$ is not an integer, that is:

$$\begin{aligned} \text{float : } & \lfloor 2^{36} - 2^{36} \cdot n \rfloor = 2^{36} + \lfloor -2^{36} \cdot n \rfloor = 2^{36} - 1 - \lfloor 2^{36} \cdot n \rfloor \\ \text{double : } & \lfloor 2^{64} - 2^{64} \cdot n \rfloor = 2^{64} + \lfloor -2^{64} \cdot n \rfloor = 2^{64} - 1 - \lfloor 2^{64} \cdot n \rfloor \end{aligned} \quad (75)$$

Therefore, the conclusion (71) is correct. Discuss the necessary and sufficient conditions for whether $\lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor n \cdot 2^{36} \rfloor$ is $2^{-1} \cdot 2^q \cdot 10^{-k-1} = n$. The same applies to double, expressed as:

$$\begin{aligned} float : 2^{-1} \cdot 2^q \cdot 10^{-k-1} = n &\Leftrightarrow \lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor n \cdot 2^{36} \rfloor \\ double : 2^{-1} \cdot 2^q \cdot 10^{-k-1} = n &\Leftrightarrow \lfloor 2^{63} \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor n \cdot 2^{64} \rfloor \end{aligned} \quad (76)$$

Similarly, the necessary and sufficient conditions for whether $\lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor 2^{36} - n \cdot 2^{36} \rfloor$ is $2^{-1} \cdot 2^q \cdot 10^{-k-1} = 1 - n$. The same applies to double, expressed as:

$$\begin{aligned} float : 2^{-1} \cdot 2^q \cdot 10^{-k-1} = 1 - n &\Leftrightarrow \lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor 2^{36} - n \cdot 2^{36} \rfloor \\ double : 2^{-1} \cdot 2^q \cdot 10^{-k-1} = 1 - n &\Leftrightarrow \lfloor 2^{63} \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor 2^{64} - n \cdot 2^{64} \rfloor \end{aligned} \quad (77)$$

The sufficient conditions of equations (76) and (77) are obviously established. Introduce the proof that equation (76) holds. For float, only the necessary conditions need to be discussed, that is, whether $2^{-1} \cdot 2^q \cdot 10^{-k-1} = n$ must hold true when $\lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor n \cdot 2^{36} \rfloor$ holds, or equivalent to $\lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor \neq \lfloor n \cdot 2^{36} \rfloor$ must hold true when $2^{-1} \cdot 2^q \cdot 10^{-k-1} \neq n$. The following is proved by proof by contradiction.

Assume that $\lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor n \cdot 2^{36} \rfloor$ holds when $2^{-1} \cdot 2^q \cdot 10^{-k-1} \neq n$. Then there is:

$$\begin{aligned} \lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor &= \lfloor n \cdot 2^{36} \rfloor \\ \Rightarrow 0 < |2^{35} \cdot 2^q \cdot 10^{-k-1} - 2^{36} \cdot n| &< 1 \\ \Rightarrow 0 < |(2c-1) \cdot 2^{q-1} \cdot 10^{-k-1} - m| &< 2^{-36} \end{aligned} \quad (78)$$

As is known from equation (57), there is:

$$m - 1 < (2c-1) \cdot 2^{q-1} \cdot 10^{-k-1} < m + 1 \quad (79)$$

Suppose the decimal part of $(2c-1) \cdot 2^{q-1} \cdot 10^{-k-1}$ is represented as n^- , thus we have:

$$|(2c-1) \cdot 2^{q-1} \cdot 10^{-k-1} - m| = \begin{cases} n^-; & \text{if } (2c-1) \cdot 2^{q-1} \cdot 10^{-k-1} > m \\ 1 - n^-; & \text{if } (2c-1) \cdot 2^{q-1} \cdot 10^{-k-1} < m \end{cases} \quad (80)$$

Substitute expression (80) into expression (78), and we get:

$$\begin{aligned} 0 < |(2c-1) \cdot 2^{q-1} \cdot 10^{-k-1} - m| &< 2^{-36} \\ \Rightarrow 0 < n^- < 2^{-36} \text{ or } 0 < 1 - n^- &< 2^{-36} \end{aligned} \quad (81)$$

Similarly, it can be known that the double range is the range of n^- . Therefore, there is:

$$\begin{aligned} float : n^- &\in (0, 2^{-36}) \cup (1 - 2^{-36}, 1) \\ double : n^- &\in (0, 2^{-64}) \cup (1 - 2^{-64}, 1) \end{aligned} \quad (82)$$

When $2^{-1} \cdot 2^q \cdot 10^{-k-1} \neq n$, it is known from equation (54) that $(2c-1) \cdot 2^{q-1} \cdot 10^{-k-1}$ is not an integer. Therefore, there is:

$$0 < n^- < 1 \quad (83)$$

It is only necessary to prove that equation (82) does not hold. Discuss the range of the decimal part n^- when $(2c - 1) \cdot 2^{q-1} \cdot 10^{-k-1}$ is not an integer. According to equation (61), there are:

$$(2c - 1) \cdot 2^{q-1} \cdot 10^{-k-1} = (2c - 1) \cdot \frac{x}{y} = \begin{cases} \frac{(2c-1) \cdot 2^{q-k-2}}{5^{k+1}}; q \geq 2 \\ \frac{(2c-1)}{2^{2+k-q} \cdot 5^{k+1}}; 1 \geq q \geq 0 \\ \frac{(2c-1) \cdot 5^{-k-1}}{2^{2+k-q}}; q < 0 \end{cases} \quad (84)$$

The maximum value of $2c - 1$ is:

$$\begin{aligned} float : (2c - 1)_{\max} &= 2^{25} - 3 \\ double : (2c - 1)_{\max} &= 2^{54} - 3 \end{aligned} \quad (85)$$

Discuss based on the denominator range in equation (84).

- $y \leq (2c - 1)_{\max}$
When $y \leq (2c - 1)_{\max}$, y_{\max} is the expression (85), the following holds true:

$$\begin{aligned} \frac{1}{y_{\max}} \leq n^- &\leq 1 - \frac{1}{y_{\max}} \\ \frac{1}{y_{\max}} \leq 1 - n^- &\leq 1 - \frac{1}{y_{\max}} \end{aligned} \quad (86)$$

Therefore, when $y \leq (2c - 1)_{\max}$, equation (82) does not hold true.

- $y > (2c - 1)_{\max}$

Call function (37) to calculate the approximation results P_* / Q_* and P^* / Q^* of all possible upper and lower limit rational numbers:

$$\left(\frac{P_*}{Q_*}, \frac{P^*}{Q^*} \right) = f((2c - 1)_{\max}, x, y) \quad (87)$$

Therefore, for n^- , the following conclusion can be drawn from formula (36).

$$n^- \in \left[\frac{(Q_*x)\%y}{y}, \frac{(Q^*x)\%y}{y} \right] \quad (88)$$

By exhausting all possibilities, we thus have (the test code file is test3.py) :

$$\begin{aligned} float : 2^{-33} < n^- < 1 - 2^{-29} \\ double : 2^{-62} < n^- < 1 - 2^{-63} \end{aligned} \quad (89)$$

$$\begin{aligned} float : \left[\frac{(Q_*x)\%y}{y}, \frac{(Q^*x)\%y}{y} \right] \cap (0, 2^{-36}) &= \emptyset \\ \left[\frac{(Q_*x)\%y}{y}, \frac{(Q^*x)\%y}{y} \right] \cap (1 - 2^{-36}, 1) &= \emptyset \\ double : \left[\frac{(Q_*x)\%y}{y}, \frac{(Q^*x)\%y}{y} \right] \cap (0, 2^{-64}) &= \emptyset \\ \left[\frac{(Q_*x)\%y}{y}, \frac{(Q^*x)\%y}{y} \right] \cap (1 - 2^{-64}, 1) &= \emptyset \end{aligned} \quad (90)$$

Therefore, when $y > (2c - 1)_{\max}$, equation (82) does not hold true.

In summary, when $2^{-1} \cdot 2^q \cdot 10^{-k-1} \neq n$, equation (82) does not hold true, that is, $\lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor \neq \lfloor n \cdot 2^{36} \rfloor$ must hold true. Therefore, when $\lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor n \cdot 2^{36} \rfloor$ holds, $2^{-1} \cdot 2^q \cdot 10^{-k-1} = n$ must hold true. Therefore, equation (76) holds.

Similarly, it can be proved that when $\lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor 2^{36} - n \cdot 2^{36} \rfloor$ holds, $2^{-1} \cdot 2^q \cdot 10^{-k-1} = 1 - n$ must hold true. The same applies to double. Similarly, by proof of contradiction, for float, it is assumed that when $2^{-1} \cdot 2^q \cdot 10^{-k-1} \neq 1 - n$ holds, $\lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor 2^{36} - 2^{36} \cdot n \rfloor$ holds. That is:

$$\begin{aligned} \lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor &= \lfloor 2^{36} - n \cdot 2^{36} \rfloor \\ \Rightarrow 0 < &\left| 2^{35} \cdot 2^q \cdot 10^{-k-1} - 2^{36} + 2^{36} \cdot n \right| < 1 \\ \Rightarrow 0 < &\left| 2^{q-1} \cdot 10^{-k-1} - 1 + n \right| < 2^{-36} \\ \Rightarrow -2^{-36} < &(2c+1) \cdot 2^{q-1} \cdot 10^{-k-1} - m - 1 < 2^{-36} \end{aligned} \quad (91)$$

As is known from equation (58), there is:

$$m < (2c+1) \cdot 2^{q-1} \cdot 10^{-k-1} < m+2 \quad (92)$$

Suppose the decimal part of $(2c+1) \cdot 2^{q-1} \cdot 10^{-k-1}$ is represented as n^+ , thus we have:

$$(2c+1) \cdot 2^{q-1} \cdot 10^{-k-1} - m - 1 = \begin{cases} n^+; & \text{if } (2c+1) \cdot 2^{q-1} \cdot 10^{-k-1} > m+1 \\ 1 - n^+; & \text{if } (2c+1) \cdot 2^{q-1} \cdot 10^{-k-1} < m+1 \end{cases} \quad (93)$$

Substitute expression (93) into expression (91), and we get:

$$\begin{aligned} 0 < &\left| (2c+1) \cdot 2^{q-1} \cdot 10^{-k-1} - m - 1 \right| < 2^{-36} \\ \Rightarrow 0 < &1 - n^+ < 2^{-36} \text{ or } 0 < n^+ < 2^{-36} \end{aligned} \quad (94)$$

Similarly, it can be known that the double range is the range of n^+ . Therefore, there is:

$$\begin{aligned} \text{float : } n^+ &\in (0, 2^{-36}) \cup (1 - 2^{-36}, 1) \\ \text{double : } n^+ &\in (0, 2^{-64}) \cup (1 - 2^{-64}, 1) \end{aligned} \quad (95)$$

When $2^{-1} \cdot 2^q \cdot 10^{-k-1} \neq 1 - n$, it is known from equation (55) that $(2c+1) \cdot 2^{q-1} \cdot 10^{-k-1}$ is not an integer.

Therefore, there is:

$$0 < n^+ < 1 \quad (96)$$

It is only necessary to prove that equation (95) does not hold. Discuss the range of the decimal part n^+ when $(2c+1) \cdot 2^{q-1} \cdot 10^{-k-1}$ is not an integer. According to equation (61), there are:

$$(2c+1) \cdot 2^{q-1} \cdot 10^{-k-1} = (2c+1) \cdot \frac{x}{y} = \begin{cases} \frac{(2c+1) \cdot 2^{q-k-2}}{5^{k+1}}; & q \geq 2 \\ \frac{(2c+1)}{2^{2+k-q} \cdot 5^{k+1}}; & 1 \geq q \geq 0 \\ \frac{(2c+1) \cdot 5^{-k-1}}{2^{2+k-q}}; & q < 0 \end{cases} \quad (97)$$

The maximum value of $2c+1$ is:

$$\begin{aligned} \text{float : } (2c+1)_{\max} &= 2^{25} - 1 \\ \text{double : } (2c+1)_{\max} &= 2^{54} - 1 \end{aligned} \quad (98)$$

Discuss based on the denominator range in equation (97). 304

- $y \leq (2c + 1)_{\max}$ 305

When $y \leq (2c + 1)_{\max}$, y_{\max} is the expression (98), the following holds true: 306

$$\begin{aligned} \frac{1}{y_{\max}} &\leq n^+ \leq 1 - \frac{1}{y_{\max}} \\ \frac{1}{y_{\max}} &\leq 1 - n^+ \leq 1 - \frac{1}{y_{\max}} \end{aligned} \quad (99)$$

Therefore, when $y \leq (2c + 1)_{\max}$, equation (95) does not hold true. 307

- $y > (2c + 1)_{\max}$ 308

Call function (37) to calculate the approximation results P_* / Q_* and P^* / Q^* of all possible upper and lower limit rational numbers: 309

$$\left(\frac{P_*}{Q_*}, \frac{P^*}{Q^*} \right) = f((2c + 1)_{\max}, x, y) \quad (100)$$

Therefore, for n^+ , the following conclusion can be drawn from formula (36). 311

$$n^+ \in \left[\frac{(Q_*x)\%y}{y}, \frac{(Q^*x)\%y}{y} \right] \quad (101)$$

By exhausting all possibilities, we thus have (the test code file is test7.py): 312

$$\begin{aligned} float : 2^{-33} &< n^+ < 1 - 2^{-29} \\ double : 2^{-62} &< n^+ < 1 - 2^{-63} \end{aligned} \quad (102)$$

$$\begin{aligned} float : \left[\frac{(Q_*x)\%y}{y}, \frac{(Q^*x)\%y}{y} \right] \cap (0, 2^{-36}) &= \emptyset \\ \left[\frac{(Q_*x)\%y}{y}, \frac{(Q^*x)\%y}{y} \right] \cap (1 - 2^{-36}, 1) &= \emptyset \\ double : \left[\frac{(Q_*x)\%y}{y}, \frac{(Q^*x)\%y}{y} \right] \cap (0, 2^{-64}) &= \emptyset \\ \left[\frac{(Q_*x)\%y}{y}, \frac{(Q^*x)\%y}{y} \right] \cap (1 - 2^{-64}, 1) &= \emptyset \end{aligned} \quad (103)$$

Therefore, when $y > (2c + 1)_{\max}$, equation (95) does not hold true. 314

In summary, when $2^{-1} \cdot 2^q \cdot 10^{-k-1} \neq 1 - n$, equation (95) does not hold true, that is, $\lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor \neq \lfloor 2^{36} - n \cdot 2^{36} \rfloor$ must hold true. Therefore, when $\lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor 2^{36} - n \cdot 2^{36} \rfloor$ holds, $2^{-1} \cdot 2^q \cdot 10^{-k-1} = 1 - n$ must hold true. Therefore, equation (77) holds. 315
316
317
318

The following conclusions hold: 319

$$\begin{aligned} float : \lfloor 2^{36} - n \cdot 2^{36} \rfloor &= \begin{cases} 2^{36} - 1 - \lfloor 2^{36} \cdot n \rfloor; & \text{if } c \cdot 2^{36+q} \cdot 10^{-k-1} \notin Z \\ 2^{36} - \lfloor 2^{36} \cdot n \rfloor; & \text{if } c \cdot 2^{36+q} \cdot 10^{-k-1} \in Z \end{cases} \\ double : \lfloor 2^{64} - n \cdot 2^{64} \rfloor &= \begin{cases} 2^{64} - 1 - \lfloor 2^{64} \cdot n \rfloor; & \text{if } c \cdot 2^{64+q} \cdot 10^{-k-1} \notin Z \\ 2^{64} - \lfloor 2^{64} \cdot n \rfloor; & \text{if } c \cdot 2^{64+q} \cdot 10^{-k-1} \in Z \end{cases} \end{aligned} \quad (104)$$

Discuss whether the following equation (105) holds when conditions (68) and (69) are met: 320

$$\begin{aligned} \text{float : } & \lfloor c \cdot \frac{2^{q+35-k}}{5^{k+1}} \rfloor = \lfloor c \cdot \frac{2^{q+35-k}}{5^{k+1}} \cdot r \rfloor \\ &= \lfloor c \cdot \frac{2^{q+35-k}}{5^{k+1}} \cdot \frac{(2^{63-\lfloor(-k-1)\cdot\log_2(10)\rfloor}/10^{k+1})+1}{10^{-k-1}} \cdot 2^{\lfloor(-k-1)\cdot\log_2(10)\rfloor-63} \rfloor \quad (105) \\ \text{double : } & \lfloor c \cdot \frac{2^{q+63-k}}{5^{k+1}} \rfloor = \lfloor c \cdot \frac{2^{q+63-k}}{5^{k+1}} \cdot r \rfloor \\ &= \lfloor c \cdot \frac{2^{q+63-k}}{5^{k+1}} \cdot \frac{(2^{127-\lfloor(-k-1)\cdot\log_2(10)\rfloor}/10^{k+1})+1}{10^{-k-1}} \cdot 2^{\lfloor(-k-1)\cdot\log_2(10)\rfloor-127} \rfloor \end{aligned}$$

There are:

$$\begin{aligned} \text{float : } & \lfloor c \cdot \frac{2^{q+35-k}}{5^{k+1}} \rfloor = \lfloor 2^{36} \cdot (m+n) \rfloor = 2^{36}m + \lfloor 2^{36}n \rfloor \quad (106) \\ \text{double : } & \lfloor c \cdot \frac{2^{q+63-k}}{5^{k+1}} \rfloor = \lfloor 2^{64} \cdot (m+n) \rfloor = 2^{64}m + \lfloor 2^{64}n \rfloor \end{aligned}$$

It has been proven earlier that m can be accurately calculated. Then, when (1-118) holds true, the values $\lfloor 2^{36}n \rfloor$ and $\lfloor 2^{64}n \rfloor$ on the right side of equations (70) and (71) can be accurately calculated. 323 324

From equation (63), we have:

$$c = \frac{t \cdot 5^{k+1} - 1}{2} \quad (107)$$

Substituting equation (107) into equation (105), we have:

$$\begin{aligned} \text{float : } & c \cdot \frac{2^{q+35-k}}{5^{k+1}} = t \cdot 2^{q+34-k} - \frac{2^{q+34-k}}{5^{k+1}} \quad (108) \\ \text{double : } & c \cdot \frac{2^{q+63-k}}{5^{k+1}} = t \cdot 2^{q+62-k} - \frac{2^{q+62-k}}{5^{k+1}} \end{aligned}$$

When conditions (68) and (69) are met, $t \cdot 2^{q+34-k}$ and $t \cdot 2^{q+62-k}$ are integers. Under the condition of meeting condition (68), the decimal part of expression (108) is represented as: 327 328

$$\begin{aligned} \text{float : } & \frac{2^{q+34-k} \% 5^{k+1}}{5^{k+1}}; 2 \leq q \leq 33 \quad (109) \\ \text{double : } & \frac{2^{q+62-k} \% 5^{k+1}}{5^{k+1}}; 2 \leq q \leq 76 \end{aligned}$$

It is only necessary to prove that the increase in the value $c \cdot \frac{2^{q+35-k}}{5^{k+1}} \cdot r$ on the right side of the expression compared to the value $c \cdot \frac{2^{q+35-k}}{5^{k+1}}$ on the left side plus the decimal part of the value on the left side is less than 1 for equation (105) to hold true. That is: 329 330 331

$$\begin{aligned} \text{float : } & \frac{2^{q+34-k} \% 5^{k+1}}{5^{k+1}} + \left(c \cdot \frac{2^{q+35-k}}{5^{k+1}} \cdot r - c \cdot \frac{2^{q+35-k}}{5^{k+1}} \right) < 1 \quad (110) \\ \text{double : } & \frac{2^{q+62-k} \% 5^{k+1}}{5^{k+1}} + \left(c \cdot \frac{2^{q+63-k}}{5^{k+1}} \cdot r - c \cdot \frac{2^{q+63-k}}{5^{k+1}} \right) < 1 \end{aligned}$$

By exhaustively calculating the maximum possible c value under each q and substituting it into equation (110), it holds. The calculation result is in test2.py. The calculation results show that for the float range and the double range, equation (110) always holds true. Therefore, equation (105) holds true, and thus the values of $\lfloor 2^{36}n \rfloor$ and $\lfloor 2^{64}n \rfloor$ on the right 332 333 334 335

side of equations (70) and (71) can be accurately calculated. The values of $\lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor$ and $\lfloor 2^{63} \cdot 2^q \cdot 10^{-k-1} \rfloor$ on the left side of equations (70) and (71) can be calculated through lookup tables.

$$\begin{aligned} float : \lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor &= pow10 \gg (28 - q - \lfloor (-k - 1) \cdot \log_2(10) \rfloor) \\ double : \lfloor 2^{64} \cdot 2^q \cdot 10^{-k-1} \rfloor &= pow10 \gg (64 - q - \lfloor (-k - 1) \cdot \log_2(10) \rfloor) \end{aligned} \quad (111)$$

The code file for verifying the validity of equation (111) is test4.py. Therefore, when conditions (68) and (69) are met, the values of both sides of equations (70) and (71) can be accurately calculated.

Discuss the relationship between the following two values within all ranges of floating-point numbers:

$$\begin{aligned} float : \lfloor c \cdot 2^{q+36} \cdot 10^{-k-1} \rfloor; \lfloor c \cdot 2^{q+36} \cdot r \cdot 10^{-k-1} \rfloor; \\ double : \lfloor c \cdot 2^{q+64} \cdot 10^{-k-1} \rfloor; \lfloor c \cdot 2^{q+64} \cdot r \cdot 10^{-k-1} \rfloor; \end{aligned} \quad (112)$$

When $r = 1$, it is obvious that the two values in expression (112) are equal. When $r \neq 1$, or equivalent to $r > 1$, has:

$$\begin{aligned} float : \\ c \cdot 2^{q+36} \cdot r \cdot 10^{-k-1} &= c \cdot 2^{q+36} \cdot 10^{-k-1} + c \cdot 2^{q+36} \cdot (r - 1) \cdot 10^{-k-1} \\ &< c \cdot 2^{q+36} \cdot 10^{-k-1} + 2^{24} \cdot 2^{36} \cdot 2^q \cdot 10^{-k-1} \cdot (r - 1) \\ &< c \cdot 2^{q+36} \cdot 10^{-k-1} + 2^{-3} \\ \lfloor c \cdot 2^{q+36} \cdot r \cdot 10^{-k-1} \rfloor &\leq \lfloor c \cdot 2^{q+36} \cdot 10^{-k-1} \rfloor + 1 \quad (113) \\ double : \\ c \cdot 2^{q+64} \cdot r \cdot 10^{-k-1} &= c \cdot 2^{q+64} \cdot 10^{-k-1} + c \cdot 2^{q+64} \cdot (r - 1) \cdot 10^{-k-1} \\ &< c \cdot 2^{q+64} \cdot 10^{-k-1} + 2^{53} \cdot 2^{64} \cdot 2^q \cdot 10^{-k-1} \cdot (r - 1) \\ &< c \cdot 2^{q+64} \cdot 10^{-k-1} + 2^{-10} \\ \lfloor c \cdot 2^{q+64} \cdot r \cdot 10^{-k-1} \rfloor &\leq \lfloor c \cdot 2^{q+64} \cdot 10^{-k-1} \rfloor + 1 \end{aligned}$$

Therefore, there is:

$$\begin{aligned} float : 0 &\leq \lfloor c \cdot 2^{q+36} \cdot r \cdot 10^{-k-1} \rfloor - \lfloor c \cdot 2^{q+36} \cdot 10^{-k-1} \rfloor \leq 1 \\ double : 0 &\leq \lfloor c \cdot 2^{q+64} \cdot r \cdot 10^{-k-1} \rfloor - \lfloor c \cdot 2^{q+64} \cdot 10^{-k-1} \rfloor \leq 1 \end{aligned} \quad (114)$$

Because there is:

$$\lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor c \cdot 2^q \cdot r \cdot 10^{-k-1} \rfloor = m \quad (115)$$

$$\begin{aligned} float : \lfloor c \cdot 2^{q+36} \cdot 10^{-k-1} \rfloor &= 2^{36}m + \lfloor 2^{36}n \rfloor \\ double : \lfloor c \cdot 2^{q+64} \cdot 10^{-k-1} \rfloor &= 2^{64}m + \lfloor 2^{64}n \rfloor \end{aligned} \quad (116)$$

Suppose:

$$n_r = c \cdot 2^q \cdot r \cdot 10^{-k-1} - m \quad (117)$$

Therefore, the following conclusion can be drawn: when condition (68) is met, from equation (105), we have:

$$\begin{aligned} float : 2 \leq q \leq 33 \&\& (2c \pm 1)\%5^{k+1} = 0 \Rightarrow \lfloor 2^{36} \cdot n_r \rfloor = \lfloor 2^{36} \cdot n_r \rfloor \\ double : 2 \leq q \leq 76 \&\& (2c \pm 1)\%5^{k+1} = 0 \Rightarrow \lfloor 2^{64} \cdot n_r \rfloor = \lfloor 2^{64} \cdot n_r \rfloor \end{aligned} \quad (118)$$

Within the range of floating-point numbers, there exists:

$$\begin{aligned} \text{float} : \lfloor 2^{36} \cdot n \rfloor &\leq \lfloor 2^{36} \cdot n_r \rfloor \leq \lfloor 2^{36} \cdot n \rfloor + 1 \\ \text{double} : \lfloor 2^{64} \cdot n \rfloor &\leq \lfloor 2^{64} \cdot n_r \rfloor \leq \lfloor 2^{64} \cdot n \rfloor + 1 \end{aligned} \quad (119)$$

To simplify the expression, *even* is used to indicate whether c is an even number:

$$\text{even} = (c + 1) \% 2 \in \{0, 1\} \quad (120)$$

The following will introduce the proof process of the float range, from equation (121) to equation (144). When $2^{-1} \cdot 2^q \cdot 10^{-k-1} = n$ or $2^{-1} \cdot 2^q \cdot 10^{-k-1} = 1 - n$, $2^{-1} \cdot 2^q \cdot 10^{-k-1} = n$ is the boundary condition for *one* = 0, and $2^{-1} \cdot 2^q \cdot 10^{-k-1} = 1 - n$ is the boundary condition for *one* = 10. Whether *one* is 0 or 10 is determined based on whether c is an even number. Therefore, the following exists:

$$\text{float} : \begin{cases} \text{one} = 0 : \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + \text{even} > \lfloor n_r \cdot 2^{36} \rfloor \\ \text{one} = 10 : \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + \text{even} > 2^{36} - 1 - \lfloor n_r \cdot 2^{36} \rfloor \end{cases} \quad (121)$$

Therefore, when $2^{-1} \cdot 2^q \cdot 10^{-k-1} = n$ or $2^{-1} \cdot 2^q \cdot 10^{-k-1} = 1 - n$, we can use the condition (122) to determine whether *one* = 0 or *one* = 10.

$$\text{float} : \begin{cases} \text{if } \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + \text{even} > \lfloor n_r \cdot 2^{36} \rfloor : \text{one} = 0 \\ \text{if } \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + \text{even} > 2^{36} - 1 - \lfloor n_r \cdot 2^{36} \rfloor : \text{one} = 10 \end{cases} \quad (122)$$

When $2^{-1} \cdot 2^q \cdot 10^{-k-1} > n$ or $2^{-1} \cdot 2^q \cdot 10^{-k-1} > 1 - n$, We can also use the above condition (122) to determine whether *one* = 0 or *one* = 10. When $2^{-1} \cdot 2^q \cdot 10^{-k-1} < n$ or $2^{-1} \cdot 2^q \cdot 10^{-k-1} < 1 - n$, we can also use the above condition (122) to determine whether *one* ≠ 0 or *one* ≠ 10. The proof is as follows:

When $2^{-1} \cdot 2^q \cdot 10^{-k-1} < n$, there must exist *one* ≠ 0, and there is:

$$2^{-1} \cdot 2^q \cdot 10^{-k-1} - n = n^- - 1 \in (2^{-33} - 1, -2^{-29}) \quad (123)$$

Therefore, the following exists:

$$2^{q+35} \cdot 10^{-k-1} - 2^{36} \cdot n \in (2^3 - 2^{36}, -2^7) \quad (124)$$

Suppose there are two real numbers a and b , and the following relationship must exist:

$$\begin{aligned} 0 &\leq b - \lfloor b \rfloor < 1 \\ a - \lfloor a \rfloor - 1 &< b - \lfloor b \rfloor < 1 + a - \lfloor a \rfloor \\ a - b - 1 &< \lfloor a \rfloor - \lfloor b \rfloor < a - b + 1 \end{aligned} \quad (125)$$

When $a = 2^{q+35} \cdot 10^{-k-1}$ and $b = 2^{36} \cdot n$, the following exists:

$$\lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor - \lfloor 2^{36} \cdot n \rfloor < 2^{q+35} \cdot 10^{-k-1} - 2^{36} \cdot n + 1 \quad (126)$$

From equation (124), we have:

$$\lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor - \lfloor 2^{36} \cdot n \rfloor < 1 - 2^7 < 0 \quad (127)$$

Therefore, there is:

$$\begin{aligned} \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + even &\leq \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + 1 \\ &< \lfloor 2^{36} \cdot n \rfloor \leq \lfloor 2^{36} \cdot n_r \rfloor \\ \Rightarrow \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + even &< \lfloor 2^{36} \cdot n_r \rfloor \end{aligned} \quad (128)$$

Therefore, when $2^{-1} \cdot 2^q \cdot 10^{-k-1} < n$, the condition (122) can be used to determine that one $\neq 0$.

When $2^{-1} \cdot 2^q \cdot 10^{-k-1} > n$, there must exist one = 0, and there is:

$$2^{-1} \cdot 2^q \cdot 10^{-k-1} - n = n^- \in (2^{-33}, 1 - 2^{-29}) \quad (129)$$

Therefore, the following exists:

$$2^{q+35} \cdot 10^{-k-1} - 2^{36} \cdot n \in (2^3, 2^{36} - 2^7) \quad (130)$$

When $a = 2^{q+35} \cdot 10^{-k-1}$ and $b = 2^{36} \cdot n$, from equation (125), the following exists:

$$\lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor - \lfloor 2^{36} \cdot n \rfloor > 2^{q+35} \cdot 10^{-k-1} - 2^{36} \cdot n - 1 \quad (131)$$

From equation (130), we have:

$$\lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor - \lfloor 2^{36} \cdot n \rfloor > 2^3 - 1 \geq 0 \quad (132)$$

Therefore, there is:

$$\begin{aligned} \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + even &\geq \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor \\ &> \lfloor 2^{36} \cdot n \rfloor + 1 \geq \lfloor 2^{36} \cdot n_r \rfloor \\ \Rightarrow \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + even &> \lfloor 2^{36} \cdot n_r \rfloor \end{aligned} \quad (133)$$

Therefore, when $2^{-1} \cdot 2^q \cdot 10^{-k-1} > n$, the condition (122) can be used to determine that one = 0.

When $2^{-1} \cdot 2^q \cdot 10^{-k-1} < 1 - n$, there must exist one $\neq 10$, and there is:

$$2^{-1} \cdot 2^q \cdot 10^{-k-1} + n = n^+ \in (2^{-33}, 1 - 2^{-29}) \quad (134)$$

Therefore, the following exists:

$$2^{q+35} \cdot 10^{-k-1} + 2^{36} \cdot n \in (2^3, 2^{36} - 2^7) \quad (135)$$

Suppose there are two real numbers a and b , and the following relationship must exist:

$$\begin{aligned} a - 1 < \lfloor a \rfloor &\leq a \\ b - 1 < \lfloor b \rfloor &\leq b \\ a + b - 2 < \lfloor a \rfloor + \lfloor b \rfloor &\leq a + b \end{aligned} \quad (136)$$

When $a = 2^{q+35} \cdot 10^{-k-1}$ and $b = 2^{36} \cdot n$, the following exists:

$$\lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + \lfloor 2^{36} \cdot n \rfloor \leq 2^{q+35} \cdot 10^{-k-1} + 2^{36} \cdot n \quad (137)$$

From equation (135), we have:

$$\lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + \lfloor 2^{36} \cdot n \rfloor < 2^{36} - 2^7 \quad (138)$$

Therefore, there is:

$$\begin{aligned} \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + even &\leq \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + 1 \\ &< 2^{36} - 2 - \lfloor 2^{36} \cdot n \rfloor \\ &\leq 2^{36} - 1 - \lfloor 2^{36} \cdot n_r \rfloor \\ \Rightarrow \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + even &< 2^{36} - 1 - \lfloor 2^{36} \cdot n_r \rfloor \end{aligned} \quad (139)$$

Therefore, when $2^{-1} \cdot 2^q \cdot 10^{-k-1} < 1 - n$, the condition (122) can be used to determine that one $\neq 10$.

When $2^{-1} \cdot 2^q \cdot 10^{-k-1} > 1 - n$, there must exist one = 10, and there is:

$$2^{-1} \cdot 2^q \cdot 10^{-k-1} + n = n^+ + 1 \in (1 + 2^{-33}, 2 - 2^{-29}) \quad (140)$$

Therefore, the following exists:

$$2^{q+35} \cdot 10^{-k-1} + 2^{36} \cdot n \in (2^3 + 2^{36}, 2^{37} - 2^7) \quad (141)$$

When $a = 2^{q+35} \cdot 10^{-k-1}$ and $b = 2^{36} \cdot n$, from equation (136), the following exists:

$$\lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + \lfloor 2^{36} \cdot n \rfloor > 2^{q+35} \cdot 10^{-k-1} + 2^{36} \cdot n - 2 \quad (142)$$

From equation (141), we have:

$$\lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + \lfloor 2^{36} \cdot n \rfloor > 2^{36} + 2^3 - 2 \geq 2^{36} \quad (143)$$

Therefore, there is:

$$\begin{aligned} \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + even &\geq \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor \\ &> 2^{36} - 2 - \lfloor 2^{36} \cdot n \rfloor \\ &\geq 2^{36} - 1 - \lfloor 2^{36} \cdot n_r \rfloor \\ \Rightarrow \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + even &> 2^{36} - 1 - \lfloor 2^{36} \cdot n_r \rfloor \end{aligned} \quad (144)$$

Therefore, when $2^{-1} \cdot 2^q \cdot 10^{-k-1} > 1 - n$, the condition (122) can be used to determine that one = 10.

From the above proof, it can be seen that when condition (68) is met, the condition (122) can be used to determine whether one = 0 or one = 10 when $2^{-1} \cdot 2^q \cdot 10^{-k-1} = n$ or $2^{-1} \cdot 2^q \cdot 10^{-k-1} = 1 - n$. When $2^{-1} \cdot 2^q \cdot 10^{-k-1} > n$ or $2^{-1} \cdot 2^q \cdot 10^{-k-1} > 1 - n$, the condition (122) can be used to determine whether one = 0 or one = 10. When $2^{-1} \cdot 2^q \cdot 10^{-k-1} < n$ or $2^{-1} \cdot 2^q \cdot 10^{-k-1} < 1 - n$, the condition (122) can be used to determine whether one $\neq 0$ or one $\neq 10$.

For the double range, the value of one can be calculated based on the following conditions.

$$double : \begin{cases} \text{if } \lfloor 2^{q+64} \cdot 10^{-k-1} \rfloor + even > \lfloor n_r \cdot 2^{64} \rfloor : \text{one} = 0 \\ \text{if } \lfloor 2^{q+64} \cdot 10^{-k-1} \rfloor + even > 2^{64} - 1 - \lfloor n_r \cdot 2^{64} \rfloor : \text{one} = 10 \end{cases} \quad (145)$$

Readers can deduce it by themselves based on the above proof process. The proof process of this section is completed. In the code implementation, the two judgment conditions can

be quickly calculated using addition and subtraction shift operations, and can be compiled by the compiler into cmov instructions, thereby reducing the impact of branch prediction failure on performance.

3.5. Determine whether $one = \lfloor 10n \rfloor$ or $one = \lfloor 10n \rfloor + 1$

Determine whether one is $\lfloor 10n \rfloor$ or $\lfloor 10n \rfloor + 1$ based on the decimal part of $10n$. There are two cases: the decimal part of $10n$ is 0.5 and it is not 0.5.

3.5.1. $10n - \lfloor 10n \rfloor = 0.5$

When the decimal part of $10n$ is 0.5, there must be:

$$\begin{aligned} 10n - \lfloor 10n \rfloor &= 0.5 \\ \Rightarrow 10 \cdot c \cdot 2^q \cdot 10^{-k-1} - \lfloor 10 \cdot c \cdot 2^q \cdot 10^{-k-1} \rfloor &= 0.5 \\ \Rightarrow c \cdot 2^q \cdot 10^{-k} - \lfloor c \cdot 2^q \cdot 10^{-k} \rfloor &= 0.5 \\ \Rightarrow c \cdot 2^q \cdot 10^{-k} &= \lfloor c \cdot 2^q \cdot 10^{-k} \rfloor + 0.5 \\ \Rightarrow 2c \cdot 2^q \cdot 10^{-k} &= 2\lfloor c \cdot 2^q \cdot 10^{-k} \rfloor + 1 \end{aligned} \quad (146)$$

So $2c \cdot 2^q \cdot 10^{-k}$ is an odd number. Then the following expression is odd:

$$c \cdot 2^{q+1} \cdot 10^{-k} = c \cdot 2^{q-k+1} \cdot 5^{-k} \quad (147)$$

According to the range of q , there are:

$$c \cdot 2^{q+1} \cdot 10^{-k} = \begin{cases} \frac{c \cdot 2^{q-k+1}}{5^k}; q \geq 0 \\ c \cdot 2 \cdot 5^{-k}; q = -1 \\ \frac{c \cdot 5^{-k}}{2^{k-q-1}}; q \leq -2 \end{cases} \quad (148)$$

According to the range of q , the following situations are discussed:

- $q \geq 0$
When $q \geq 0$, it can be concluded that $q - k + 1 \geq 1$, the numerator $c \cdot 2^{q-k+1}$ is even and the denominator 5^k is odd, which does not meet the condition.
- $q = -1$
When $q = -1$, it can be concluded that $c \cdot 2 \cdot 5^{-k}$ is even, which does not meet the condition.
- $q \leq -2$
 5^{-k} is an odd number. c is an odd multiple of 2^{k-q-1} . So:

$$\begin{aligned} float : c \geq 2^{k-q-1} \Rightarrow k - q - 1 &\leq 22 \Rightarrow q \geq -34 \\ double : c \geq 2^{k-q-1} \Rightarrow k - q - 1 &\leq 51 \Rightarrow q \geq -75 \end{aligned} \quad (149)$$

Therefore, when q meets the above conditions, c must be an odd multiple of 2^{k-q-1} to meet the condition. Therefore, when the following conditions are met, expression (147) is an odd number:

$$\begin{aligned} float : -34 \leq q \leq -2 \&& c \% 2^{k-q} = 2^{k-q-1} \\ double : -75 \leq q \leq -2 \&& c \% 2^{k-q} = 2^{k-q-1} \end{aligned} \quad (150)$$

When q is within the above range (150), $r = 1$ is derived from equation (30). Therefore, there is:

$$n_r = n \quad (151)$$

The following equation holds:

$$20m + 20n = c \cdot 2^q \cdot 10^{-k+1} = c \cdot 2^{q-k+1} \cdot 5^{-k} = \frac{c}{2^{k-q-1}} \cdot 5^{-k} \quad (152)$$

Since $-k \geq 1$, 5^{-k} is multiple of 5 and is an odd number. Since $\frac{c}{2^{k-q-1}}$ and 5^{-k} are both odd numbers, $20m$ is an even number, $20n$ is multiple of 5 and is an odd number. Therefore, there is:

$$\begin{aligned} 20n &\in \{5, 15\} \\ \Rightarrow n &\in \{0.25, 0.75\} \\ \Rightarrow n_r &\in \{0.25, 0.75\} \end{aligned} \quad (153)$$

The result of *one* is an even number between $\lfloor 10n \rfloor$ and $\lfloor 10n \rfloor + 1$. Therefore, when the following conditions are met:

$$one = \begin{cases} \lfloor 10n \rfloor = 2, \text{if } n = 0.25 \\ \lfloor 10n \rfloor + 1 = 8, \text{if } n = 0.75 \end{cases} \Rightarrow one = \lfloor 20n + 1 \rfloor // 2 - (n = 0.25 ? 1 : 0) \quad (154)$$

3.5.2. $10n - \lfloor 10n \rfloor \neq 0.5$

When the decimal part of $10n$ is not 0.5, round to the nearest integer value based on the decimal part of $10n$. Therefore, there is:

$$one = \begin{cases} \lfloor 10n \rfloor, \text{if } 10n - \lfloor 10n \rfloor < 0.5 \\ \lfloor 10n \rfloor + 1, \text{if } 10n - \lfloor 10n \rfloor > 0.5 \end{cases} \Rightarrow one = \lfloor 10n + 0.5 \rfloor = \lfloor 20n + 1 \rfloor // 2 \quad (155)$$

Since $\lfloor 20n + 1 \rfloor = \lfloor 20n \rfloor + 1$, it is only necessary to accurately calculate the value of $\lfloor 20n \rfloor$. And, there is:

$$\begin{aligned} d &= ten + one \\ &= 10m + \lfloor 20n + 1 \rfloor // 2 \\ &= (\lfloor 20m + 20n \rfloor + 1) // 2 \end{aligned} \quad (156)$$

Suppose there are:

$$20m + 20n = c \cdot 2^{q+1} \cdot 10^{-k} = c \cdot 2^{q-k+1} \cdot 5^{-k} = c \cdot \frac{x}{y} \quad (157)$$

Suppose the decimal part of $20n$ is n_{20} .

When $y \leq c_{\max} = C$, the range of the decimal part must include:

$$\begin{aligned} float: \frac{1}{2^{24}-1} &= \frac{1}{C} \leq n_{20} \leq 1 - \frac{1}{C} = \frac{2^{24}-2}{2^{24}-1} \\ double: \frac{1}{2^{53}-1} &= \frac{1}{C} \leq n_{20} \leq 1 - \frac{1}{C} = \frac{2^{53}-2}{2^{53}-1} \end{aligned} \quad (158)$$

When $y > c_{\max} = C$, the range of the decimal part must include(the test file is test5.py):

$$\begin{aligned} float: 2^{-32} &< n_{20} < 1 - 2^{-30} \\ double: 2^{-64} &< n_{20} < 1 - 2^{-62} \end{aligned} \quad (159)$$

Therefore, the range of n_{20} satisfies equation (159). In the code implementation, for float, only the high 36 bits of n_r are retained, and for double, only the high 70 bits of n_r are

retained. Suppose the discarded part of a float is represented as n_{36} , and similarly, the discarded part of a double is represented as n_{70} . Therefore, there is:

446

$$\begin{aligned} \text{float : } n_{36} &\in [0, 2^{-36}) \\ \text{double : } n_{70} &\in [0, 2^{-70}) \end{aligned} \quad (160)$$

Calculate the boundary conditions of the following expression:

448

$$\begin{aligned} \text{float : } F &= 20 \cdot (c \cdot 2^q \cdot r \cdot 10^{-k-1} - n_{36}) \\ \text{double : } F &= 20 \cdot (c \cdot 2^q \cdot r \cdot 10^{-k-1} - n_{70}) \end{aligned} \quad (161)$$

Therefore, there is:

449

$$\begin{aligned} \text{float : } F_{\min} &> 20 \cdot (c \cdot 2^q \cdot 10^{-k-1} - 2^{-36}) \\ &= 20m + 20n - 20 \cdot 2^{-36} \\ F_{\max} &< 20 \cdot (c \cdot 2^q \cdot (1 + 2^{-63}) \cdot 10^{-k-1} - 0) \\ &< 20m + 20n + 20 \cdot 2^{-63} \cdot c \\ &< 20m + \lfloor 20n \rfloor + 1 \\ \text{double : } F_{\min} &> 20 \cdot (c \cdot 2^q \cdot 10^{-k-1} - 2^{-70}) \\ &= 20m + 20n - 20 \cdot 2^{-70} \\ &> 20m + \lfloor 20n \rfloor \\ F_{\max} &< 20 \cdot (c \cdot 2^q \cdot (1 + 2^{-127}) \cdot 10^{-k-1} - 0) \\ &< 20m + 20n + 20 \cdot 2^{-127} \cdot c \\ &< 20m + \lfloor 20n \rfloor + 1 \end{aligned} \quad (162)$$

Therefore, there is:

450

$$\begin{aligned} \text{float : } \lfloor F \rfloor &= 20m + \lfloor 20n \rfloor \\ \text{double : } \lfloor F \rfloor &= 20m + \lfloor 20n \rfloor \end{aligned} \quad (163)$$

In fact, in the above proof process, for float, $\lfloor F_{\min} \rfloor \neq 20m + \lfloor 20n \rfloor$ may exist, but the code implementation has passed the exhaustive test, so this not-so-perfect proof process can be ignored. Therefore, the calculation of d can be simplified as follows:

451

452

453

$$\begin{aligned} d &= \text{ten} + \text{one} \\ &= (\lfloor F \rfloor + 1) // 2 \\ &= (\lfloor 20 \cdot (c \cdot 2^q \cdot r \cdot 10^{-k-1} - n_x) \rfloor + 1) // 2 \end{aligned} \quad (164)$$

For the float range, $n_x = n_{36}$; for the double range, $n_x = n_{70}$.

454

For double, quickly determine that $n == 0.25$ in equation (154).

455

When $n = 0.25, \lfloor 2^{64} \cdot n_r \rfloor = \lfloor 2^{64} \cdot n \rfloor = 2^{62}$. Therefore, the following condition can be used to quickly determine whether $n = 0.25$:

456

457

$$\text{double : } n = 0.25 \text{ if } \lfloor 2^{64} \cdot n_r \rfloor = 2^{62} \quad (165)$$

When $n \neq 0.25$, calculate the range of the decimal part of the following expression:

458

$$4m + 4n = c \cdot 2^{q+2} \cdot 10^{-k-1} \quad (166)$$

Therefore, when equation (166) is not an integer, we have:(test6.py):

459

$$2^{-62} < 4n - \lfloor 4n \rfloor < 1 - 2^{-62} \quad (167)$$

Calculate the two boundary cases of $4n$ that are closest to 1:

460

$$\begin{aligned} \lfloor 4n \rfloor = 0 &\Rightarrow 4n - 0 < 1 - 2^{-62} \Rightarrow \lfloor 2^{64} \cdot n \rfloor \leq 2^{62} - 2 \\ \lfloor 4n \rfloor = 1 &\Rightarrow 4n - 1 > 2^{-62} \Rightarrow \lfloor 2^{64} \cdot n \rfloor \geq 2^{62} + 1 \end{aligned} \quad (168)$$

Then there are:

461

$$\begin{aligned} \lfloor 2^{64} \cdot n \rfloor &\neq 2^{62} \&& \lfloor 2^{64} \cdot n \rfloor + 1 \neq 2^{62} \\ &\Rightarrow \lfloor 2^{64} \cdot n_r \rfloor \neq 2^{62} \end{aligned} \quad (169)$$

Therefore, the following condition can be used to quickly determine whether $n \neq 0.25$:

462

$$double : n \neq 0.25 \text{ if } \lfloor 2^{64} \cdot n_r \rfloor \neq 2^{62} \quad (170)$$

In summary, for double, the following condition can be used to quickly determine whether $n = 0.25$:

463

$$\begin{aligned} double : n &= 0.25 \text{ if } \lfloor 2^{64} \cdot n_r \rfloor = 2^{62} \\ double : n &\neq 0.25 \text{ if } \lfloor 2^{64} \cdot n_r \rfloor \neq 2^{62} \end{aligned} \quad (171)$$

464

In the double range, introduce a fast way to calculate one:

465

$$double : one = \lfloor \frac{\lfloor 2^{64} n_r \rfloor}{2^{64}} \cdot 10 + (n = 0.25) ? 0 : \left(2^{-1} + \frac{6}{2^{64}} \right) \rfloor \quad (172)$$

The proof of equation (172) is as follows:

466

when $n = 0.25$, $\lfloor \frac{\lfloor 2^{64} n_r \rfloor}{2^{64}} \cdot 10 \rfloor = \lfloor 10n \rfloor = 2$;

467

when $n \neq 0.25$, equation (172) can be equivalent to the following:

468

$$double : one = \lfloor \frac{\lfloor 2^{64} n_r \rfloor}{2^{64}} \cdot 10 + 2^{-1} + \frac{6}{2^{64}} \rfloor \quad (173)$$

According to the $10n - \lfloor 10n \rfloor$ range, one is represented as:

469

$$double : one = \begin{cases} \lfloor 10n \rfloor, & \text{if } 10n - \lfloor 10n \rfloor < 0.5 \\ 8, & \text{if } 10n - \lfloor 10n \rfloor = 0.5 \\ \lfloor 10n \rfloor + 1, & \text{if } 10n - \lfloor 10n \rfloor > 0.5 \end{cases} = \lfloor 20n + 1 \rfloor // 2 \quad (174)$$

Therefore, when $n \neq 0.25$, we need to prove that the following equation holds:

470

$$\lfloor \frac{\lfloor 2^{64} n_r \rfloor}{2^{64}} \cdot 10 + 2^{-1} + \frac{6}{2^{64}} \rfloor = \begin{cases} \lfloor 10n \rfloor, & \text{if } 10n - \lfloor 10n \rfloor < 0.5 \\ 8, & \text{if } 10n - \lfloor 10n \rfloor = 0.5 \\ \lfloor 10n \rfloor + 1, & \text{if } 10n - \lfloor 10n \rfloor > 0.5 \end{cases} = \lfloor 20n + 1 \rfloor // 2 \quad (175)$$

From the range of n , there is:

471

$$\frac{\lfloor 2^{64} n_r \rfloor}{2^{64}} \in (n_r - 2^{-64}, n_r] \quad (176)$$

Because the following conditions exist:

472

$$\begin{aligned} c \cdot 2^q \cdot 10^{-k-1} &= m + n \\ c \cdot 2^q \cdot r \cdot 10^{-k-1} &= m + n, \end{aligned} \quad (177)$$

473

Therefore, the following relationship can be concluded:

$$\begin{aligned}
 n_r - n &= (r - 1) \cdot c \cdot 2^q \cdot 10^{-k-1} \\
 n_r &= (r - 1) \cdot (m + n) + n \\
 \Rightarrow n &\leq n_r < 2^{-127} \cdot c + n \\
 n &\leq n_r < 2^{-127} \cdot 2^{53} + n \\
 n &\leq n_r < 2^{-74} + n
 \end{aligned} \tag{178}$$

From equation (176) and (178), it can be concluded that:

$$\begin{aligned}
 \frac{\lfloor 2^{64}n_r \rfloor}{2^{64}} &\in (n - 2^{-64}, n + 2^{-74}) \\
 \Rightarrow \frac{\lfloor 2^{64}n_r \rfloor}{2^{64}} \cdot 10 &\in (10n - 10 \cdot 2^{-64}, 10n + 10 \cdot 2^{-74}) \\
 \Rightarrow \frac{\lfloor 2^{64}n_r \rfloor}{2^{64}} \cdot 20 &\in (20n - 20 \cdot 2^{-64}, 20n + 20 \cdot 2^{-74}) \\
 \Rightarrow \frac{\lfloor 2^{64}n_r \rfloor}{2^{64}} \cdot 20 &\in ([20n] + n_{20} - 20 \cdot 2^{-64}, [20n] + n_{20} + 20 \cdot 2^{-74})
 \end{aligned} \tag{179}$$

Discuss the range of values of x when the following conditions are met.

$$\lfloor \frac{\lfloor 2^{64}n_r \rfloor}{2^{64}} \cdot 20 + 1 + x \rfloor // 2 = \lfloor 20n + 1 \rfloor // 2 = one \tag{180}$$

Therefore, the following conclusions can be drawn:

$$\begin{aligned}
 \lfloor 20n \rfloor + n_{20} - 20 \cdot 2^{-64} + 1 + x &\geq \lfloor 20n + 1 \rfloor \Rightarrow x \geq 20 \cdot 2^{-64} - n_{20} \\
 \lfloor 20n \rfloor + n_{20} + 20 \cdot 2^{-74} + 1 + x &< \lfloor 20n + 2 \rfloor \Rightarrow x < 1 - 20 \cdot 2^{-74} - n_{20}
 \end{aligned} \tag{181}$$

Suppose $x = 12 \cdot 2^{-64}$. Through the exhaustive method, all floating-point numbers that do not meet the following conditions can be obtained.

$$x = 12 \cdot 2^{-64} \geq 20 \cdot 2^{-64} - n_{20} \tag{182}$$

All floating-point numbers that do not meet condition (182) are as follows (in hexadecimal) :

$$\begin{aligned}
 &0xd17c0747bd76fa1, \\
 &0xd27c0747bd76fa1, \\
 &0x4d73de005bd620df, \\
 &0x4d83de005bd620df, \\
 &0x4d93de005bd620df,
 \end{aligned} \tag{183}$$

Through the exhaustive method, all floating-point numbers that do not meet the following conditions can be obtained.

$$x = 12 \cdot 2^{-64} < 1 - 20 \cdot 2^{-74} - n_{20} \tag{184}$$

All floating-point numbers that do not meet condition (184) are as follows (in hexadecimal)

:

483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509

$$\begin{aligned} & 0x612491daad0ba280, \\ & 0x6159b651584e8b20, \\ & 0x619011f2d73116f4, \\ & 0x61c4166f8cf8cb1, \\ & 0x61d4166f8cf8cb1, \end{aligned} \quad (185)$$

There are:

$$2\left(\frac{\lfloor 2^{64}n_r \rfloor}{2^{64}} \cdot 10 + 2^{-1} + \frac{6}{2^{64}}\right) = \frac{\lfloor 2^{64}n_r \rfloor}{2^{64}} \cdot 20 + 1 + x \quad (186)$$

When the floating-point number is not within the above range (183) and (185), the condition (181) is satisfied. We have tested all floating-point numbers within the above-mentioned range (183) and (185), and the algorithm implementation code has output the correct result, that is, it satisfies the SW principle.

In summary, equation (175) and equation (172) holds. Therefore, equation (172) can be used to quickly calculate *one*.

3.6. Irregular Number

Due to the limited and small number of irregular floating-point numbers, there are a total of 2046 double floating-point numbers and 254 float floating-point numbers. The correctness of the algorithm code in this paper can be proved by the exhaustive method. Therefore, it is not introduced in this article. For the specific implementation process, please refer to the source code.

4. Benchmark result

The link to the benchmark project is https://github.com/xjb714/f2dec_bench. In the performance test comparison, we compared the time spent by the following several different algorithms converting floating-point numbers to decimal results, as shown in Table (1). Test process: Generate 10^7 random numbers without 0, NaN, and Inf, measure the total time spent converting all floating-point numbers to decimal results, and obtain the average time for converting a single floating-point number to decimal. We conducted benchmark tests on two processors, and the test results are shown in table (2) and (3). The algorithm proposed in this paper supports parallel computing and can be compiled into parallelized code by a compatible compiler. However, empirical evaluation indicates that, as of now, only the Clang compiler generates AVX-512 or AVX-2 instruction sets from this algorithm when targeting the x86-64 architecture.

Table 2. Benchmark results on AMD R7 7840H and Ubuntu 24.04. The unit is nanosecond(ns).

algorithm	float			double		
	gcc 13.3	icpx 2025.0.4	clang 18.1.3	gcc 13.3	icpx 2025.0.4	clang 18.1.3
schubfach	11.72	11.26	11.26	11.96	12.07	11.97
schubfach_xjb	7.50	5.38	5.45	8.70	7.18	7.48
ryu	14.23	14.21	14.48	13.73	13.53	13.59
dragonbox				10.44	7.94	8.95
yy_fast				6.14	5.59	5.03
yy_full				7.66	4.58	4.94
teju_jagua	13.73	14.90	14.39	13.38	15.23	13.59
xjb	7.84	2.88	4.19 1.88(AVX2) 1.68(AVX512)	6.60	3.90	3.60 3.13(AVX2) 3.47(AVX512)

Table 3. Benchmark results on Apple M1 and MacOS 26.0.1. The unit is nanosecond(ns).

algorithm	float	double
	apple clang 17.0.0	apple clang 17.0.0
schubfach	11.11	12.62
schubfach_xjb	6.06	7.26
ryu	15.47	14.25
dragonbox		10.15
yy_fast		4.74
yy_full		6.54
teju_jagua	19.35	18.77
xjb	3.24	3.54

Table 1. All algorithms in the benchmark test.

algorithm	float	double	description
schubfach	schubfach32	schubfach64	Raffaello Giulietti, https://github.com/c4f7fcc09cb06515/Schubfach .
schubfach_xjb	schubfach32_xjb	schubfach64_xjb	It is improved by schubfach and has the same output result.
ryu	ryu32	ryu64	Ulf Adams, https://github.com/ulfjack/ryu .
dragonbox		dragonbox	Junekey Jeon, https://github.com/jk-jeon/dragonbox .
yy_fast		yy_fast	yy's algorithm, Some calculation steps have been modified. https://github.com/ibireme/c_numconv_benchmark/blob/master/vendor/yy_double/yy_double.c .
yy_full		yy_full	Same as yy_fast.
teju_jagua	teju32	teju64	Cassio Neri, https://github.com/cassioneri/teju_jagua .
xjb	xjb32	xjb64	this paper, https://github.com/xjb714/xjb .

* The blank Spaces in the table indicate no-code implementations.

The test results of AVX2 and AVX512 are for reference only. We will continue to optimize them in the future. From the benchmark results in Table 2 and Table 3, it can be seen that the performance of the algorithm in this paper is better than other algorithms in most cases.

510
511
512
513

5. Conclusions and future work

This paper proposes a new floating-point number to decimal conversion algorithm. The algorithm improves the calculation process of schubfach[2] algorithms, reduces the number of multiplication operations, and optimizes some calculation steps. The algorithm has been implemented in C/C++ language and passed exhaustive tests. The benchmark results show that the performance of the algorithm is better than most existing algorithms in most cases. Future work includes further optimization of the algorithm to improve performance, especially for parallel computing on x86-64 and arm64 architecture.

References

1. G. L. Steel Jr. and J. L. White. How to Print Floating-Point Numbers Accurately. In *Proceedings of the ACM SIGPLAN 1990 Conference on Programming Language Design and Implementation, PLDI 1990*. ACM, New York, NY, USA, 112-126. <https://doi.org/10.1145/93542.93559>
2. R. Giulietti. The Schubfach Way to Render Doubles. 2020. https://drive.google.com/file/d/1KLtG_LalbK9ETXI290zqCxvBW94dj058/view (Sep. 2020)
3. Ulf Adams. 2018. Ryū: Fast Float-to-String Conversion. In *Proceedings of 39th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI'18)*. ACM, New York, NY, USA, 13 pages. <https://doi.org/10.1145/3192366.3192369>
4. Ulf Adams. 2019. Ryū Revisited: Printf Floating Point Conversion. *Proc. ACM Program. Lang.* 3, OOPSLA, Article 169 (October 2019), 23 pages. <https://doi.org/10.1145/3360595>
5. Junekey Jeon. 2024. Dragonbox: A New Floating-Point Binary-to-Decimal Conversion Algorithm. <https://github.com/jk-jeon/dragonbox>

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.