

Article

Not peer-reviewed version

---

# xjb: Fast Float to String Algorithm

---

[Junbo Xiang](#) and [Tiejun Wang](#)\*

Posted Date: 1 April 2026

doi: 10.20944/preprints202511.1698.v2

Keywords: floating-point; printing; algorithm; performance



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a [Creative Commons CC BY 4.0 license](#), which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

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.

Article

# xjb: Fast Float to String Algorithm

Junbo Xiang  and Tiejun Wang \* 

Chengdu University of Information Technology

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

## Abstract

Efficiently and accurately converting floating-point numbers to decimal strings is a critical challenge in numerical computation and data exchange. While existing algorithms like Ryū, Dragonbox, and Schubfach satisfy the Steele-White (SW) principle for accuracy, they often suffer from performance bottlenecks due to branch prediction failures and high-precision multiplication overhead. This paper presents a novel floating-point to string conversion algorithm called "xjb", an optimized variant of the Schubfach algorithm designed to deliver superior performance for IEEE754 single-precision (binary32) and double-precision (binary64) floating-point numbers. By minimizing instruction dependencies, reducing multiplication operations, mitigating branch prediction penalties and by utilizing the simd instruction set, xjb achieves significant performance gains. The algorithm features concise core implementation. Extensive benchmarking across diverse platforms, including AMD R7-7840H and Apple M1, demonstrates that xjb outperforms state-of-the-art algorithms in most scenarios while maintaining full compliance with the SW principle.

**Keywords:** floating-point; printing; algorithm; performance

## 1. Introduction

Floating-point to decimal string conversion is a fundamental operation in computer systems with widespread applications across numerous domains. From scientific computing and financial systems to web services and database management, the ability to efficiently and accurately convert binary floating-point representations into human-readable decimal strings is essential. Despite its apparent simplicity, this conversion problem presents significant challenges in balancing the competing demands of correctness, performance, and output compactness.

### 1.1. Background and Motivation

In 1990, Steele and White [1,2] published the seminal paper *How to Print Floating-Point Numbers Accurately*, which established the foundational principles for optimal floating-point printing algorithms, now widely known as the Steele-White (SW) principle. The SW principle comprises four key requirements:

- **Information preservation:** The printed result must be parseable back to the original floating-point number without loss of precision.
- **Minimum length:** The output string should be as short as possible while maintaining information preservation.
- **Correct rounding:** When multiple representations satisfy the first two criteria, the algorithm must correctly round to the nearest value, with ties broken by selecting the even value.
- **Left-to-right generation:** The output digits should be generated sequentially from the most significant to the least significant digit.

These principles ensure that floating-point numbers have a unique, well-defined decimal representation that is both human-readable and machine-parseable. Algorithms satisfying the SW principle

guarantee that the conversion process is reversible and produces the shortest possible output, which is crucial for data exchange, serialization, and user interface display.

Over the past three decades, significant research efforts have been devoted to developing efficient algorithms that satisfy the SW principle. Early approaches, such as the **Dragon4** algorithm, provided correct results but suffered from performance limitations due to the use of arbitrary-precision arithmetic. Subsequent innovations introduced various optimizations:

- **dtoa.c** [3] is an improved version of Dragon4, written by Robert G. Burger and R. Kent Dybvig in 1996. The source code link is <https://www.netlib.org/fp/dtoa.c>.
- **Grisu3** [4] pioneered the use of precomputed powers of ten to avoid expensive arbitrary-precision operations, though it occasionally falls back to slower methods for certain inputs.
- **Errol** [5] improved upon Grisu3 by reducing the fallback rate through more precise error analysis.
- **Ryū** [6,7] achieved significant performance gains through careful instruction scheduling and lookup table optimizations, establishing a new baseline for high-performance conversion.
- **Schubfach** [8] introduced a novel approach based on the concept of "Schubfach" (pigeonhole principle), providing a compact and elegant solution with competitive performance.
- **Grisu-Exact** [9] eliminated the need for fallbacks entirely while maintaining high performance.
- **Dragonbox** [10] reduces the number of multiplications by sacrificing more branches.
- **yy\_double** [11] explored alternative computational strategies to minimize the multiplication cost, but there were still a few branches.
- **uscale** [12–14] was proposed by Russ Cox and is used to enhance the performance of the floating-point number printing algorithm in the Go programming language.
- **zmij** [15] is based on yy\_double and has undergone extensive code optimization.

Despite these advances, existing algorithms still face several performance challenges that limit their effectiveness in high-throughput scenarios:

1. **Branch prediction penalties:** Many algorithms rely heavily on conditional branches to handle different cases, leading to frequent branch mispredictions on modern pipelined processors.
2. **High-precision multiplication overhead:** The conversion process requires high-precision arithmetic operations, particularly multiplications involving large precomputed constants, which can be expensive on standard hardware.
3. **Instruction dependency chains:** Sequential dependencies between operations limit instruction-level parallelism and prevent efficient utilization of modern superscalar processors.
4. **Limited SIMD utilization:** Most existing algorithms do not exploit vector instruction sets (SIMD) that are now ubiquitous in contemporary processors.

These limitations motivate the need for a new approach that addresses these performance bottlenecks while maintaining full compliance with the SW principle.

## 1.2. Contributions

This paper presents **xjb**, a novel floating-point to string conversion algorithm that achieves superior performance through systematic optimization of the underlying computational structure. The xjb algorithm is derived from the Schubfach algorithm and incorporates insights from yy\_double and Dragonbox, but introduces several key innovations:

1. **Reduced instruction dependencies:** By carefully restructuring the computation, xjb minimizes data dependencies between operations, enabling better instruction-level parallelism and improved pipeline utilization.
2. **Minimized multiplication operations:** The algorithm reduces the number of expensive high-precision multiplications required during conversion, significantly decreasing computational overhead.
3. **Mitigated branch prediction penalties:** Through branchless programming techniques and careful case analysis, xjb minimizes conditional branches that could lead to prediction failures.

4. **SIMD instruction utilization:** The algorithm is designed to leverage SIMD instructions where applicable.
5. **Concise core implementation:** Despite its sophisticated optimizations, xjb maintains a compact and readable core implementation, facilitating adoption and maintenance.

The xjb algorithm supports both IEEE 754 single-precision (binary32) and double-precision (binary64) floating-point formats, which are the most widely used floating-point representations in modern computing. For simplicity, this paper uses *float* to refer to IEEE 754 binary32 and *double* to refer to IEEE 754 binary64.

### 1.3. Evaluation Overview

We conducted extensive benchmarking of xjb across diverse hardware platforms, including AMD R7-7840H and Apple M1 processors. Our evaluation demonstrates that xjb outperforms state-of-the-art algorithms in most scenarios while maintaining full compliance with the SW principle. The algorithm exhibits excellent portability and scalability, making it suitable for deployment across a wide range of systems from embedded devices to high-performance servers.

The implementation code for the xjb algorithm and all benchmark tools used in this paper are available at <https://github.com/xjb714/xjb>.

The remainder of this paper is organized as follows. Section 2 reviews the IEEE 754 floating-point representation and establishes the mathematical foundation for the conversion problem. Section 3 presents the core principles and derivation of the xjb algorithm and describes the implementation details and optimizations. Section 4 presents experimental results comparing xjb against existing algorithms.

### 1.4. Explanation of Special Symbols in This Article

Provide special explanations for the special symbols used in the formulas of this article.

**Table 1.** Explanation of Special Symbols in This Article.

symbol	brief explanation	example
%	Integer modulus operation	$2 = 8\%3$
//	Integer division operation	$1 = 5//3$
<< or >>	Left or right shift of binary values	$8 = 1<<3$
? :	Similar to the ternary operator in C syntax	$a = 1?a:b$

## 2. IEEE754 Floating-Point Number Representation

This section establishes the mathematical foundation for floating-point number representation and defines the notation used throughout this paper. We focus on the IEEE 754 standard, which is the most widely adopted floating-point arithmetic standard in modern computing systems.

### 2.1. Scope and Assumptions

For clarity of presentation, we make the following simplifying assumptions:

- We consider only positive floating-point numbers, as negative numbers differ only by a leading minus sign.
- We exclude special values (zero, NaN, and infinity) from our analysis, as these are handled separately in practice.

These assumptions are standard in the literature and do not affect the generality of our algorithm, as the excluded cases can be handled with straightforward special-case logic.

### 2.2. Binary Representation

The IEEE 754 standard [16,17] defines two primary floating-point formats relevant to this work:

**Double-precision (binary64):** A 64-bit format consisting of:

- 1 sign bit ( $s$ ): indicates positive ( $s = 0$ ) or negative ( $s = 1$ ).
- 11 exponent bits ( $e$ ): biased exponent in the range  $[0, 2047]$ .
- 52 fraction bits ( $f$ ): significand fraction in the range  $[0, 2^{52} - 1]$ .

**Single-precision (binary32):** A 32-bit format consisting of:

- 1 sign bit ( $s$ ): indicates positive ( $s = 0$ ) or negative ( $s = 1$ ).
- 8 exponent bits ( $e$ ): biased exponent in the range  $[0, 255]$ .
- 23 fraction bits ( $f$ ): significand fraction in the range  $[0, 2^{23} - 1]$ .

### 2.3. Classification of Floating-Point Numbers

We classify floating-point numbers into three categories based on their exponent and fraction fields:

1. **Subnormal numbers** ( $e = 0$  and  $f \neq 0$ ): These represent very small values close to zero, where the implicit leading bit of the significand is 0 instead of 1.
2. **Normal numbers** ( $e \neq 0$  and  $f \neq 0$ ): The standard case where the implicit leading bit of the significand is 1.
3. **Irregular numbers** ( $e \neq 0$  and  $f = 0$ ): Numbers with zero fraction field, representing powers of two.

We use the term **regular** to refer to both subnormal and normal numbers (i.e., all cases where  $f \neq 0$ ).

### 2.4. Value Representation

The real value  $v$  of a positive floating-point number can be expressed in the unified form  $v = c \cdot 2^q$ , where  $c$  is the integer significand and  $q$  is the exponent. The general formula covering all cases is:

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

For each category, the values decompose as follows:

**Subnormal numbers** ( $e = 0$ ):

$$\text{subnormal : } \begin{cases} \text{double : } v = f \cdot 2^{-1074} \\ \text{float : } v = f \cdot 2^{-149} \end{cases} \quad (2)$$

**Irregular numbers** ( $f = 0, e \neq 0$ ):

$$\text{irregular : } \begin{cases} \text{double : } v = 2^{52} \cdot 2^{e-1075} \\ \text{float : } v = 2^{23} \cdot 2^{e-150} \end{cases} \quad (3)$$

**Normal numbers** ( $e \neq 0, f \neq 0$ ):

$$\text{normal : } \begin{cases} \text{double : } v = (f + 2^{52}) \cdot 2^{e-1075} \\ \text{float : } v = (f + 2^{23}) \cdot 2^{e-150} \end{cases} \quad (4)$$

### 2.5. Rounding Intervals

A critical concept for accurate floating-point printing is the **rounding interval**  $R_v$ , which defines the range of real numbers that round to the given floating-point value  $v$  when parsed. The rounding interval is bounded by:

$$\begin{aligned}
v_l &= \begin{cases} \left(c - \frac{1}{2}\right) \cdot 2^q & \text{if } f \neq 0 \text{ or } e \leq 1 \\ \left(c - \frac{1}{4}\right) \cdot 2^q & \text{if } f = 0 \end{cases} \\
v_r &= \left(c + \frac{1}{2}\right) \cdot 2^q \\
R_v &= \begin{cases} [v_l, v_r] & \text{if } f \bmod 2 = 0 \text{ (even significand)} \\ (v_l, v_r) & \text{if } f \bmod 2 = 1 \text{ (odd significand)} \end{cases}
\end{aligned} \tag{5}$$

The rounding radius for regular floating-point numbers is  $2^{q-1}$ . The distinction between closed and open intervals at the boundaries depends on the parity of the significand, ensuring correct rounding according to the round-to-even rule specified in the IEEE 754 standard. Any decimal number within  $R_v$  will parse back to the original floating-point value  $v$ , which is essential for ensuring the information preservation property of the SW principle.

Table 2 summarizes the valid ranges for  $c$  and  $q$  across different categories.

**Table 2.** Valid ranges for significand  $c$  and exponent  $q$

Category	Float (binary32)	Double (binary64)
Subnormal	$1 \leq c \leq 2^{23} - 1, q = -149$	$1 \leq c \leq 2^{52} - 1, q = -1074$
Normal	$2^{23} + 1 \leq c \leq 2^{24} - 1$ $-148 \leq q \leq 104$	$2^{52} + 1 \leq c \leq 2^{53} - 1$ $-1073 \leq q \leq 971$
Irregular	$c = 2^{23}, -149 \leq q \leq 104$	$c = 2^{52}, -1074 \leq q \leq 971$

### 3. Principle of Algorithm

This section presents the algorithmic principles and mathematical foundation of the xjb floating-point to string conversion algorithm. We first introduce the overall architecture and design goals, followed by the mathematical formulation of the conversion problem.

#### 3.1. Design Overview

This paper focuses on converting float (single-precision) and double (double-precision) floating-point numbers to decimal strings. The conversion process consists of two stages:

1. **Float-to-Decimal Conversion:** Converting binary floating-point values to decimal significand-exponent pairs  $(d, k)$ .
2. **Decimal-to-String Conversion:** Formatting  $(d, k)$  into human-readable strings.

We focus primarily on the first stage, as it is the computationally intensive component that benefits most from our optimizations.

The xjb algorithm is designed with four key performance objectives:

- Minimize branch mispredictions through unlikely branches and branchless techniques
- Reduce expensive high-precision multiplication operations
- Optimize instruction dependencies to enable better parallelism
- Maintain core implementation simplicity and readability

#### 3.2. Mathematical Foundation

Before presenting the algorithm details, we establish the mathematical framework for the float-to-decimal conversion problem.

Recall from Section 2 that any floating-point value  $v$  can be expressed in the form  $v = c \cdot 2^q$ , where  $c$  is the integer significand and  $q$  is the exponent. Our goal is to find the optimal decimal representation  $opt = d \cdot 10^k$  that satisfies the SW principle.

As established in Section 2, regular floating-point numbers (which include all normal and sub-normal numbers with non-zero fraction fields) account for the vast majority of possible floating-point values. For the purposes of algorithm derivation, we focus primarily on regular numbers, as special cases can be handled with minimal additional logic.

The valid ranges for the significand  $c$  and exponent  $q$  for regular floating-point numbers are:

$$\begin{aligned} \text{float} : & \begin{cases} 1 \leq c \leq 2^{24} - 1, c \neq 2^{23}, & q = -149 \\ 2^{23} + 1 \leq c \leq 2^{24} - 1, & -148 \leq q \leq 104 \end{cases} \\ \text{double} : & \begin{cases} 1 \leq c \leq 2^{53} - 1, c \neq 2^{52}, & q = -1074 \\ 2^{52} + 1 \leq c \leq 2^{53} - 1, & -1073 \leq q \leq 971 \end{cases} \end{aligned} \quad (6)$$

For irregular floating-point numbers (powers of two), the ranges are:

$$\begin{aligned} \text{float} : & c = 2^{23}, \quad -149 \leq q \leq 104 \\ \text{double} : & c = 2^{52}, \quad -1074 \leq q \leq 971 \end{aligned} \quad (7)$$

For subnormal numbers:

$$\begin{aligned} \text{float} : & c \leq 2^{23} - 1, \quad q = -149 \\ \text{double} : & c \leq 2^{52} - 1, \quad q = -1074 \end{aligned} \quad (8)$$

The conversion problem can now be formally stated as: given a floating-point value  $v = c \cdot 2^q$ , find the optimal decimal representation  $opt = d \cdot 10^k$  such that:

$$\begin{aligned} v = c \cdot 2^q & \rightarrow opt = d \cdot 10^k \\ \text{subject to: } & opt \in R_v, \quad d \in \mathbb{Z}^+, \quad k \in \mathbb{Z} \end{aligned} \quad (9)$$

where  $R_v$  is the rounding interval of  $v$  as defined in Section 2.

For example, consider the IEEE 754 binary64 floating-point number representing 1.3. Its actual value is 1.3000000000000000444089209850062616169452667236328125, with hexadecimal representation 3ff4cccccccccd. The optimal decimal representation  $opt$  satisfying the SW principle is simply 1.3. Similarly, for the binary32 representation of 1.3, with actual value 1.2999999523162841796875 and hexadecimal representation 3fa66666, the optimal representation is also 1.3.

### 3.3. Review of the Schubfach Algorithm and Our Derivation

This section reviews the Schubfach algorithm and presents our derivation of an optimized variant. We begin by establishing the mathematical foundation for determining the optimal decimal representation.

#### 3.3.1. Candidate Values for the Significand $d$

The Schubfach algorithm identifies four candidate values for the decimal significand  $d$ :

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

The exponent  $k$  is computed as:

$$k = \begin{cases} \lfloor q \cdot \lg 2 \rfloor & \text{if } v \in \text{regular} \\ \lfloor q \cdot \lg 2 - \lg(4/3) \rfloor & \text{otherwise} \end{cases} \quad (11)$$

For efficient computation on modern processors, Equation (11) can be implemented using integer arithmetic:

$$\begin{aligned} \text{double : } k &= (q \cdot 315653 - (v \in \text{regular ? 0 : 131072})) \gg 20 \\ \text{float : } k &= (q \cdot 1233 - (v \in \text{regular ? 0 : 512})) \gg 12 \end{aligned} \quad (12)$$

### 3.3.2. Decomposition into Integer and Fractional Parts

Let  $m = \lfloor v \cdot 10^{-k-1} \rfloor$  denote the integer part and  $n = v \cdot 10^{-k-1} - m$  the fractional part, where  $0 \leq n < 1$ . Substituting into Equation (10):

$$d \in \{10m, \lfloor 10(m+n) \rfloor, \lfloor 10(m+n) \rfloor + 1, 10m + 10\} \quad (13)$$

Since  $\lfloor 10(m+n) \rfloor = 10m + \lfloor 10n \rfloor$ , the candidates simplify to:

$$d \in \{10m, 10m + \lfloor 10n \rfloor, 10m + \lfloor 10n \rfloor + 1, 10m + 10\} \quad (14)$$

Based on the Schubfach algorithm and the SW principle, if the value  $10m$  or  $10m + 10$  falls within the range  $R_v$ , it is selected as the optimal solution  $d$ . In cases where neither  $10m$  nor  $10m + 10$  lies within  $R_v$ , the optimal solution  $d$  is determined as either  $10m + \lfloor 10n \rfloor$  or  $10m + \lfloor 10n \rfloor + 1$  in accordance with the rules of correct rounding. We decompose  $d$  as  $d = ten + one$ , where  $ten = 10m$  and  $one \in \{0, \lfloor 10n \rfloor, \lfloor 10n \rfloor + 1, 10\}$ . The problem thus reduces to determining the appropriate value of  $one$ .

$$one = \begin{cases} 0; & \text{if } 10m \in R_v \\ 10; & \text{else if } 10m + 10 \in R_v \\ \lfloor 10n \rfloor \text{ or } \lfloor 10n \rfloor + 1; & \text{else apply correct rounding} \end{cases} \quad (15)$$

### 3.3.3. Selection Criteria for *One*

The selection of  $one$  depends on the relationship between the rounding interval and the candidate values. Recall that the rounding interval for a regular floating-point number  $v = c \cdot 2^q$  is  $R_v = [v - 2^{q-1}, v + 2^{q-1}]$ , where  $2^{q-1}$  represents the half-unit in the last place (ulp) of  $v$ .

- **Case  $one = 0$  (i.e.,  $d = 10m$ ):** This case applies when  $10m \cdot 10^k$  is the closest representable decimal to  $v$  within the rounding interval. The condition is derived as follows:  
The lower bound of the rounding interval must be less than  $10m \cdot 10^k$ :

$$\begin{aligned} c \cdot 2^q - 10m \cdot 10^k &< 2^{q-1} \\ c \cdot 2^q - \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor \cdot 10^{k+1} &< 2^{q-1} \\ c \cdot 2^q \cdot 10^{-k-1} - \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor &< 2^{-1} \cdot 2^q \cdot 10^{-k-1} \\ n &< 2^{q-1} \cdot 10^{-k-1} \end{aligned} \quad (16)$$

When equality holds ( $n = 2^{q-1} \cdot 10^{-k-1}$ ), we apply the round-to-even rule, requiring  $c$  to be even:

$$one = 0 \quad \text{if } 2^{q-1} \cdot 10^{-k-1} > n \quad \text{or} \quad (2^{q-1} \cdot 10^{-k-1} = n \text{ and } c \bmod 2 = 0) \quad (17)$$

- **Case  $one = 10$  (i.e.,  $d = 10m + 10$ ):** This case applies when  $(10m + 10) \cdot 10^k$  is the closest representable decimal to  $v$ . The condition is derived similarly:  
The upper bound of the rounding interval must be greater than  $(10m + 10) \cdot 10^k$ :

$$\begin{aligned} (10m + 10) \cdot 10^k - c \cdot 2^q &< 2^{q-1} \\ \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor \cdot 10^{k+1} + 10^{k+1} - c \cdot 2^q &< 2^{q-1} \\ 1 - (c \cdot 2^q \cdot 10^{-k-1} - \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor) &< 2^{-1} \cdot 2^q \cdot 10^{-k-1} \\ 1 - n &< 2^{q-1} \cdot 10^{-k-1} \end{aligned} \quad (18)$$

When equality holds ( $1 - n = 2^{q-1} \cdot 10^{-k-1}$ ), we again apply round-to-even:

$$one = 10 \quad \text{if} \quad 2^{q-1} \cdot 10^{-k-1} > 1 - n \quad \text{or} \quad \left( 2^{q-1} \cdot 10^{-k-1} = 1 - n \text{ and } c \bmod 2 = 0 \right) \quad (19)$$

- **Case**  $one \in \{\lfloor 10n \rfloor, \lfloor 10n \rfloor + 1\}$ : When neither boundary condition applies, the optimal value lies between  $10m + \lfloor 10n \rfloor$  and  $10m + \lfloor 10n \rfloor + 1$ . We determine  $one$  by rounding  $10n$  to the nearest integer:
  - If the fractional part  $\{10n\} < 0.5$ :  $one = \lfloor 10n \rfloor$
  - If the fractional part  $\{10n\} > 0.5$ :  $one = \lfloor 10n \rfloor + 1$
  - If the fractional part  $\{10n\} = 0.5$ : apply round-to-even

For irregular floating-point numbers (powers of two), additional verification is required to ensure the selected value lies within the rounding interval  $R_v$ , as the interval boundaries differ for these special cases.

### 3.3.4. Algorithm Overview

Algorithm 1 summarizes our optimized variant of the Schubfach algorithm (xjb32 for float, xjb64 for double). Given inputs  $c$  and  $q$ , the algorithm returns  $d$  and  $k$  such that  $d \cdot 10^k$  satisfies the SW principle. The computation of  $k$  follows Equation (12); the remainder of this section focuses on efficient computation of  $d$ .

The subsequent sections detail:

1. Lookup table precomputation
2. Efficient computation of  $m$
3. Fast boundary condition testing for  $one \in \{0, 10\}$
4. Efficient computation of  $\lfloor 10n \rfloor$  and rounding
5. Handling of irregular floating-point numbers
6. implementation of pseudocode

The following will discuss the above content in detail from Section 3.4 to Section 3.9.

### 3.4. Lookup Table Precomputation

The algorithm in this paper employs a lookup table to store pre-computed values of  $10^{-k-1}$  for different ranges of  $q$ :  $[-149, 104]$  for float and  $[-1074, 971]$  for double. These lookup tables use extended precision: 64-bit for float and 128-bit for double. The reference implementation is available in [gen.py](#).

#### 3.4.1. Fundamental Calculation

Let  $B$  denote the bit length of each entry in the lookup table, with  $B = 64$  for float and  $B = 128$  for double. For any integer  $e_{10}$  (representing a power of 10), we aim to represent  $10^{e_{10}}$  in the form  $f \cdot 2^{\lfloor e_2 \rfloor}$ , where  $1 \leq f < 2$  and  $e_2$  is a real number. This gives:

$$f \cdot 2^{\lfloor e_2 \rfloor} = 2^{e_2} = 10^{e_{10}} \quad (20)$$

Taking the logarithm base 2 of both sides, we get:

$$\lfloor e_2 \rfloor = \lfloor e_{10} \cdot \log_2(10) \rfloor \quad (21)$$

Solving for  $f$  gives:

$$f = \frac{10^{e_{10}}}{2^{\lfloor e_{10} \cdot \log_2(10) \rfloor}} \quad (22)$$

The lookup table entries are computed using upward rounding:

---

**Algorithm 1** The xjb Algorithm for Float-to-Decimal Conversion
 

---

**Require:** Floating-point components  $c$  (significand) and  $q$  (exponent)

**Ensure:** Decimal representation  $d \cdot 10^k$  satisfying the SW principle

```

1:  $c \cdot 2^q \leftarrow v$ 
2: if  $v \in \text{regular}$  then
3:    $k \leftarrow \lfloor q \cdot \lg 2 \rfloor$ 
4: else
5:    $k \leftarrow \lfloor q \cdot \lg 2 - \lg(4/3) \rfloor$ 
6: end if
7:  $m \leftarrow \lfloor v \cdot 10^{-k-1} \rfloor$ 
8:  $n \leftarrow v \cdot 10^{-k-1} - m$ 
9:  $ten \leftarrow 10m$ 
10:  $\delta \leftarrow 10n - \lfloor 10n \rfloor$  {fractional part of  $10n$ }
11: if  $\delta = 0.5$  then
12:   if  $\lfloor 10n \rfloor \bmod 2 = 0$  then
13:      $one \leftarrow \lfloor 10n \rfloor$  {round to even}
14:   else
15:      $one \leftarrow \lfloor 10n \rfloor + 1$ 
16:   end if
17: else if  $\delta < 0.5$  then
18:    $one \leftarrow \lfloor 10n \rfloor$ 
19: else
20:    $one \leftarrow \lfloor 10n \rfloor + 1$ 
21: end if
22: if  $v \in \text{irregular}$  then
23:   if  $\delta > 2^{q-2} \cdot 10^{-k}$  then
24:      $one \leftarrow \lfloor 10n \rfloor + 1$ 
25:   end if
26:   if  $2^{q-2} \cdot 10^{-k-1} \geq n$  then
27:      $one \leftarrow 0$ 
28:   end if
29: else
30:   if  $2^{q-1} \cdot 10^{-k-1} > n$  OR  $(2^{q-1} \cdot 10^{-k-1} = n$  and  $c \bmod 2 = 0)$  then
31:      $one \leftarrow 0$ 
32:   end if
33:   if  $2^{q-1} \cdot 10^{-k-1} > 1 - n$  OR  $(2^{q-1} \cdot 10^{-k-1} = 1 - n$  and  $c \bmod 2 = 0)$  then
34:      $one \leftarrow 10$ 
35:   end if
36: end if
37:  $d \leftarrow ten + one$ 
38:
39: return  $d, k$ 

```

---

$$\begin{aligned}
lookup[e_{10}] &= \lceil f \cdot 2^{B-1} \rceil \\
&= \lceil \frac{10^{e_{10}}}{2^{\lfloor e_{10} \cdot \log_2(10) \rfloor}} \cdot 2^{B-1} \rceil \\
&= \lceil 10^{e_{10}} \cdot 2^{B-1-\lfloor e_{10} \cdot \log_2(10) \rfloor} \rceil
\end{aligned} \tag{23}$$

Notably,  $f \cdot 2^{B-1}$  becomes an integer for certain ranges of  $e_{10}$ :  $0 \leq e_{10} \leq 27$  for float and  $0 \leq e_{10} \leq 55$  for double.

### 3.4.2. Detailed Calculation Process

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 (6), 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:

$$\begin{aligned}
-32 &\leq e_{10} \leq 44 \\
e_2 &= |\lfloor e_{10} \cdot \log_2(10) \rfloor - 63| \\
pow10t &= \begin{cases} 2^{e_2} / 10^{|e_{10}|}; & \text{if } e_{10} < 0 \\ 10^{|e_{10}|} / 2^{e_2}; & \text{if } e_{10} \geq 20 \\ 10^{|e_{10}|} \cdot 2^{e_2}; & \text{if } 1 \leq e_{10} \leq 19 \end{cases} \\
f_{1,e_{10}} &= pow10 = pow10t + (e_{10} \geq 0 \&\& e_{10} \leq 27 ? 0 : 1)
\end{aligned} \tag{24}$$

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:

$$\begin{aligned}
r_{1,e_{10}} &= \frac{f_{1,e_{10}} \cdot 2^{\lfloor e_{10} \cdot \log_2(10) \rfloor - 63}}{10^{e_{10}}} \\
&\in \begin{cases} 1; & \text{if } 0 \leq e_{10} \leq 27 \\ (1, 1 + 2^{-63}); & \text{if } e_{10} < 0 \text{ or } e_{10} > 27 \end{cases}
\end{aligned} \tag{25}$$

- Double

The range of  $-k - 1$  is calculated to be  $[-293, 323]$  through the  $q$  value range in Equation (6), 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:

$$\begin{aligned}
-293 &\leq e_{10} \leq 323 \\
e_2 &= |\lfloor e_{10} \cdot \log_2(10) \rfloor - 127| \\
pow10t &= \begin{cases} 2^{e_2} / 10^{|e_{10}|}; & \text{if } e_{10} < 0 \\ 10^{|e_{10}|} / 2^{e_2}; & \text{if } e_{10} \geq 39 \\ 10^{|e_{10}|} \cdot 2^{e_2}; & \text{if } 1 \leq e_{10} \leq 38 \end{cases} \\
f_{1,e_{10}} &= pow10 = pow10t + (e_{10} \geq 0 \&\& e_{10} \leq 55 ? 0 : 1)
\end{aligned} \tag{26}$$

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:

$$\begin{aligned}
r_{1,e_{10}} &= \frac{f_{1,e_{10}} \cdot 2^{\lfloor e_{10} \cdot \log_2(10) \rfloor - 127}}{10^{e_{10}}} \\
&\in \begin{cases} 1; & \text{if } 0 \leq e_{10} \leq 55 \\ (1, 1 + 2^{-127}); & \text{if } e_{10} < 0 \text{ or } e_{10} > 55 \end{cases}
\end{aligned} \tag{27}$$

Let  $r_1$  denote the error for float lookup table entries,  $r_2$  for double entries, and  $r$  for both. In Algorithm (1), we retrieve  $10^{-k-1}$  from the lookup table. The lookup table provides exact values when  $q$  falls within specific ranges:

$$\begin{aligned} \text{float} : 0 \leq -k - 1 \leq 27 &\Rightarrow -93 \leq q \leq -1 \\ \text{double} : 0 \leq -k - 1 \leq 55 &\Rightarrow -186 \leq q \leq -1 \end{aligned} \quad (28)$$

For  $q$  outside these ranges, the lookup table entries have bounded relative errors:

$$\begin{aligned} \text{float} : 0 < r_1 - 1 < 2^{-63} \\ \text{double} : 0 < r_2 - 1 < 2^{-127} \end{aligned} \quad (29)$$

### 3.4.3. Storage Requirements

The float lookup table requires 616 bytes of storage, calculated as  $(44 - (-32) + 1) \times 8$  bytes (77 entries  $\times$  8 bytes each). The double lookup table requires 9872 bytes, calculated as  $(323 - (-293) + 1) \times 16$  bytes (617 entries  $\times$  16 bytes each).

### 3.4.4. Implementation Notes

The lookup table pre-computation uses efficient integer arithmetic to avoid precision loss during calculations. The conditional logic in Equations (24) and (26) optimizes the computation based on the sign and magnitude of  $e_{10}$ , ensuring efficient generation of accurate lookup table entries.

## 3.5. Efficient Computation of $m$

This section presents an efficient method for calculating  $m$  in Algorithm (1), which is defined as  $m = \lfloor v \cdot 10^{-k-1} \rfloor$ .

### 3.5.1. Relevant Theorems

We start with some theorems from the Dragonbox algorithm paper:

Let  $n$ ,  $P$ , and  $Q$  be positive integers where  $P$  and  $Q$  are coprime,  $P < Q$ ,  $1 \leq n \leq n_{max}$ , and  $Q > n_{max}$ . Let  $P^*/Q^*$  be the best rational approximation of  $P/Q$  with  $Q^* \leq n_{max}$  and  $P_*/Q_*$  be the best rational approximation with  $Q_* \leq n_{max}$ .

If  $n \cdot P$  does not divide  $Q$  evenly:

$$\lfloor n \cdot \frac{P}{Q} \rfloor + 1 = \lceil n \cdot \frac{P}{Q} \rceil \quad (30)$$

Suppose the following holds:

$$\lfloor n \cdot \frac{P}{Q} \rfloor = \lfloor n \cdot \zeta \rfloor \quad (31)$$

Then:

$$\frac{P_*}{Q_*} = \max_{1 \leq n \leq n_{max}} \frac{\lfloor n \cdot \frac{P}{Q} \rfloor}{n} \leq \zeta < \min_{1 \leq n \leq n_{max}} \frac{\lfloor n \cdot \frac{P}{Q} \rfloor + 1}{n} = \min_{1 \leq n \leq n_{max}} \frac{\lceil n \cdot \frac{P}{Q} \rceil}{n} = \frac{P^*}{Q^*} \quad (32)$$

Thus, the range of  $\zeta$  is:

$$\frac{P_*}{Q_*} \leq \zeta < \frac{P^*}{Q^*} \quad (33)$$

The range of the fractional part of  $n \cdot \frac{P}{Q}$  is:

$$\left[ \frac{(Q_*P) \% Q}{Q}, \frac{(Q^*P) \% Q}{Q} \right] \quad (34)$$

Here, the fractional part is minimized when  $n = Q_*$  and maximized when  $n = Q^*$ .

### 3.5.2. Best Rational Approximation Function

We define the best rational approximation function (implemented in `test1.py`):

$$\left(\frac{P_*}{Q_*}, \frac{P^*}{Q^*}\right) = (DN, UP) = f(C, P, Q) \quad (35)$$

This function  $f(C, P, Q)$  calculates the best rational approximation with denominator not exceeding  $C$  using the Farey sequence median theorem.  $DN$  and  $UP$  are adjacent terms in the  $C$ -th order Farey sequence  $F_C$ . Therefore, we can calculate Equation (34) very quickly.

### 3.5.3. Key Proof

In Algorithm (1), we need to compute  $m = \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor$ . We aim to prove:

$$m = \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor c \cdot 2^q \cdot r \cdot 10^{-k-1} \rfloor \quad (36)$$

Where  $r$  is the lookup table error defined in Equations (25) and (27). When condition (28) is met,  $r = 1$ , and the equation holds trivially. For  $r \neq 1$ :

$$\begin{aligned} \text{float} : 1 < r < 1 + 2^{-63} \\ \text{double} : 1 < r < 1 + 2^{-127} \end{aligned} \quad (37)$$

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

$$2^q \cdot 10^{-k-1} = 10^{-1} \cdot \left(2^q \cdot 10^{-\lfloor q \cdot \lg(2) \rfloor}\right) = 10^{-1} \cdot \left(10^{q \cdot \lg(2) - \lfloor q \cdot \lg(2) \rfloor}\right) \quad (38)$$

When  $q$  is not 0, Equation (38) exists:

$$\begin{aligned} q \cdot \lg(2) &\neq \lfloor q \cdot \lg(2) \rfloor \\ 0 < q \cdot \lg(2) - \lfloor q \cdot \lg(2) \rfloor &< 1 \end{aligned} \quad (39)$$

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

$$10^{-1} \leq 2^q \cdot 10^{-k-1} < 1 \quad (40)$$

Because there is:

$$c \cdot 2^q \cdot 10^{-k-1} = c \cdot \frac{2^{q-k-1}}{5^{k+1}} \in [0.1c, c) \quad (41)$$

Therefore:

$$c \cdot 2^q \cdot 10^{-k-1} = \begin{cases} \frac{c \cdot 2^{q-k-1}}{5^{k+1}}; q \geq 1 \\ \frac{c}{2^{1+k-q} \cdot 5^{k+1}} = \frac{c}{10}; q = 0 \\ \frac{c \cdot 5^{-k-1}}{2^{1+k-q}}; q < 0 \end{cases} \quad (42)$$

Suppose:

$$c \cdot 2^q \cdot 10^{-k-1} = c \cdot \frac{x}{y} < c \quad (43)$$

Then there are:

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

### 3.5.4. Bit Width Calculation

Define maximum values for  $c$ :

$$\begin{aligned} \text{float} : c &\leq c_{max} = C_1 = 2^{24} - 1 \\ \text{double} : c &\leq c_{max} = C_2 = 2^{53} - 1 \end{aligned} \quad (45)$$

Let  $C$  denote either  $C_1$  or  $C_2$  depending on the precision.

For  $y > C$ , compute  $P^*$  and  $Q^*$  for each  $q$  using  $f(C, x, y)$  and find the minimum  $BIT$  such that:

$$\frac{x}{y}(1 + 2^{-BIT}) < \frac{P^*}{Q^*} \quad (46)$$

For  $y \leq C$ :

$$c \cdot \frac{x}{y} \left(1 + \frac{1}{Cy}\right) = \frac{cx + \frac{c}{C} \cdot \frac{x}{y}}{y} < \frac{cx + 1}{y} \quad (47)$$

Thus:

$$\lfloor c \cdot \frac{x}{y} \rfloor = \lfloor c \cdot \frac{x}{y} \left(1 + \frac{1}{Cy}\right) \rfloor \quad (48)$$

Similarly, find the minimum  $BIT$  such that:

$$\frac{x}{y}(1 + 2^{-BIT}) < \frac{x}{y} \left(1 + \frac{1}{Cy}\right) \quad (49)$$

### 3.5.5. Results

The maximum of the minimum  $BIT$  values for all  $q$  (calculated in [test1.py](#) in about 1-2 seconds) is:

$$\begin{aligned} \text{float} : BIT_{max} &= 52 \\ \text{double} : BIT_{max} &= 113 \end{aligned} \quad (50)$$

Thus:

$$\begin{aligned} \text{float} : \lfloor c \cdot \frac{x}{y} \rfloor &= \lfloor c \cdot \frac{x}{y} \cdot (1 + 2^{-52}) \rfloor = \lfloor c \cdot \frac{x}{y} \cdot r_1 \rfloor \\ \text{double} : \lfloor c \cdot \frac{x}{y} \rfloor &= \lfloor c \cdot \frac{x}{y} \cdot (1 + 2^{-113}) \rfloor = \lfloor c \cdot \frac{x}{y} \cdot r_2 \rfloor \end{aligned} \quad (51)$$

This result confirms that  $m$  can be calculated efficiently using the lookup table values, even with their inherent errors. Once  $m$  is determined,  $ten = 10m$  can be computed quickly.

### 3.6. Fast Boundary Condition Testing for $one = 0$ and $one = 10$

In Algorithm 1, the conditions for determining  $one = 0$  and  $one = 10$  appear on lines 31 and 34, respectively. This section introduces an optimized method to quickly test these boundary conditions using equivalent mathematical formulations.

#### 3.6.1. Equivalent Conditions for Boundary Testing

We start by deriving equivalent mathematical conditions for testing  $one = 0$  and  $one = 10$ .

##### Case 1: Testing $one = 0$

When  $2^{-1} \cdot 2^q \cdot 10^{-k-1} = n$ , this is equivalent to:

$$\begin{aligned} c \cdot 2^q \cdot 10^{-k-1} - \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor &= 2^{-1} \cdot 2^q \cdot 10^{-k-1} \\ (2c - 1) \cdot 2^{q-1} \cdot 10^{-k-1} &= \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor \end{aligned} \quad (52)$$

Case 2: Testing  $one = 10$

When  $2^{-1} \cdot 2^q \cdot 10^{-k-1} = 1 - n$ , this is equivalent to:

$$\begin{aligned} \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor - c \cdot 2^q \cdot 10^{-k-1} + 1 &= 2^{-1} \cdot 2^q \cdot 10^{-k-1} \\ (2c + 1) \cdot 2^{q-1} \cdot 10^{-k-1} &= \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor + 1 \end{aligned} \quad (53)$$

### 3.6.2. Integer Testing Analysis

To further analyze these conditions, we start with the range of  $2^{q-1} \cdot 10^{-k-1}$  from Equation (40):

$$2^{q-1} \cdot 10^{-k-1} \in [0.05, 0.5) \quad (54)$$

Analysis for  $one = 0$

For the  $one = 0$  case, we can derive:

$$\begin{aligned} \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor - 1 &< c \cdot 2^q \cdot 10^{-k-1} - 0.5 \\ &< (2c - 1) \cdot 2^{q-1} \cdot 10^{-k-1} \\ &\leq c \cdot 2^q \cdot 10^{-k-1} - 0.05 < \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor + 1 \end{aligned} \quad (55)$$

This implies that when  $(2c - 1) \cdot 2^{q-1} \cdot 10^{-k-1}$  is an integer, it must equal  $\lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor$ .

Analysis for  $one = 10$

Similarly, for the  $one = 10$  case:

$$\begin{aligned} \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor &< c \cdot 2^q \cdot 10^{-k-1} + 0.05 \\ &\leq (2c + 1) \cdot 2^{q-1} \cdot 10^{-k-1} \\ &< c \cdot 2^q \cdot 10^{-k-1} + 0.5 < \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor + 2 \end{aligned} \quad (56)$$

This implies that when  $(2c + 1) \cdot 2^{q-1} \cdot 10^{-k-1}$  is an integer, it must equal  $\lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor + 1$ .

### 3.6.3. Key Insight: Integer Divisibility Test

The key insight is that testing for  $one = 0$  or  $one = 10$  is equivalent to checking whether  $(2c \pm 1) \cdot 2^{q-1} \cdot 10^{-k-1}$  is an integer. This can be rewritten as:

$$(2c \pm 1) \cdot 2^{q-1} \cdot 10^{-k-1} = (2c \pm 1) \cdot 2^{q-k-2} \cdot 5^{-k-1} \quad (57)$$

We analyze different ranges of  $q$  to simplify this condition:

- **Case  $q \geq 2$ :** From  $q \geq 2$ , we get  $k \geq 0$ . The expression simplifies to checking if  $(2c \pm 1) \cdot 2^{q-k-2}$  is divisible by  $5^{k+1}$ . Since 2 and 5 are coprime, this reduces to checking if  $(2c \pm 1)$  is divisible by  $5^{k+1}$ :

$$(2c \pm 1) \pmod{5^{k+1}} = 0 \quad (58)$$

Let  $t$  be a positive integer such that  $2c \pm 1 = t \cdot 5^{k+1}$ . Since  $2c \pm 1$  is odd,  $t$  must also be odd. Considering the ranges of  $c$  for float and double:

$$\begin{aligned} \text{float : } 2c - 1 &\in [2^{24} + 1, 2^{25} - 3]; \quad 2c + 1 \in [2^{24} + 3, 2^{25} - 1]; \\ \text{double : } 2c - 1 &\in [2^{53} + 1, 2^{54} - 3]; \quad 2c + 1 \in [2^{53} + 3, 2^{54} - 1]; \end{aligned} \quad (59)$$

This gives us the range for  $t$ :

$$\begin{aligned} \text{float} : \frac{2^{24} + 1}{5^{k+1}} &\leq t \leq \frac{2^{25} - 1}{5^{k+1}}; \\ \text{double} : \frac{2^{53} + 1}{5^{k+1}} &\leq t \leq \frac{2^{54} - 1}{5^{k+1}}; \end{aligned} \quad (60)$$

The maximum values of  $k$  where  $t$  can be at least one odd integer are:

$$\begin{aligned} \text{float} : k_{\max} = 9 \Rightarrow q_{\max} = 33, t = 3 \\ \text{double} : k_{\max} = 22 \Rightarrow q_{\max} = 76, t = 1 \end{aligned} \quad (61)$$

- **Case  $1 \geq q \geq 0$ :** The denominator  $2^{2+k-q} \cdot 5^{k+1}$  is even, while the numerator  $(2c \pm 1)$  is odd, so no solution exists.
- **Case  $q < 0$ :** The denominator  $2^{2+k-q}$  is even, while the numerator  $(2c \pm 1) \cdot 5^{-k-1}$  is odd, so no solution exists.

### 3.6.4. Summary of Boundary Conditions

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

$$\begin{aligned} \text{float} : 2 \leq q \leq 33 \ \&\& \ (2c \pm 1) \pmod{5^{k+1}} = 0; \\ \text{double} : 2 \leq q \leq 76 \ \&\& \ (2c \pm 1) \pmod{5^{k+1}} = 0; \end{aligned} \quad (62)$$

The range of  $-k - 1$  is:

$$\begin{aligned} \text{float} : -10 \leq -k - 1 \leq -1 \\ \text{double} : -23 \leq -k - 1 \leq -1 \end{aligned} \quad (63)$$

### 3.6.5. Efficient Implementation

We can further simplify the testing conditions using bitwise operations. For  $one = 0$ :

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

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

$$\begin{aligned} \text{float} : \begin{cases} 2^{35} \cdot 2^q \cdot 10^{-k-1} = 2^{36} - 2^{36} \cdot n \Rightarrow \\ \lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor 2^{36} - 2^{36} \cdot n \rfloor = 2^{36} - 1 - \lfloor 2^{36} \cdot n \rfloor \end{cases} \\ \text{double} : \begin{cases} 2^{63} \cdot 2^q \cdot 10^{-k-1} = 2^{64} - 2^{64} \cdot n \Rightarrow \\ \lfloor 2^{63} \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor 2^{64} - 2^{64} \cdot n \rfloor = 2^{64} - 1 - \lfloor 2^{64} \cdot n \rfloor \end{cases} \end{aligned} \quad (65)$$

The discussion on whether  $\lfloor 2^{36} - 2^{36} \cdot n \rfloor = 2^{36} - 1 - \lfloor 2^{36} \cdot n \rfloor$  in Equation (65) holds true, that is, whether  $2^{36} \cdot n$  in Equation (65) is an integer, or equivalent to discussing whether the following values are integers when Equation (62) holds true (the same applies to double) :

$$\begin{aligned} \text{float} : 2^{36} \cdot (m + n) &= c \cdot 2^{q+36} \cdot 10^{-k-1} = c \cdot 2^{q-k+35} \cdot 5^{-k-1} = c \cdot \frac{2^{q-k+35}}{5^{k+1}} \\ \text{double} : 2^{64} \cdot (m + n) &= c \cdot 2^{q+64} \cdot 10^{-k-1} = c \cdot 2^{q-k+63} \cdot 5^{-k-1} = c \cdot \frac{2^{q-k+63}}{5^{k+1}} \end{aligned} \quad (66)$$

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

$$c = t \cdot 5^{k+1}; t \geq 1 \quad (67)$$

Therefore, when Equation (67) was established, there were:

$$2c \pm 1 = 2 \cdot t \cdot 5^{k+1} \pm 1 \quad (68)$$

Expression (68) cannot divide  $5^{k+1}$  evenly, which contradicts Equation (62), so  $c$  cannot divide  $5^{k+1}$  evenly. Therefore, for float,  $c \cdot 2^{q+36} \cdot 10^{-k-1}$  and  $2^{36} \cdot n$  are not integers; For double,  $c \cdot 2^{64+q} \cdot 10^{-k-1}$  and  $2^{64} \cdot n$  are not integers, that is:

$$\begin{aligned} \text{float} : [2^{36} - 2^{36} \cdot n] &= 2^{36} + [-2^{36} \cdot n] = 2^{36} - 1 - [2^{36} \cdot n] \\ \text{double} : [2^{64} - 2^{64} \cdot n] &= 2^{64} + [-2^{64} \cdot n] = 2^{64} - 1 - [2^{64} \cdot n] \end{aligned} \quad (69)$$

Therefore, the conclusion (65) is correct. Discuss the necessary and sufficient conditions for whether  $[2^{35} \cdot 2^q \cdot 10^{-k-1}] = [2^{36} \cdot n]$  is  $2^{-1} \cdot 2^q \cdot 10^{-k-1} = n$ . The same applies to double, expressed as:

$$\begin{aligned} \text{float} : 2^{-1} \cdot 2^q \cdot 10^{-k-1} = n &\Leftrightarrow [2^{35} \cdot 2^q \cdot 10^{-k-1}] = [2^{36} \cdot n] \\ \text{double} : 2^{-1} \cdot 2^q \cdot 10^{-k-1} = n &\Leftrightarrow [2^{63} \cdot 2^q \cdot 10^{-k-1}] = [2^{64} \cdot n] \end{aligned} \quad (70)$$

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

$$\begin{aligned} \text{float} : 2^{-1} \cdot 2^q \cdot 10^{-k-1} = 1 - n &\Leftrightarrow [2^{35} \cdot 2^q \cdot 10^{-k-1}] = [2^{36} - 2^{36} \cdot n] \\ \text{double} : 2^{-1} \cdot 2^q \cdot 10^{-k-1} = 1 - n &\Leftrightarrow [2^{63} \cdot 2^q \cdot 10^{-k-1}] = [2^{64} - 2^{64} \cdot n] \end{aligned} \quad (71)$$

The sufficient conditions of Equations (70) and (71) are obviously established. Introduce the proof that Equation (70) 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  $[2^{35} \cdot 2^q \cdot 10^{-k-1}] = [2^{36} \cdot n]$  holds, or equivalent to  $[2^{35} \cdot 2^q \cdot 10^{-k-1}] \neq [2^{36} \cdot n]$  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  $[2^{35} \cdot 2^q \cdot 10^{-k-1}] = [2^{36} \cdot n]$  holds when  $2^{-1} \cdot 2^q \cdot 10^{-k-1} \neq n$ . Then there is:

$$\begin{aligned} [2^{35} \cdot 2^q \cdot 10^{-k-1}] &= [2^{36} \cdot n] \\ \Rightarrow 0 &< |2^{35} \cdot 2^q \cdot 10^{-k-1} - 2^{36} \cdot n| < 1 \\ \Rightarrow 0 &< |(2c - 1) \cdot 2^{q-1} \cdot 10^{-k-1} - m| < 2^{-36} \end{aligned} \quad (72)$$

As is known from Equation (55), there is:

$$m - 1 < (2c - 1) \cdot 2^{q-1} \cdot 10^{-k-1} < m + 1 \quad (73)$$

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

$$|(2c - 1) \cdot 2^{q-1} \cdot 10^{-k-1} - m| = \begin{cases} n^-; & \text{if } (2c - 1) \cdot 2^{q-1} \cdot 10^{-k-1} > m \\ 1 - n^-; & \text{if } (2c - 1) \cdot 2^{q-1} \cdot 10^{-k-1} < m \end{cases} \quad (74)$$

Substitute expression (74) into expression (72), and we get:

$$\begin{aligned} 0 &< |(2c - 1) \cdot 2^{q-1} \cdot 10^{-k-1} - m| < 2^{-36} \\ \Rightarrow 0 &< n^- < 2^{-36} \text{ or } 0 < 1 - n^- < 2^{-36} \end{aligned} \quad (75)$$

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

$$\begin{aligned} \text{float} : n^- &\in (0, 2^{-36}) \cup (1 - 2^{-36}, 1) \\ \text{double} : n^- &\in (0, 2^{-64}) \cup (1 - 2^{-64}, 1) \end{aligned} \quad (76)$$

When  $2^{-1} \cdot 2^q \cdot 10^{-k-1} \neq n$ , it is known from Equation (52) that  $(2c - 1) \cdot 2^{q-1} \cdot 10^{-k-1}$  is not an integer. Therefore, there is:

$$0 < n^- < 1 \quad (77)$$

It is only necessary to prove that Equation (76) 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 (57), there are:

$$(2c - 1) \cdot 2^{q-1} \cdot 10^{-k-1} = (2c - 1) \cdot \frac{x}{y} = \begin{cases} \frac{(2c-1) \cdot 2^{q-k-2}}{5^{k+1}}; q \geq 2 \\ \frac{(2c-1)}{2^{2+k-q} \cdot 5^{k+1}}; 1 \geq q \geq 0 \\ \frac{(2c-1) \cdot 5^{-k-1}}{2^{2+k-q}}; q < 0 \end{cases} \quad (78)$$

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

$$\begin{aligned} \text{float} : (2c - 1)_{\max} &= 2^{25} - 3 \\ \text{double} : (2c - 1)_{\max} &= 2^{54} - 3 \end{aligned} \quad (79)$$

Discuss based on the denominator range in Equation (78).

- $y \leq (2c - 1)_{\max}$   
When  $y \leq (2c - 1)_{\max}$ ,  $y_{\max}$  is the expression (79), the following holds true:

$$\begin{aligned} \frac{1}{y_{\max}} &\leq n^- \leq 1 - \frac{1}{y_{\max}} \\ \frac{1}{y_{\max}} &\leq 1 - n^- \leq 1 - \frac{1}{y_{\max}} \end{aligned} \quad (80)$$

Therefore, when  $y \leq (2c - 1)_{\max}$ , Equation (76) does not hold true.

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

Call function (35) to calculate the approximation results  $P_*/Q_*$  and  $P^*/Q^*$  of all possible upper and lower limit rational numbers:

$$\left( \frac{P_*}{Q_*}, \frac{P^*}{Q^*} \right) = f((2c - 1)_{\max}, x, y) \quad (81)$$

Therefore, for  $n^-$ , the following conclusion can be drawn from formula (34).

$$n^- \in \left[ \frac{(Q_* x) \% y}{y}, \frac{(Q^* x) \% y}{y} \right] \quad (82)$$

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

$$\begin{aligned} \text{float} : 2^{-33} &< n^- < 1 - 2^{-29} \\ \text{double} : 2^{-62} &< n^- < 1 - 2^{-63} \end{aligned} \quad (83)$$

$$\begin{aligned}
\text{float} : & \left[ \frac{(Q_*x)\%y}{y}, \frac{(Q^*x)\%y}{y} \right] \cap (0, 2^{-36}) = \emptyset \\
& \left[ \frac{(Q_*x)\%y}{y}, \frac{(Q^*x)\%y}{y} \right] \cap (1 - 2^{-36}, 1) = \emptyset \\
\text{double} : & \left[ \frac{(Q_*x)\%y}{y}, \frac{(Q^*x)\%y}{y} \right] \cap (0, 2^{-64}) = \emptyset \\
& \left[ \frac{(Q_*x)\%y}{y}, \frac{(Q^*x)\%y}{y} \right] \cap (1 - 2^{-64}, 1) = \emptyset
\end{aligned} \tag{84}$$

Therefore, when  $y > (2c - 1)_{\max}$ , Equation (76) does not hold true.

In summary, when  $2^{-1} \cdot 2^q \cdot 10^{-k-1} \neq n$ , Equation (76) does not hold true, that is,  $\lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor \neq \lfloor 2^{36} \cdot n \rfloor$  must hold true. Therefore, when  $\lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor 2^{36} \cdot n \rfloor$  holds,  $2^{-1} \cdot 2^q \cdot 10^{-k-1} = n$  must hold true. Therefore, Equation (70) holds.

Similarly, it can be proved that when  $\lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor 2^{36} - 2^{36} \cdot n \rfloor$  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:

$$\begin{aligned}
& \lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor 2^{36} - 2^{36} \cdot n \rfloor \\
& \Rightarrow 0 < \left| 2^{35} \cdot 2^q \cdot 10^{-k-1} - 2^{36} + 2^{36} \cdot n \right| < 1 \\
& \Rightarrow 0 < \left| 2^{q-1} \cdot 10^{-k-1} - 1 + n \right| < 2^{-36} \\
& \Rightarrow -2^{-36} < (2c + 1) \cdot 2^{q-1} \cdot 10^{-k-1} - m - 1 < 2^{-36}
\end{aligned} \tag{85}$$

As is known from Equation (56), there is:

$$m < (2c + 1) \cdot 2^{q-1} \cdot 10^{-k-1} < m + 2 \tag{86}$$

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

$$(2c + 1) \cdot 2^{q-1} \cdot 10^{-k-1} - m - 1 = \begin{cases} n^+; & \text{if } (2c + 1) \cdot 2^{q-1} \cdot 10^{-k-1} > m + 1 \\ 1 - n^+; & \text{if } (2c + 1) \cdot 2^{q-1} \cdot 10^{-k-1} < m + 1 \end{cases} \tag{87}$$

Substitute expression (87) into expression (85), and we get:

$$\begin{aligned}
0 < \left| (2c + 1) \cdot 2^{q-1} \cdot 10^{-k-1} - m - 1 \right| < 2^{-36} \\
\Rightarrow 0 < 1 - n^+ < 2^{-36} \text{ or } 0 < n^+ < 2^{-36}
\end{aligned} \tag{88}$$

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

$$\begin{aligned}
\text{float} : n^+ & \in (0, 2^{-36}) \cup (1 - 2^{-36}, 1) \\
\text{double} : n^+ & \in (0, 2^{-64}) \cup (1 - 2^{-64}, 1)
\end{aligned} \tag{89}$$

When  $2^{-1} \cdot 2^q \cdot 10^{-k-1} \neq 1 - n$ , it is known from Equation (53) that  $(2c + 1) \cdot 2^{q-1} \cdot 10^{-k-1}$  is not an integer. Therefore, there is:

$$0 < n^+ < 1 \tag{90}$$

It is only necessary to prove that Equation (89) 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 (57), there are:

$$(2c + 1) \cdot 2^{q-1} \cdot 10^{-k-1} = (2c + 1) \cdot \frac{x}{y} = \begin{cases} \frac{(2c+1) \cdot 2^{q-k-2}}{5^{k+1}}; q \geq 2 \\ \frac{(2c+1)}{2^{2+k-q} \cdot 5^{k+1}}; 1 \geq q \geq 0 \\ \frac{(2c+1) \cdot 5^{-k-1}}{2^{2+k-q}}; q < 0 \end{cases} \quad (91)$$

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

$$\begin{aligned} \text{float} : (2c + 1)_{\max} &= 2^{25} - 1 \\ \text{double} : (2c + 1)_{\max} &= 2^{54} - 1 \end{aligned} \quad (92)$$

Discuss based on the denominator range in Equation (91).

- $y \leq (2c + 1)_{\max}$   
When  $y \leq (2c + 1)_{\max}$ ,  $y_{\max}$  is the expression (92), the following holds true:

$$\begin{aligned} \frac{1}{y_{\max}} &\leq n^+ \leq 1 - \frac{1}{y_{\max}} \\ \frac{1}{y_{\max}} &\leq 1 - n^+ \leq 1 - \frac{1}{y_{\max}} \end{aligned} \quad (93)$$

Therefore, when  $y \leq (2c + 1)_{\max}$ , Equation (89) does not hold true.

- $y > (2c + 1)_{\max}$   
Call function (35) to calculate the approximation results  $\frac{P_*}{Q_*}$  and  $\frac{P^*}{Q^*}$  of all possible upper and lower limit rational numbers:

$$\left( \frac{P_*}{Q_*}, \frac{P^*}{Q^*} \right) = f((2c + 1)_{\max}, x, y) \quad (94)$$

Therefore, for  $n^+$ , the following conclusion can be drawn from formula (34).

$$n^+ \in \left[ \frac{(Q_* x) \% y}{y}, \frac{(Q^* x) \% y}{y} \right] \quad (95)$$

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

$$\begin{aligned} \text{float} : 2^{-33} < n^+ < 1 - 2^{-29} \\ \text{double} : 2^{-62} < n^+ < 1 - 2^{-63} \end{aligned} \quad (96)$$

$$\begin{aligned} \text{float} : \left[ \frac{(Q_* x) \% y}{y}, \frac{(Q^* x) \% y}{y} \right] \cap (0, 2^{-36}) &= \emptyset \\ \left[ \frac{(Q_* x) \% y}{y}, \frac{(Q^* x) \% y}{y} \right] \cap (1 - 2^{-36}, 1) &= \emptyset \\ \text{double} : \left[ \frac{(Q_* x) \% y}{y}, \frac{(Q^* x) \% y}{y} \right] \cap (0, 2^{-64}) &= \emptyset \\ \left[ \frac{(Q_* x) \% y}{y}, \frac{(Q^* x) \% y}{y} \right] \cap (1 - 2^{-64}, 1) &= \emptyset \end{aligned} \quad (97)$$

Therefore, when  $y > (2c + 1)_{\max}$ , Equation (89) does not hold true.

In summary, when  $2^{-1} \cdot 2^q \cdot 10^{-k-1} \neq 1 - n$ , Equation (89) does not hold true, that is,  $\lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor \neq \lfloor 2^{36} - 2^{36} \cdot n \rfloor$  must hold true. Therefore, when  $\lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor 2^{36} - 2^{36} \cdot n \rfloor$  holds,  $2^{-1} \cdot 2^q \cdot 10^{-k-1} = 1 - n$  must hold true. The same is true for double. Therefore, Equation (71) holds.

The following conclusions hold:

$$\begin{aligned} \text{float} : [2^{36} - 2^{36} \cdot n] &= \begin{cases} 2^{36} - 1 - \lfloor 2^{36} \cdot n \rfloor; & \text{if } c \cdot 2^{36+q} \cdot 10^{-k-1} \notin \mathbb{Z} \\ 2^{36} - \lfloor 2^{36} \cdot n \rfloor; & \text{if } c \cdot 2^{36+q} \cdot 10^{-k-1} \in \mathbb{Z} \end{cases} \\ \text{double} : [2^{64} - 2^{64} \cdot n] &= \begin{cases} 2^{64} - 1 - \lfloor 2^{64} \cdot n \rfloor; & \text{if } c \cdot 2^{64+q} \cdot 10^{-k-1} \notin \mathbb{Z} \\ 2^{64} - \lfloor 2^{64} \cdot n \rfloor; & \text{if } c \cdot 2^{64+q} \cdot 10^{-k-1} \in \mathbb{Z} \end{cases} \end{aligned} \quad (98)$$

Discuss whether the following Equation (99) holds when conditions (62) and (63) are met:

$$\begin{aligned} \text{float} : [c \cdot \frac{2^{q+35-k}}{5^{k+1}}] &= \lfloor c \cdot \frac{2^{q+35-k}}{5^{k+1}} \cdot r \rfloor \\ &= \lfloor c \cdot \frac{2^{q+35-k}}{5^{k+1}} \cdot \frac{(2^{63 - \lfloor (-k-1) \cdot \log_2(10) \rfloor} // 10^{k+1}) + 1}{10^{-k-1}} \cdot 2^{\lfloor (-k-1) \cdot \log_2(10) \rfloor - 63} \rfloor \\ \text{double} : [c \cdot \frac{2^{q+63-k}}{5^{k+1}}] &= \lfloor c \cdot \frac{2^{q+63-k}}{5^{k+1}} \cdot r \rfloor \\ &= \lfloor c \cdot \frac{2^{q+63-k}}{5^{k+1}} \cdot \frac{(2^{127 - \lfloor (-k-1) \cdot \log_2(10) \rfloor} // 10^{k+1}) + 1}{10^{-k-1}} \cdot 2^{\lfloor (-k-1) \cdot \log_2(10) \rfloor - 127} \rfloor \end{aligned} \quad (99)$$

There are:

$$\begin{aligned} \text{float} : [c \cdot \frac{2^{q+35-k}}{5^{k+1}}] &= \lfloor 2^{36} \cdot (m+n) \rfloor = 2^{36} \cdot m + \lfloor 2^{36} \cdot n \rfloor \\ \text{double} : [c \cdot \frac{2^{q+63-k}}{5^{k+1}}] &= \lfloor 2^{64} \cdot (m+n) \rfloor = 2^{64} \cdot m + \lfloor 2^{64} \cdot n \rfloor \end{aligned} \quad (100)$$

It has been proven earlier that  $m$  can be accurately calculated. Then, when (99) holds true, the values  $\lfloor 2^{36} \cdot n \rfloor$  and  $\lfloor 2^{64} \cdot n \rfloor$  on the right side of Equations (64) and (65) can be accurately calculated.

From Equation (57), we have:

$$c = \frac{t \cdot 5^{k+1} - 1}{2} \quad (101)$$

Substituting Equation (101) into Equation (99), we have:

$$\begin{aligned} \text{float} : c \cdot \frac{2^{q+35-k}}{5^{k+1}} &= t \cdot 2^{q+34-k} - \frac{2^{q+34-k}}{5^{k+1}} \\ \text{double} : c \cdot \frac{2^{q+63-k}}{5^{k+1}} &= t \cdot 2^{q+62-k} - \frac{2^{q+62-k}}{5^{k+1}} \end{aligned} \quad (102)$$

When conditions (62) and (63) are met,  $t \cdot 2^{q+34-k}$  and  $t \cdot 2^{q+62-k}$  are integers. Under the condition of meeting condition (62), the decimal part of expression (102) is represented as:

$$\begin{aligned} \text{float} : \frac{2^{q+34-k} \% 5^{k+1}}{5^{k+1}}; & 2 \leq q \leq 33 \\ \text{double} : \frac{2^{q+62-k} \% 5^{k+1}}{5^{k+1}}; & 2 \leq q \leq 76 \end{aligned} \quad (103)$$

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 (99) to hold true. That is:

$$\begin{aligned} \text{float} : \frac{2^{q+34-k} \% 5^{k+1}}{5^{k+1}} + \left( c \cdot \frac{2^{q+35-k}}{5^{k+1}} \cdot r - c \cdot \frac{2^{q+35-k}}{5^{k+1}} \right) &< 1 \\ \text{double} : \frac{2^{q+62-k} \% 5^{k+1}}{5^{k+1}} + \left( c \cdot \frac{2^{q+63-k}}{5^{k+1}} \cdot r - c \cdot \frac{2^{q+63-k}}{5^{k+1}} \right) &< 1 \end{aligned} \quad (104)$$

By exhaustively calculating the maximum possible  $c$  value under each  $q$  and substituting it into Equation (104), it holds. The calculation result is in `test2.py`. The calculation results show that for the float range and the double range, Equation (104) always holds true. Therefore, Equation (99) holds true, and thus the values of  $\lfloor 2^{36} \cdot n \rfloor$  and  $\lfloor 2^{64} \cdot n \rfloor$  on the right side of Equations (64) and (65) 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 (64) and (65) can be calculated through lookup tables.

$$\begin{aligned} \text{float} : \lfloor 2^{35} \cdot 2^q \cdot 10^{-k-1} \rfloor &= \text{pow}10 \gg (28 - q - \lfloor (-k - 1) \cdot \log_2(10) \rfloor) \\ \text{double} : \lfloor 2^{63} \cdot 2^q \cdot 10^{-k-1} \rfloor &= \text{pow}10 \gg (64 - q - \lfloor (-k - 1) \cdot \log_2(10) \rfloor) \end{aligned} \quad (105)$$

The code file for verifying the validity of Equation (105) is `test4.py`. Therefore, when conditions (62) and (63) are met, the values of both sides of Equations (64) and (65) can be accurately calculated.

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

$$\begin{aligned} \text{float} : \lfloor c \cdot 2^{q+36} \cdot 10^{-k-1} \rfloor; \lfloor c \cdot 2^{q+36} \cdot r \cdot 10^{-k-1} \rfloor; \\ \text{double} : \lfloor c \cdot 2^{q+64} \cdot 10^{-k-1} \rfloor; \lfloor c \cdot 2^{q+64} \cdot r \cdot 10^{-k-1} \rfloor; \end{aligned} \quad (106)$$

When  $r = 1$ , it is obvious that the two values in expression (106) are equal. When  $r \neq 1$ , or equivalent to  $r > 1$ , has:

$$\begin{aligned} \text{float} : c \cdot 2^{q+36} \cdot r \cdot 10^{-k-1} &= c \cdot 2^{q+36} \cdot 10^{-k-1} + c \cdot 2^{q+36} \cdot (r - 1) \cdot 10^{-k-1} \\ &< c \cdot 2^{q+36} \cdot 10^{-k-1} + 2^{24} \cdot 2^{36} \cdot 2^q \cdot 10^{-k-1} \cdot (r - 1) \\ &< c \cdot 2^{q+36} \cdot 10^{-k-1} + 2^{-3} \\ \lfloor c \cdot 2^{q+36} \cdot r \cdot 10^{-k-1} \rfloor &\leq \lfloor c \cdot 2^{q+36} \cdot 10^{-k-1} \rfloor + 1 \\ \text{double} : c \cdot 2^{q+64} \cdot r \cdot 10^{-k-1} &= c \cdot 2^{q+64} \cdot 10^{-k-1} + c \cdot 2^{q+64} \cdot (r - 1) \cdot 10^{-k-1} \\ &< c \cdot 2^{q+64} \cdot 10^{-k-1} + 2^{53} \cdot 2^{64} \cdot 2^q \cdot 10^{-k-1} \cdot (r - 1) \\ &< c \cdot 2^{q+64} \cdot 10^{-k-1} + 2^{-10} \\ \lfloor c \cdot 2^{q+64} \cdot r \cdot 10^{-k-1} \rfloor &\leq \lfloor c \cdot 2^{q+64} \cdot 10^{-k-1} \rfloor + 1 \end{aligned} \quad (107)$$

Therefore, there is:

$$\begin{aligned} \text{float} : 0 \leq \lfloor c \cdot 2^{q+36} \cdot r \cdot 10^{-k-1} \rfloor - \lfloor c \cdot 2^{q+36} \cdot 10^{-k-1} \rfloor &\leq 1 \\ \text{double} : 0 \leq \lfloor c \cdot 2^{q+64} \cdot r \cdot 10^{-k-1} \rfloor - \lfloor c \cdot 2^{q+64} \cdot 10^{-k-1} \rfloor &\leq 1 \end{aligned} \quad (108)$$

Because there is:

$$\lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor = \lfloor c \cdot 2^q \cdot r \cdot 10^{-k-1} \rfloor = m \quad (109)$$

$$\begin{aligned} \text{float} : \lfloor c \cdot 2^{q+36} \cdot 10^{-k-1} \rfloor &= 2^{36} \cdot m + \lfloor 2^{36} \cdot n \rfloor \\ \text{double} : \lfloor c \cdot 2^{q+64} \cdot 10^{-k-1} \rfloor &= 2^{64} \cdot m + \lfloor 2^{64} \cdot n \rfloor \end{aligned} \quad (110)$$

Suppose:

$$n_r = c \cdot 2^q \cdot r \cdot 10^{-k-1} - m \quad (111)$$

Therefore, the following conclusion can be drawn: when condition (62) is met, from Equation (99), we have:

$$\begin{aligned} \text{float} : 2 \leq q \leq 33 \ \&\& (2c \pm 1) \% 5^{k+1} = 0 \Rightarrow \lfloor 2^{36} \cdot n \rfloor = \lfloor 2^{36} \cdot n_r \rfloor \\ \text{double} : 2 \leq q \leq 76 \ \&\& (2c \pm 1) \% 5^{k+1} = 0 \Rightarrow \lfloor 2^{64} \cdot n \rfloor = \lfloor 2^{64} \cdot n_r \rfloor \end{aligned} \quad (112)$$

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

$$\begin{aligned} \text{float} : \lfloor 2^{36} \cdot n \rfloor &\leq \lfloor 2^{36} \cdot n_r \rfloor \leq \lfloor 2^{36} \cdot n \rfloor + 1 \\ \text{double} : \lfloor 2^{64} \cdot n \rfloor &\leq \lfloor 2^{64} \cdot n_r \rfloor \leq \lfloor 2^{64} \cdot n \rfloor + 1 \end{aligned} \quad (113)$$

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

$$even = (c + 1)\%2 \in \{0, 1\} \quad (114)$$

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:

$$\begin{aligned} \text{float} : \begin{cases} one = 0 : \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + even > \lfloor 2^{36} \cdot n_r \rfloor \\ one = 10 : \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + even > 2^{36} - 1 - \lfloor 2^{36} \cdot n_r \rfloor \end{cases} \\ \text{double} : \begin{cases} one = 0 : \lfloor 2^{q+63} \cdot 10^{-k-1} \rfloor + even > \lfloor 2^{64} \cdot n_r \rfloor \\ one = 10 : \lfloor 2^{q+63} \cdot 10^{-k-1} \rfloor + even > 2^{64} - 1 - \lfloor 2^{64} \cdot n_r \rfloor \end{cases} \end{aligned} \quad (115)$$

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 (116) to determine whether  $one = 0$  or  $one = 10$ .

$$\begin{aligned} \text{float} : \begin{cases} \text{if } \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + even > \lfloor 2^{36} \cdot n_r \rfloor : one = 0 \\ \text{if } \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + even > 2^{36} - 1 - \lfloor 2^{36} \cdot n_r \rfloor : one = 10 \end{cases} \\ \text{double} : \begin{cases} \text{if } \lfloor 2^{q+63} \cdot 10^{-k-1} \rfloor + even > \lfloor 2^{64} \cdot n_r \rfloor : one = 0 \\ \text{if } \lfloor 2^{q+63} \cdot 10^{-k-1} \rfloor + even > 2^{64} - 1 - \lfloor 2^{64} \cdot n_r \rfloor : one = 10 \end{cases} \end{aligned} \quad (116)$$

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 (116) 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 (116) to determine whether  $one \neq 0$  or  $one \neq 10$ . There are a total of four situations. The proof is as follows:

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

$$\begin{aligned} \text{float} : 2^{-1} \cdot 2^q \cdot 10^{-k-1} - n = n^- - 1 \in (2^{-33} - 1, -2^{-29}) \\ \text{double} : 2^{-1} \cdot 2^q \cdot 10^{-k-1} - n = n^- - 1 \in (2^{-62} - 1, -2^{-63}) \end{aligned} \quad (117)$$

Therefore, the following exists:

$$\begin{aligned} \text{float} : 2^{q+35} \cdot 10^{-k-1} - 2^{36} \cdot n \in (2^3 - 2^{36}, -2^7) \\ \text{double} : 2^{q+63} \cdot 10^{-k-1} - 2^{64} \cdot n \in (4 - 2^{64}, -2) \end{aligned} \quad (118)$$

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

$$\begin{aligned} 0 \leq b - \lfloor b \rfloor < 1 \\ a - \lfloor a \rfloor - 1 < b - \lfloor b \rfloor < 1 + a - \lfloor a \rfloor \\ a - b - 1 < \lfloor a \rfloor - \lfloor b \rfloor < a - b + 1 \end{aligned} \quad (119)$$

When  $a = 2^{q+35} \cdot 10^{-k-1}$  and  $b = 2^{36} \cdot n$  or  $a = 2^{q+63} \cdot 10^{-k-1}$  and  $b = 2^{64} \cdot n$ , the following exists:

$$\begin{aligned} \text{float} : \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor - \lfloor 2^{36} \cdot n \rfloor < 2^{q+35} \cdot 10^{-k-1} - 2^{36} \cdot n + 1 \\ \text{double} : \lfloor 2^{q+63} \cdot 10^{-k-1} \rfloor - \lfloor 2^{64} \cdot n \rfloor < 2^{q+63} \cdot 10^{-k-1} - 2^{64} \cdot n + 1 \end{aligned} \quad (120)$$

From Equation (118), we have:

$$\begin{aligned} \text{float} : \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor - \lfloor 2^{36} \cdot n \rfloor < 1 - 2^7 < 0 \\ \text{double} : \lfloor 2^{q+63} \cdot 10^{-k-1} \rfloor - \lfloor 2^{64} \cdot n \rfloor < 1 - 2 < 0 \end{aligned} \quad (121)$$

Therefore, there is:

$$\begin{aligned}
 \text{float} : [2^{q+35} \cdot 10^{-k-1}] + \text{even} &\leq [2^{q+35} \cdot 10^{-k-1}] + 1 \\
 &< [2^{36} \cdot n] \leq [2^{36} \cdot n_r] \\
 \Rightarrow [2^{q+35} \cdot 10^{-k-1}] + \text{even} &< [2^{36} \cdot n_r] \\
 \text{double} : [2^{q+63} \cdot 10^{-k-1}] + \text{even} &\leq [2^{q+63} \cdot 10^{-k-1}] + 1 \\
 &< [2^{64} \cdot n] \leq [2^{64} \cdot n_r] \\
 \Rightarrow [2^{q+63} \cdot 10^{-k-1}] + \text{even} &< [2^{64} \cdot n_r]
 \end{aligned} \tag{122}$$

Therefore, when  $2^{-1} \cdot 2^q \cdot 10^{-k-1} < n$ , the condition (116) can be used to determine that  $\text{one} \neq 0$ .

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

$$\begin{aligned}
 \text{float} : 2^{-1} \cdot 2^q \cdot 10^{-k-1} - n &= n^- \in (2^{-33}, 1 - 2^{-29}) \\
 \text{double} : 2^{-1} \cdot 2^q \cdot 10^{-k-1} - n &= n^- \in (2^{-62}, 1 - 2^{-63})
 \end{aligned} \tag{123}$$

Therefore, the following exists:

$$\begin{aligned}
 \text{float} : 2^{q+35} \cdot 10^{-k-1} - 2^{36} \cdot n &\in (2^3, 2^{36} - 2^7) \\
 \text{double} : 2^{q+63} \cdot 10^{-k-1} - 2^{64} \cdot n &\in (4, 2^{64} - 2)
 \end{aligned} \tag{124}$$

When  $a = 2^{q+35} \cdot 10^{-k-1}$  and  $b = 2^{36} \cdot n$  or  $a = 2^{q+63} \cdot 10^{-k-1}$  and  $b = 2^{64} \cdot n$ , from Equation (119), the following exists:

$$\begin{aligned}
 \text{float} : [2^{q+35} \cdot 10^{-k-1}] - [2^{36} \cdot n] &> 2^{q+35} \cdot 10^{-k-1} - 2^{36} \cdot n - 1 \\
 \text{double} : [2^{q+63} \cdot 10^{-k-1}] - [2^{64} \cdot n] &> 2^{q+63} \cdot 10^{-k-1} - 2^{64} \cdot n - 1
 \end{aligned} \tag{125}$$

From Equation (124), we have:

$$\begin{aligned}
 \text{float} : [2^{q+35} \cdot 10^{-k-1}] - [2^{36} \cdot n] &> 2^3 - 1 \geq 0 \\
 \text{double} : [2^{q+63} \cdot 10^{-k-1}] - [2^{64} \cdot n] &> 4 - 1 \geq 0
 \end{aligned} \tag{126}$$

Therefore, there is:

$$\begin{aligned}
 \text{float} : [2^{q+35} \cdot 10^{-k-1}] + \text{even} &\geq [2^{q+35} \cdot 10^{-k-1}] \\
 &> [2^{36} \cdot n] + 1 \geq [2^{36} \cdot n_r] \\
 \Rightarrow [2^{q+35} \cdot 10^{-k-1}] + \text{even} &> [2^{36} \cdot n_r] \\
 \text{double} : [2^{q+63} \cdot 10^{-k-1}] + \text{even} &\geq [2^{q+63} \cdot 10^{-k-1}] \\
 &> [2^{64} \cdot n] + 1 \geq [2^{64} \cdot n_r] \\
 \Rightarrow [2^{q+63} \cdot 10^{-k-1}] + \text{even} &> [2^{64} \cdot n_r]
 \end{aligned} \tag{127}$$

Therefore, when  $2^{-1} \cdot 2^q \cdot 10^{-k-1} > n$ , the condition (116) can be used to determine that  $\text{one} = 0$ .

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

$$\begin{aligned}
 \text{float} : 2^{-1} \cdot 2^q \cdot 10^{-k-1} + n &= n^+ \in (2^{-33}, 1 - 2^{-29}) \\
 \text{double} : 2^{-1} \cdot 2^q \cdot 10^{-k-1} + n &= n^+ \in (2^{-62}, 1 - 2^{-63})
 \end{aligned} \tag{128}$$

Therefore, the following exists:

$$\begin{aligned} \text{float} : 2^{q+35} \cdot 10^{-k-1} + 2^{36} \cdot n &\in (2^3, 2^{36} - 2^7) \\ \text{double} : 2^{q+63} \cdot 10^{-k-1} + 2^{64} \cdot n &\in (4, 2^{64} - 2) \end{aligned} \quad (129)$$

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

$$\begin{aligned} a - 1 &< \lfloor a \rfloor \leq a \\ b - 1 &< \lfloor b \rfloor \leq b \\ a + b - 2 &< \lfloor a \rfloor + \lfloor b \rfloor \leq a + b \end{aligned} \quad (130)$$

When  $a = 2^{q+35} \cdot 10^{-k-1}$  and  $b = 2^{36} \cdot n$  or  $a = 2^{q+63} \cdot 10^{-k-1}$  and  $b = 2^{64} \cdot n$ , the following exists:

$$\begin{aligned} \text{float} : \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + \lfloor 2^{36} \cdot n \rfloor &\leq 2^{q+35} \cdot 10^{-k-1} + 2^{36} \cdot n \\ \text{double} : \lfloor 2^{q+63} \cdot 10^{-k-1} \rfloor + \lfloor 2^{64} \cdot n \rfloor &\leq 2^{q+63} \cdot 10^{-k-1} + 2^{64} \cdot n \end{aligned} \quad (131)$$

From Equation (129), we have:

$$\begin{aligned} \text{float} : \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + \lfloor 2^{36} \cdot n \rfloor &< 2^{36} - 2^7 \\ \text{double} : \lfloor 2^{q+63} \cdot 10^{-k-1} \rfloor + \lfloor 2^{64} \cdot n \rfloor &< 2^{64} - 2 \end{aligned} \quad (132)$$

Therefore, there is:

$$\begin{aligned} \text{float} : \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + \text{even} &\leq \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + 1 \\ &< 2^{36} - 2 - \lfloor 2^{36} \cdot n \rfloor \\ &< 2^{36} - 1 - \lfloor 2^{36} \cdot n_r \rfloor \\ \Rightarrow \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + \text{even} &< 2^{36} - 1 - \lfloor 2^{36} \cdot n_r \rfloor \\ \text{double} : \lfloor 2^{q+63} \cdot 10^{-k-1} \rfloor + \text{even} &\leq \lfloor 2^{q+63} \cdot 10^{-k-1} \rfloor + 1 \\ &< 2^{64} - 2 - \lfloor 2^{64} \cdot n \rfloor \\ &< 2^{64} - 1 - \lfloor 2^{64} \cdot n_r \rfloor \\ \Rightarrow \lfloor 2^{q+63} \cdot 10^{-k-1} \rfloor + \text{even} &< 2^{64} - 1 - \lfloor 2^{64} \cdot n_r \rfloor \end{aligned} \quad (133)$$

Therefore, when  $2^{-1} \cdot 2^q \cdot 10^{-k-1} < 1 - n$ , the condition (116) can be used to determine that  $one \neq 10$ .  
(4) When  $2^{-1} \cdot 2^q \cdot 10^{-k-1} > 1 - n$ , there must exist  $one = 10$ , and there is:

$$\begin{aligned} \text{float} : 2^{-1} \cdot 2^q \cdot 10^{-k-1} + n = n^+ + 1 &\in (1 + 2^{-33}, 2 - 2^{-29}) \\ \text{double} : 2^{-1} \cdot 2^q \cdot 10^{-k-1} + n = n^+ + 1 &\in (1 + 2^{-62}, 2 - 2^{-63}) \end{aligned} \quad (134)$$

Therefore, the following exists:

$$\begin{aligned} \text{float} : 2^{q+35} \cdot 10^{-k-1} + 2^{36} \cdot n &\in (2^3 + 2^{36}, 2^{37} - 2^7) \\ \text{double} : 2^{q+63} \cdot 10^{-k-1} + 2^{64} \cdot n &\in (4 + 2^{64}, 2^{65} - 2) \end{aligned} \quad (135)$$

When  $a = 2^{q+35} \cdot 10^{-k-1}$  and  $b = 2^{36} \cdot n$  or  $a = 2^{q+63} \cdot 10^{-k-1}$  and  $b = 2^{64} \cdot n$ , from Equation (130), the following exists:

$$\begin{aligned} \text{float} : \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + \lfloor 2^{36} \cdot n \rfloor &> 2^{q+35} \cdot 10^{-k-1} + 2^{36} \cdot n - 2 \\ \text{double} : \lfloor 2^{q+63} \cdot 10^{-k-1} \rfloor + \lfloor 2^{64} \cdot n \rfloor &> 2^{q+63} \cdot 10^{-k-1} + 2^{64} \cdot n - 2 \end{aligned} \quad (136)$$

From Equation (135), we have:

$$\begin{aligned} \text{float} : [2^{q+35} \cdot 10^{-k-1}] + [2^{36} \cdot n] &> 2^{36} + 2^3 - 2 > 2^{36} \\ \text{double} : [2^{q+63} \cdot 10^{-k-1}] + [2^{64} \cdot n] &> 2^{64} + 2 - 2 > 2^{64} \end{aligned} \quad (137)$$

Therefore, there is:

$$\begin{aligned} \text{float} : [2^{q+35} \cdot 10^{-k-1}] + \text{even} &\geq [2^{q+35} \cdot 10^{-k-1}] \\ &> 2^{36} - [2^{36} \cdot n] \\ &> 2^{36} - 1 - [2^{36} \cdot n_r] \\ \Rightarrow [2^{q+35} \cdot 10^{-k-1}] + \text{even} &> 2^{36} - 1 - [2^{36} \cdot n_r] \\ \text{double} : [2^{q+63} \cdot 10^{-k-1}] + \text{even} &\geq [2^{q+63} \cdot 10^{-k-1}] \\ &> 2^{64} - [2^{64} \cdot n] \\ &> 2^{64} - 1 - [2^{64} \cdot n_r] \\ \Rightarrow [2^{q+63} \cdot 10^{-k-1}] + \text{even} &> 2^{64} - 1 - [2^{64} \cdot n_r] \end{aligned} \quad (138)$$

Therefore, when  $2^{-1} \cdot 2^q \cdot 10^{-k-1} > 1 - n$ , the condition (116) can be used to determine that  $\text{one} = 10$ .

From the above proof, it can be seen that when condition (62) is met, the condition (116) can be used to determine whether  $\text{one} = 0$  or  $\text{one} = 10$  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 (116) can be used to determine whether  $\text{one} = 0$  or  $\text{one} = 10$ . 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 (116) can be used to determine whether  $\text{one} \neq 0$  or  $\text{one} \neq 10$ .

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.

Ultimately, we reached the following conclusion:  $\text{one}$  can quickly determine whether  $\text{one}$  equals 0 or 10 by using the following method.

$$\begin{aligned} \text{float} : [2^{q+35} \cdot 10^{-k-1}] + \text{even} &> [2^{36} \cdot n_r] \Rightarrow \text{one} = 0 \\ [2^{q+35} \cdot 10^{-k-1}] + \text{even} &> 2^{36} - 1 - [2^{36} \cdot n_r] \Rightarrow \text{one} = 10 \\ \text{double} : [2^{q+63} \cdot 10^{-k-1}] + \text{even} &> [2^{64} \cdot n_r] \Rightarrow \text{one} = 0 \\ [2^{q+63} \cdot 10^{-k-1}] + \text{even} &> 2^{64} - 1 - [2^{64} \cdot n_r] \Rightarrow \text{one} = 10 \end{aligned} \quad (139)$$

### 3.7. Efficient Computation of $\lfloor 10n \rfloor$ and Rounding

Determine whether  $\text{one}$  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.

#### 3.7.1. $10n - \lfloor 10n \rfloor = 0.5$

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

$$\begin{aligned} 10n - \lfloor 10n \rfloor &= 0.5 \\ \Rightarrow 10 \cdot c \cdot 2^q \cdot 10^{-k-1} - \lfloor 10 \cdot c \cdot 2^q \cdot 10^{-k-1} \rfloor &= 0.5 \\ \Rightarrow c \cdot 2^q \cdot 10^{-k} - \lfloor c \cdot 2^q \cdot 10^{-k} \rfloor &= 0.5 \\ \Rightarrow c \cdot 2^q \cdot 10^{-k} &= \lfloor c \cdot 2^q \cdot 10^{-k} \rfloor + 0.5 \\ \Rightarrow 2c \cdot 2^q \cdot 10^{-k} &= 2\lfloor c \cdot 2^q \cdot 10^{-k} \rfloor + 1 \end{aligned} \quad (140)$$

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

$$c \cdot 2^{q+1} \cdot 10^{-k} = c \cdot 2^{q-k+1} \cdot 5^{-k} \quad (141)$$

According to the range of  $q$ , there are:

$$c \cdot 2^{q+1} \cdot 10^{-k} = \begin{cases} \frac{c \cdot 2^{q-k+1}}{5^k}; q \geq 0 \\ c \cdot 2 \cdot 5^{-k}; q = -1 \\ \frac{c \cdot 5^{-k}}{2^{k-q-1}}; q \leq -2 \end{cases} \quad (142)$$

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:

$$\begin{aligned} \text{float} : c \geq 2^{k-q-1} &\Rightarrow k - q - 1 \leq 22 \Rightarrow q \geq -34 \\ \text{double} : c \geq 2^{k-q-1} &\Rightarrow k - q - 1 \leq 51 \Rightarrow q \geq -75 \end{aligned} \quad (143)$$

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 (141) is an odd number:

$$\begin{aligned} \text{float} : -34 \leq q \leq -2 \ \&\& \ c \% 2^{k-q} = 2^{k-q-1} \\ \text{double} : -75 \leq q \leq -2 \ \&\& \ c \% 2^{k-q} = 2^{k-q-1} \end{aligned} \quad (144)$$

When  $q$  is within the above range (144),  $r = 1$  is derived from Equation (28). Therefore, there is:

$$n_r = n \quad (145)$$

The following equation holds:

$$20m + 20n = c \cdot 2^q \cdot 10^{-k+1} = c \cdot 2^{q-k+1} \cdot 5^{-k} = \frac{c}{2^{k-q-1}} \cdot 5^{-k} \quad (146)$$

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:

$$\begin{aligned} 20n &\in \{5, 15\} \\ \Rightarrow n &\in \{0.25, 0.75\} \\ \Rightarrow n_r &\in \{0.25, 0.75\} \end{aligned} \quad (147)$$

The result of *one* is an even number between  $\lfloor 10n \rfloor$  and  $\lfloor 10n \rfloor + 1$ . Therefore:

$$\text{one} = \begin{cases} \lfloor 10n \rfloor = 2, \text{ if } n = 0.25 \\ \lfloor 10n \rfloor + 1 = 8, \text{ if } n = 0.75 \end{cases} \Rightarrow \text{one} = \lfloor 20n + 1 \rfloor // 2 - (n = 0.25 ? 1 : 0) \quad (148)$$

### 3.7.2. $10n - \lfloor 10n \rfloor \neq 0.5$

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:

$$\text{one} = \begin{cases} \lfloor 10n \rfloor, \text{ if } 10n - \lfloor 10n \rfloor < 0.5 \\ \lfloor 10n \rfloor + 1, \text{ if } 10n - \lfloor 10n \rfloor > 0.5 \end{cases} \Rightarrow \text{one} = \lfloor 10n + 0.5 \rfloor = \lfloor 20n + 1 \rfloor // 2 \quad (149)$$

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:

$$\begin{aligned} d &= ten + one \\ &= 10m + \lfloor 20n + 1 \rfloor // 2 \\ &= (\lfloor 20m + 20n \rfloor + 1) // 2 \end{aligned} \quad (150)$$

Suppose there are:

$$20m + 20n = c \cdot 2^{q+1} \cdot 10^{-k} = c \cdot 2^{q-k+1} \cdot 5^{-k} = c \cdot \frac{x}{y} \quad (151)$$

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

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

$$\begin{aligned} \text{float} : \frac{1}{2^{24} - 1} = \frac{1}{C} \leq n_{20} \leq 1 - \frac{1}{C} = \frac{2^{24} - 2}{2^{24} - 1} \\ \text{double} : \frac{1}{2^{53} - 1} = \frac{1}{C} \leq n_{20} \leq 1 - \frac{1}{C} = \frac{2^{53} - 2}{2^{53} - 1} \end{aligned} \quad (152)$$

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

$$\begin{aligned} \text{float} : 2^{-32} < n_{20} < 1 - 2^{-30} \\ \text{double} : 2^{-64} < n_{20} < 1 - 2^{-62} \end{aligned} \quad (153)$$

Therefore, the range of  $n_{20}$  satisfies Equation (153). 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:

$$\begin{aligned} \text{float} : n_{36} \in [0, 2^{-36}) \\ \text{double} : n_{70} \in [0, 2^{-70}) \end{aligned} \quad (154)$$

Calculate the boundary conditions of the following expression:

$$\begin{aligned} \text{float} : F = 20 \cdot (c \cdot 2^q \cdot r \cdot 10^{-k-1} - n_{36}) \\ \text{double} : F = 20 \cdot (c \cdot 2^q \cdot r \cdot 10^{-k-1} - n_{70}) \end{aligned} \quad (155)$$

Therefore, there is:

$$\begin{aligned} \text{float} : F_{\min} &> 20 \cdot (c \cdot 2^q \cdot 10^{-k-1} - 2^{-36}) \\ &= 20m + 20n - 20 \cdot 2^{-36} \\ F_{\max} &< 20 \cdot (c \cdot 2^q \cdot (1 + 2^{-63}) \cdot 10^{-k-1} - 0) \\ &< 20m + 20n + 20 \cdot 2^{-63} \cdot c \\ &< 20m + \lfloor 20n \rfloor + 1 \\ \text{double} : F_{\min} &> 20 \cdot (c \cdot 2^q \cdot 10^{-k-1} - 2^{-70}) \\ &= 20m + 20n - 20 \cdot 2^{-70} \\ &> 20m + \lfloor 20n \rfloor \\ F_{\max} &< 20 \cdot (c \cdot 2^q \cdot (1 + 2^{-127}) \cdot 10^{-k-1} - 0) \\ &< 20m + 20n + 20 \cdot 2^{-127} \cdot c \\ &< 20m + \lfloor 20n \rfloor + 1 \end{aligned} \quad (156)$$

Therefore, there is:

$$\begin{aligned} \text{float} : \lfloor F \rfloor &= 20m + \lfloor 20n \rfloor \\ \text{double} : \lfloor F \rfloor &= 20m + \lfloor 20n \rfloor \end{aligned} \quad (157)$$

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:

$$\begin{aligned} d &= \text{ten} + \text{one} \\ &= (\lfloor F \rfloor + 1) // 2 \\ &= (\lfloor 20 \cdot (c \cdot 2^q \cdot r \cdot 10^{-k-1} - n_x) \rfloor + 1) // 2 \end{aligned} \quad (158)$$

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

### 3.7.3. Efficient Implementation of $n = 0.25$ for Double

For double, quickly determine that  $n = 0.25$  in Equation (148). 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  $n = 0.25$ :

$$\text{double} : n = 0.25 \text{ if } \lfloor 2^{64} \cdot n_r \rfloor = 2^{62} \quad (159)$$

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

$$4m + 4n = c \cdot 2^{q+2} \cdot 10^{-k-1} \quad (160)$$

Therefore, when Equation (160) is not an integer, we have (the test file is [test6.py](#)):

$$2^{-62} < 4n - \lfloor 4n \rfloor < 1 - 2^{-62} \quad (161)$$

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

$$\begin{aligned} \lfloor 4n \rfloor = 0 &\Rightarrow 4n - 0 < 1 - 2^{-62} \Rightarrow \lfloor 2^{64} \cdot n \rfloor \leq 2^{62} - 2 \\ \lfloor 4n \rfloor = 1 &\Rightarrow 4n - 1 > 2^{-62} \Rightarrow \lfloor 2^{64} \cdot n \rfloor \geq 2^{62} + 1 \end{aligned} \quad (162)$$

Then there are:

$$\begin{aligned} \lfloor 2^{64} \cdot n \rfloor &\neq 2^{62} \ \&\& \ \lfloor 2^{64} \cdot n \rfloor + 1 \neq 2^{62} \\ &\Rightarrow \lfloor 2^{64} \cdot n_r \rfloor \neq 2^{62} \end{aligned} \quad (163)$$

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

$$\text{double} : n \neq 0.25 \text{ if } \lfloor 2^{64} \cdot n_r \rfloor \neq 2^{62} \quad (164)$$

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

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

### 3.7.4. Efficient Calculation of *one* for Double

In the double range, introduce another faster way to calculate *one*:

$$\text{double} : \text{one} = \lfloor \frac{\lfloor 2^{64} \cdot n_r \rfloor}{2^{64}} \cdot 10 + \left( (n = 0.25) ? 0 : \left( 2^{-1} + \frac{6}{2^{64}} \right) \right) \rfloor \quad (166)$$

The proof of Equation (166) is as follows:

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

when  $n \neq 0.25$ , Equation (166) can be equivalent to the following:

$$\text{double :one} = \lfloor \frac{\lfloor 2^{64} \cdot n_r \rfloor}{2^{64}} \cdot 10 + 2^{-1} + \frac{6}{2^{64}} \rfloor \quad (167)$$

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

$$\text{double :one} = \begin{cases} \lfloor 10n \rfloor, & \text{if } 10n - \lfloor 10n \rfloor < 0.5 \\ 8, & \text{if } 10n - \lfloor 10n \rfloor = 0.5 \\ \lfloor 10n \rfloor + 1, & \text{if } 10n - \lfloor 10n \rfloor > 0.5 \end{cases} = \lfloor 20n + 1 \rfloor // 2 \quad (168)$$

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

$$\lfloor \frac{\lfloor 2^{64} \cdot n_r \rfloor}{2^{64}} \cdot 10 + 2^{-1} + \frac{6}{2^{64}} \rfloor = \begin{cases} \lfloor 10n \rfloor, & \text{if } 10n - \lfloor 10n \rfloor < 0.5 \\ 8, & \text{if } 10n - \lfloor 10n \rfloor = 0.5 \\ \lfloor 10n \rfloor + 1, & \text{if } 10n - \lfloor 10n \rfloor > 0.5 \end{cases} = \lfloor 20n + 1 \rfloor // 2 \quad (169)$$

From the range of  $n$ , there is:

$$\frac{\lfloor 2^{64} \cdot n_r \rfloor}{2^{64}} \in (n_r - 2^{-64}, n_r] \quad (170)$$

Because the following conditions exist:

$$\begin{aligned} c \cdot 2^q \cdot 10^{-k-1} &= m + n \\ c \cdot 2^q \cdot r \cdot 10^{-k-1} &= m + n_r \end{aligned} \quad (171)$$

Therefore, the following relationship can be concluded:

$$\begin{aligned} n_r - n &= (r - 1) \cdot c \cdot 2^q \cdot 10^{-k-1} \\ n_r &= (r - 1) \cdot (m + n) + n \\ \Rightarrow n &\leq n_r < 2^{-127} \cdot c + n \\ n &\leq n_r < 2^{-127} \cdot 2^{53} + n \\ n &\leq n_r < 2^{-74} + n \end{aligned} \quad (172)$$

From Equation (170) and (172), it can be concluded that:

$$\begin{aligned} \frac{\lfloor 2^{64} \cdot n_r \rfloor}{2^{64}} &\in (n - 2^{-64}, n + 2^{-74}) \\ \Rightarrow \frac{\lfloor 2^{64} \cdot n_r \rfloor}{2^{64}} \cdot 10 &\in (10n - 10 \cdot 2^{-64}, 10n + 10 \cdot 2^{-74}) \\ \Rightarrow \frac{\lfloor 2^{64} \cdot n_r \rfloor}{2^{64}} \cdot 20 &\in (20n - 20 \cdot 2^{-64}, 20n + 20 \cdot 2^{-74}) \\ \Rightarrow \frac{\lfloor 2^{64} \cdot n_r \rfloor}{2^{64}} \cdot 20 &\in (\lfloor 20n \rfloor + n_{20} - 20 \cdot 2^{-64}, \lfloor 20n \rfloor + n_{20} + 20 \cdot 2^{-74}) \end{aligned} \quad (173)$$

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

$$\lfloor \frac{\lfloor 2^{64} \cdot n_r \rfloor}{2^{64}} \cdot 20 + 1 + x \rfloor // 2 = \lfloor 20n + 1 \rfloor // 2 = \text{one} \quad (174)$$

Therefore, the following conclusions can be drawn:

$$\begin{aligned} \lfloor 20n \rfloor + n_{20} - 20 \cdot 2^{-64} + 1 + x &\geq \lfloor 20n + 1 \rfloor \Rightarrow x \geq 20 \cdot 2^{-64} - n_{20} \\ \lfloor 20n \rfloor + n_{20} + 20 \cdot 2^{-74} + 1 + x &< \lfloor 20n + 2 \rfloor \Rightarrow x < 1 - 20 \cdot 2^{-74} - n_{20} \end{aligned} \quad (175)$$

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

$$x = 12 \cdot 2^{-64} \geq 20 \cdot 2^{-64} - n_{20} \quad (176)$$

All floating-point numbers that do not meet condition (176) are as follows (in hexadecimal) :

$$\begin{aligned} &0xd17c0747bd76fa1, \\ &0xd27c0747bd76fa1, \\ &0x4d73de005bd620df, \\ &0x4d83de005bd620df, \\ &0x4d93de005bd620df, \end{aligned} \quad (177)$$

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

$$x = 12 \cdot 2^{-64} < 1 - 20 \cdot 2^{-74} - n_{20} \quad (178)$$

All floating-point numbers that do not meet condition (178) are as follows (in hexadecimal) :

$$\begin{aligned} &0x612491daad0ba280, \\ &0x6159b651584e8b20, \\ &0x619011f2d73116f4, \\ &0x61c4166f8cfd5cb1, \\ &0x61d4166f8cfd5cb1, \end{aligned} \quad (179)$$

There are:

$$2\left(\frac{\lfloor 2^{64} \cdot n_r \rfloor}{2^{64}} \cdot 10 + 2^{-1} + \frac{6}{2^{64}}\right) = \frac{\lfloor 2^{64} \cdot n_r \rfloor}{2^{64}} \cdot 20 + 1 + x \quad (180)$$

When the floating-point number is not within the above range (177) and (179), the condition (175) is satisfied. We have tested all floating-point numbers within the above-mentioned range (177) and (179), and the algorithm implementation code has output the correct result, that is, it satisfies the SW principle. The test process file is [test8.py](#).

In summary, Equation (169) and Equation (166) holds. Therefore, Equation (166) can be used to quickly calculate *one*.

### 3.8. 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.

### 3.9. Implementation of Pseudocode

This subsection details the pseudocode implementation for converting regular floating-point numbers to decimal representation. The handling of subnormal (irregular) floating-point numbers is omitted here due to space constraints; interested readers may refer to the source code for the complete implementation.

## 3.9.1. Single-Precision Floating-Point Numbers

$$\begin{aligned}
& \text{uint32 } vi = \text{bit\_copy\_from\_f32}(v) \\
& (\text{uint64 } c, \text{int32 } q) = \text{extract}(vi) \\
& \text{const int BIT} = 36 \\
& \text{const uint64 offset} = (1ULL \ll (\text{BIT} - 2)) - 7 \\
& \quad \text{int32 } k = (q \cdot 1233) \gg 12 \\
& \quad \text{int32 } h = q + ((1701 \cdot (-k - 1)) \gg 9) \\
& \text{uint64 pow10} = \text{get\_pow10}(-k - 1) \\
& \quad \text{uint64 } cb = c \ll (h + \text{BIT} + 1) \\
& \quad \text{uint64 hi64} = \text{umul\_hi}_{64 \times 64}(cb, \text{pow10}) \\
& \quad \text{bool even} = (vi + 1) \& 1 \\
& \quad \text{uint64 half} = (\text{pow10} \gg (65 - (h + \text{BIT} + 1))) + \text{even} \\
& \text{uint32 shorter} = (\text{hi64} + \text{half}) \gg \text{BIT} \\
& \quad \text{uint64 bias} = (\text{hi64} \gg (\text{BIT} - 4)) \& 15 \\
& \text{uint32 longer} = (5 \cdot \text{hi64} + \text{offset} + \text{bias}) \gg (\text{BIT} - 1) \\
& \text{bool updown} = \text{shorter} > ((\text{hi64} - \text{half}) \gg \text{BIT}) \\
& \quad \text{uint32 } d = \text{updown} ? \text{shorter} \cdot 10 : \text{longer}
\end{aligned} \tag{181}$$

Algorithm 181 outlines the procedure for computing  $d$  and  $k$  for single-precision floating-point numbers.

Since  $c$  contains at most 23 significant bits and  $\text{pow10}$  has 64 significant bits, their product contains at most 87 significant bits. Given that the range of  $h$  is  $[-4, -1]$ , the value  $cb = c \ll (h + \text{BIT} + 1)$  does not overflow, preserving all bits of  $c$ . When  $e_{10} = -k - 1$ ,  $\text{pow10}$  is computed using Equation (24). From Equation (25), we derive:

$$\text{pow10} \cdot 2^{\lfloor (-k-1) \cdot \log_2 10 \rfloor - 63} = 10^{-k-1} \cdot r_{1,-k-1} \tag{182}$$

From Equation (51):

$$\begin{aligned}
m &= \lfloor v \cdot 10^{-k-1} \rfloor \\
&= \lfloor v \cdot 10^{-k-1} \cdot r_{1,-k-1} \rfloor \\
&= \lfloor v \cdot \text{pow10} \cdot 2^{\lfloor (-k-1) \cdot \log_2 10 \rfloor - 63} \rfloor \\
&= \lfloor c \cdot \text{pow10} \cdot 2^{q + \lfloor (-k-1) \cdot \log_2 10 \rfloor - 63} \rfloor \\
&= (c \cdot \text{pow10}) \gg (63 - q - \lfloor (-k - 1) \cdot \log_2 10 \rfloor) \\
&= (c \cdot \text{pow10}) \gg (63 - h) \\
&= ((c \ll (37 + h)) \cdot \text{pow10}) \gg (36 + 64) \\
&= \text{hi64} \gg 36
\end{aligned} \tag{183}$$

Let  $up$  indicate whether  $one = 10$ . From Equations (138) and (105):

$$\begin{aligned}
\text{bool up} &= \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + \text{even} > 2^{36} - 1 - \lfloor 2^{36} \cdot n_r \rfloor \\
&= \lfloor 2^{q+35} \cdot 10^{-k-1} \rfloor + \text{even} + \lfloor 2^{36} \cdot n_r \rfloor \geq 2^{36} \\
&= (\text{pow10} \gg (28 - h)) + \text{even} + \lfloor 2^{36} \cdot n_r \rfloor \geq 2^{36} \\
&= (\text{half} + \lfloor 2^{36} \cdot n_r \rfloor) \gg 36
\end{aligned} \tag{184}$$

When  $one \neq 10$ , we have  $d // 10 = m$ ; when  $one = 10$ , we have  $d // 10 = m + 1$ . Therefore:

$$d // 10 = m + up \quad (185)$$

The upper 28 bits of  $hi64$  represent  $m$ , while the lower 36 bits represent  $\lfloor 2^{36} \cdot n_r \rfloor$ :

$$hi64 = (m \ll 36) + \lfloor 2^{36} \cdot n_r \rfloor \quad (186)$$

Consequently:

$$shorter = d // 10 = (hi64 + half) \gg 36 \quad (187)$$

Let  $updown$  indicate whether  $one \in \{0, 10\}$ , equivalently whether the last digit of  $d$  is zero:

$$updown = (d \bmod 10 = 0) \quad (188)$$

For brevity, let  $dot\_one = \lfloor 2^{36} \cdot n_r \rfloor$ . From Equations (134) and (140):

$$\begin{aligned} one = 0 : half &> dot\_one \\ one = 10 : half &> 2^{36} - 1 - dot\_one \end{aligned} \quad (189)$$

The following equations hold:

$$\begin{aligned} one = 0 : (hi64 - half) \gg 36 &= m - 1 \\ one = 10 : (hi64 - half) \gg 36 &= m \end{aligned} \quad (190)$$

Therefore:

$$\begin{aligned} one = 0 : shorter &= m > (hi64 - half) \gg 36 = m - 1 \\ one = 10 : shorter &= m + 1 > (hi64 - half) \gg 36 = m \end{aligned} \quad (191)$$

Thus, the condition for  $one \in \{0, 10\}$  is:

$$shorter > (hi64 - half) \gg 36 \implies d \bmod 10 = 0 \quad (192)$$

When  $updown$  is true,  $d = 10 \cdot shorter$ . Otherwise, we compute  $d$  using Equation (158). When  $10n - \lfloor 10n \rfloor \neq 0.5$ :

$$\begin{aligned} d &= (\lfloor 20 \cdot (v \cdot 10^{-k-1} \cdot r - n_{36}) \rfloor + 1) // 2 \\ &= (\lfloor 20 \cdot (hi64 \cdot 2^{-36}) \rfloor + 1) // 2 \\ &= (\lfloor 5 \cdot hi64 \cdot 2^{-34} \rfloor + 1) // 2 \\ &= (\lfloor (5 \cdot hi64 + 2^{34}) \cdot 2^{-34} \rfloor) // 2 \\ &= ((5 \cdot hi64 + 2^{34}) \gg 34) // 2 \\ &= (5 \cdot hi64 + 2^{34}) \gg 35 \end{aligned} \quad (193)$$

Adding  $2^{34}$  serves as a rounding operation. When  $10n - \lfloor 10n \rfloor = 0.5$ , Equation (148) requires checking whether  $n = 0.25$ . The derivation of longer involves careful handling of edge cases; the correctness of this computation has been verified through exhaustive testing.

## 3.9.2. Double-Precision Floating-Point Numbers

$$\begin{aligned}
& \text{uint64 } vi = \text{bit\_copy\_from\_f64}(v) \\
& (\text{uint64 } c, \text{int } q) = \text{extract}(vi) \\
& \text{const int BIT} = 6 \\
& \text{int } k = (q \cdot 78913) \gg 18 \\
& \text{int } h = q + ((217707 \cdot (-k - 1)) \gg 16) \\
& \text{uint128 } \text{pow10} = \text{get\_pow10}(-k - 1) \\
& \text{uint64 } cb = c \ll (h + \text{BIT} + 1) \\
& \text{uint128 } hi128 = \text{umul\_hi}_{64 \times 128}(cb, \text{pow10}) \\
& \text{bool } \text{even} = (vi + 1) \& 1 \\
& \text{uint64 } \text{half} = (\text{pow10} \gg (64 - h)) + \text{even} \\
& \text{uint64 } \text{dot\_one} = (\text{uint64})(hi128 \gg \text{BIT}) \\
& \text{uint64 } \text{ten} = 10 \cdot (hi128 \gg (\text{BIT} + 64)) \\
& \text{uint64 } \text{offset\_num} = (\text{dot\_one} == 2^{62}) ? 0 : 2^{63} + 6 \\
& \text{uint64 } \text{one} = ((\text{uint128})10 \cdot \text{dot\_one} + \text{offset\_num}) \gg 64 \\
& \text{one} = (\text{half} > \text{dot\_one}) ? 0 : \text{one} \\
& \text{one} = (\text{half} > 2^{64} - 1 - \text{dot\_one}) ? 10 : \text{one} \\
& \text{uint64 } d = \text{ten} + \text{one}
\end{aligned} \tag{194}$$

Algorithm 194 presents the corresponding procedure for double-precision floating-point numbers.  $\text{pow10}$  is computed using Equation (26), satisfying:

$$\text{pow10} \cdot 2^{\lfloor (-k-1) \cdot \log_2 10 \rfloor - 127} = 10^{-k-1} \cdot r_{1,-k-1} \tag{195}$$

From Equation (51):

$$\begin{aligned}
m &= \lfloor v \cdot 10^{-k-1} \rfloor \\
&= \lfloor v \cdot 10^{-k-1} \cdot r_{1,-k-1} \rfloor \\
&= \lfloor v \cdot \text{pow10} \cdot 2^{\lfloor (-k-1) \cdot \log_2 10 \rfloor - 127} \rfloor \\
&= \lfloor c \cdot \text{pow10} \cdot 2^{q + \lfloor (-k-1) \cdot \log_2 10 \rfloor - 127} \rfloor \\
&= (c \cdot \text{pow10}) \gg (127 - q - \lfloor (-k-1) \cdot \log_2 10 \rfloor) \\
&= (c \cdot \text{pow10}) \gg (127 - h) \\
&= ((c \ll (7 + h)) \cdot \text{pow10}) \gg (64 + 70) \\
&= hi128 \gg 70
\end{aligned} \tag{196}$$

Given  $h \in [-4, -1]$  and  $c$  having at most 53 significant bits,  $cb$  does not overflow, ensuring accurate computation of  $\text{ten} = 10 \cdot m$ . By definition,  $\text{dot\_one} = \lfloor 2^{64} \cdot n_r \rfloor$ . From Equation (105):

$$\begin{aligned}
\text{half} &= (\text{pow10} \gg (64 - h)) + \text{even} \\
&= \lfloor 2^{63} \cdot 2^q \cdot 10^{-k-1} \rfloor + \text{even}
\end{aligned} \tag{197}$$

The value of  $\text{one}$  is first computed using Equation (166), then adjusted if the conditions for  $\text{one} = 0$  or  $\text{one} = 10$  are met:

$$\begin{aligned}
\text{one} &= \left\lfloor \frac{\lfloor 2^{64} \cdot n_r \rfloor}{2^{64}} \cdot 10 + \left( (n = 0.25) ? 0 : \left( 2^{-1} + \frac{6}{2^{64}} \right) \right) \right\rfloor \\
&= (\text{dot\_one} \cdot 10 + (\text{dot\_one} == 2^{62} ? 0 : 2^{63} + 6)) \gg 64
\end{aligned} \tag{198}$$

An equivalent formulation is:

$$one = (\text{dot\_one} == 2^{62}) ? 2 : (\text{dot\_one} \cdot 10 + 2^{63} + 6) \gg 64 \quad (199)$$

The conditions  $one = 0$  and  $one = 10$  are determined by Equations (127) and (138), respectively.

In the C/C++ implementation, the only branch occurs when computing  $c$  and  $q$  for subnormal floating-point numbers (marked as unlikely). The conversion of normal floating-point numbers is branch-free, eliminating branch misprediction penalties. As shown in Algorithms 181 and 194, the core algorithm is remarkably concise. Irregular numbers are handled in a separate branch—a worthwhile trade-off given their rarity.

### 3.10. Decimal to String Conversion

The floating-point number printing algorithm consists of two main phases. While the first phase, which computes the decimal digits  $d$  and exponent  $k$ , has been discussed in previous sections, this section focuses on the second phase: converting the computed decimal representation into a string.

Floating-point numbers are typically printed in two formats: fixed-point notation and scientific notation. Table (3) illustrates examples of both formats.

**Table 3.** Example of floating-point number printing result.

float number	fixed-point	scientific
2.34	"2.34"	"2.34e+00"
12	"12.0"	"1.2e+01"
120	"120.0"	"1.2e+02"
0.012	"0.012"	"1.2e-02"

Having computed  $d$  and  $k$ , we now convert the floating-point number to a string based on these values. According to the Schubfach algorithm,  $d$  may contain trailing zeros, meaning  $d \bmod 10$  could be zero. In order to reduce instruction dependencies. We decompose  $d$  into two parts:  $d // 10$  and  $d \bmod 10$ . Let  $up$  indicate whether  $one$  equals 10, and  $updown$  represent the cases where  $one$  is either 0 or 10. The following relationships hold:

$$\begin{aligned} d // 10 &= m + up \\ updown &= (one = 0 \vee one = 10) = (d \bmod 10 = 0) \\ d &= 10 \cdot (m + up) + (updown ? 0 : one) \end{aligned} \quad (200)$$

By calculating the approximate range of  $m$ , we obtain:

$$\begin{aligned} m &= \lfloor c \cdot 2^q \cdot 10^{-k-1} \rfloor \\ c \cdot 2^q \cdot 10^{-k-1} &\in [0.1 \cdot c, c) \\ \Rightarrow [0.1 \cdot c_{\min}] &\leq [0.1 \cdot c] \leq m \leq \lfloor c \rfloor \leq [c_{\max}] \end{aligned} \quad (201)$$

Based on the range of  $c$ , we derive:

$$\begin{aligned} \text{float} : & \begin{cases} \text{normal} : [\frac{2^{23}+1}{10}, 2^{24} - 1] = [838860.9, 16777215] \\ \text{subnormal} : [\frac{1}{10}, 2^{23} - 1] = [0.1, 8388607] \end{cases} \\ \text{double} : & \begin{cases} \text{normal} : [\frac{2^{52}+1}{10}, 2^{53} - 1] = [4.5 \times 10^{14}, 9 \times 10^{15}] \\ \text{subnormal} : [\frac{1}{10}, 2^{52} - 1] = [0.1, 4.5 \times 10^{15}] \end{cases} \end{aligned} \quad (202)$$

Let  $len10(x)$  denote the number of decimal digits in  $x$ . For  $x > 0$ :

$$len10(x) = \lfloor \log_{10}(x) \rfloor + 1 \quad (203)$$

We define  $len10(0) = 0$ . For example,  $len10(12345) = 5$ . Consequently:

$$\begin{aligned} \text{float} : & \begin{cases} \text{normal} : len10(m + up) \in [6, 8] \\ \text{subnormal} : len10(m + up) \in [0, 7] \end{cases} \\ \text{double} : & \begin{cases} \text{normal} : len10(m + up) \in [15, 16] \\ \text{subnormal} : len10(m + up) \in [0, 16] \end{cases} \end{aligned} \quad (204)$$

Since  $d // 10 = m + up$ , we have  $len10(d) = len10(m + up) + 1$ . Let  $tz10(x)$  denote the number of trailing zeros in the decimal representation of  $x$ ; for example,  $tz10(12300) = 2$ . The number of significant digits  $sig10(x)$  equals  $len10(x) - tz10(x)$ . When  $updown = 0$ ,  $one \neq 0$ , thus  $d \bmod 10 = one \neq 0$  and  $tz10(d) = 0$ . Therefore:

$$\begin{aligned} sig10(d) &= updown ? len10(d) - tz10(d) : len10(d) \\ \text{float} : sig10(d) &\in [1, 9] \\ \text{double} : sig10(d) &\in [1, 17] \end{aligned} \quad (205)$$

For normal floating-point numbers,  $len10(d)$  can be computed as:

$$\begin{aligned} \text{float} : len10(d) &= len(m + up) + 1 = 9 - [m + up < 10^7] - [m + up < 10^6] \\ \text{double} : len10(d) &= len(m + up) + 1 = 16 + [(m + up) \geq 10^{15}] \end{aligned} \quad (206)$$

where  $[P]$  denotes the Iverson bracket, which evaluates to 1 if predicate  $P$  is true and 0 otherwise.

Computing the ASCII representation of  $d$  reduces to computing the ASCII codes of  $m + up$  and  $one$ . From the above, we have:

$$\begin{aligned} \text{float} : m + up &< 10^8 \\ \text{double} : m + up &< 10^{16} \end{aligned} \quad (207)$$

These bounds allow efficient computation of the ASCII code for  $m + up$  using CPU SIMD instruction sets.

### Scalar Implementation

We first present a scalar method that does not rely on SIMD instructions.

Let  $x \in [0, 10^8)$  and  $dec\_to\_ascii8(x)$  be a function that computes the ASCII representation of  $x$ . Algorithm (2) describes this process. This version is an optimized version of the itoa algorithm written by Paul Khuong [18].

---

#### Algorithm 2 Convert an 8-digit decimal number to ASCII: $dec\_to\_ascii8(x)$

---

**Input:**  $X$  (type: uint64, range:  $[0, 10^8)$ )

**Output:**  $Y$  (type: uint64),  $tz$  (type: uint)

- 1: uint64 abcd\_efgh =  $X + (0x100000000 - 10000) * ((X * 0x68db8bbULL) >> 40)$
  - 2: uint64 ab\_cd\_ef\_gh =  $abcd\_efgh + (0x10000 - 100) * (((abcd\_efgh * 0x147b) >> 19) \& 0x7f0000007f)$
  - 3: uint64 a\_b\_c\_d\_e\_f\_g\_h =  $ab\_cd\_ef\_gh + (0x100 - 10) * (((ab\_cd\_ef\_gh * 0x67) >> 10) \& 0xf000f000f000f)$
  - 4: uint tz =  $\_\_builtin\_ctzll(a\_b\_c\_d\_e\_f\_g\_h) >> 3$
  - 5: uint64 BCD =  $CPU\_IS\_LITTLE\_ENDIAN ? byteswap(a\_b\_c\_d\_e\_f\_g\_h) : a\_b\_c\_d\_e\_f\_g\_h$
  - 6: uint64  $Y = BCD + 0x3030303030303030$
  - 7: **return**  $Y, tz$
-

The function simultaneously computes  $tz10(x)$  as the return value  $tz$ . Example:  $X = 12345600$ ,  $Y = "12345600"$ ,  $tz = 2$ .

---

**Algorithm 3** Convert a 16-digit decimal number to ASCII:  $dec\_to\_ascii16(x)$ 


---

**Input:**  $X$  (type: uint64, range:  $[0, 10^{16})$ )

**Output:**  $Y$  (type: uint128),  $tz$  (type: uint)

- 1: uint64  $abcdefgh = X // 10^8$
  - 2: uint64  $ijklmnop = X \bmod 10^8$
  - 3: (uint64  $abcdefgh\_ascii$ , uint  $tz1$ ) =  $dec\_to\_ascii8(abcdefgh)$
  - 4: (uint64  $ijklmnop\_ascii$ , uint  $tz2$ ) =  $dec\_to\_ascii8(ijklmnop)$
  - 5: uint128  $Y = CPU\_IS\_LITTLE\_ENDIAN ? (ijklmnop\_ascii \ll 64) + abcdefgh\_ascii : (abcdefgh\_ascii \ll 64) + ijklmnop\_ascii$
  - 6: uint  $tz = (ijklmnop == 0) ? 8 + tz1 : tz2$
  - 7: **return**  $Y, tz$
- 

For  $x \in [0, 10^{16})$ , let  $dec\_to\_ascii16(x)$  compute the ASCII representation of  $x$ . Algorithm (3) describes this process. Example:  $X = 123456001234500$ ,  $Y = "1234560012345600"$ ,  $tz = 2$ .

### Handling Undefined Behavior

Since  $__builtin_ctzll(0)$  is undefined behavior in C, special handling is required. For single-precision floating-point numbers, when  $m + up = 0$  (the six smallest subnormal numbers),  $updown = 0$ , avoiding the undefined case. For double-precision numbers,  $m + up = 0$  only occurs for  $5 \times 10^{-324}$  (the smallest subnormal number). In our implementation, we handle this special case separately at the function entry.

### SIMD Implementation

The SIMD implementations of  $dec\_to\_ascii8$  and  $dec\_to\_ascii16$  follow the same principles as the scalar version. For ARM64 architecture, we use the NEON instruction set; for x64 architecture, we support three variants: AVX512, SSE4.1, and SSE2. Table (4) summarizes these implementations.

**Table 4.** SIMD implementations of  $dec\_to\_ascii8$  and  $dec\_to\_ascii16$

SIMD implementation	Description
NEON [19]	Original author: Dougall Johnson. Runs on ARM processors with NEON instruction set.
SSE2 [20]	Based on scalar version; requires only SSE2 instruction set.
SSE4.1	Nearly identical to SSE2 implementation; requires SSE4.1 instruction set.
AVX512 [21]	Original author: Daniel Lemire. Requires AVX512IFMA and AVX512VBMI instruction sets.

Using the methods above, we compute the ASCII representation of  $m + up$  and the number of significant digits. Printing  $d$  is equivalent to printing  $m + up$  and  $one$ . Based on  $sig10(d)$ , we determine which buffer contents to retain. In our implementation, we adopt different formats for floating-point numbers in different ranges to enhance readability, as shown in Table (5).

**Table 5.** Printing formats for different ranges.

Type	Fixed-point	Scientific
float	$[10^{-3}, 10^7)$	other ranges
double	$[10^{-4}, 10^{16})$	other ranges

In scientific notation, the result includes an exponent. For example, in "1.2e-02", "e-02" represents the exponent. Converting  $d \cdot 10^k$  to standard decimal scientific notation yields:

$$d \cdot 10^k = \left( d \cdot 10^{-\lfloor \log_{10}(d) \rfloor} \right) \cdot 10^{k + \lfloor \log_{10}(d) \rfloor} \quad (208)$$

Since  $d \cdot 10^{-\lfloor \log_{10}(d) \rfloor} \in [1, 10)$ , the exponent is determined by  $k + \lfloor \log_{10}(d) \rfloor$ . Let  $E_{10} = k + \lfloor \log_{10}(d) \rfloor$ . Then:

$$E_{10} = k + \lfloor \log_{10}(d) \rfloor = k + \lfloor \log_{10}(m + up) \rfloor + 1 \quad (209)$$

For single-precision floating-point numbers, fixed-point notation is used when  $E_{10} \in [-3, 6]$ . For double-precision numbers, fixed-point notation is used when  $E_{10} \in [-4, 15]$ .

The floating-point number printing algorithm proposed in this paper is illustrated as Algorithm 4. The key distinction between our algorithm and others is the use of SIMD instruction sets for converting  $m + up$  to ASCII values. For single-precision numbers, we use `dec_to_ascii8`; for double-precision, we use `dec_to_ascii16`. The detailed implementation is available in the source code.

---

#### Algorithm 4 Floating-point number printing algorithm

---

**Input:**  $v$  (type: float/double),  $buffer$  (type: char\*)

**Output:**  $buffer$  (type: char\*)

- 1: compute  $m + up$ ,  $updown$ ,  $one$ ,  $k$  from  $v$
  - 2: compute  $sig_{10}(d) = updown ? len_{10}(m + up) - tz_{10}(m + up) : len_{10}(m + up) + 1$
  - 3: compute  $E_{10} = k + \lfloor \log_{10}(m + up) \rfloor + 1$
  - 4: convert  $m + up$ ,  $one$  to ASCII codes and store in  $buffer$
  - 5: based on  $sig_{10}(d)$  and  $E_{10}$ , determine which buffer contents to retain
  - 6: if result is in scientific notation, print  $E_{10}$
  - 7: **return**  $buffer$
- 

In our C implementation, all branches are designed as unlikely branches to minimize the impact of branch prediction failures. By transforming the computation of  $d$  into computing  $m + up$  and  $one$ , we reduce instruction dependencies and enhance instruction-level parallelism. Unlike other algorithms that compute the exact value of  $d$  before printing, we avoid computing  $d$  directly and instead quickly compute  $d / 10 = m + up$ .

### 3.11. Summary

This chapter explains how to quickly calculate  $d$  and  $k$ , as well as the printing optimization. Since  $d = 10m + one$ , the process of calculating  $d$  is transformed into calculating  $m$  and  $one$ .

- Section 3.5 introduces the method for quickly calculating  $m$ .
- Section 3.6 and Section 3.7 introduce the methods for quickly calculating  $one$ .
- Section 3.9 provides a detailed description of the pseudo-code implementation.
- Section 3.10 discusses print optimization.

**Table 6.** All algorithms in the benchmark test.

algorithm	float	double	description
Schubfach	Schubfach32	Schubfach64	author:Raffaello Giuliatti, <a href="https://github.com/c4f7fccc9cb06515/Schubfach">https://github.com/c4f7fccc9cb06515/Schubfach</a> .
Schubfach_xjb	Schubfach32_xjb	Schubfach64_xjb	It is improved by Schubfach and has the same output result.
Ryu	Ryu32	Ryu64	author:Ulf Adams, <a href="https://github.com/ulfjack/ryu">https://github.com/ulfjack/ryu</a> .
Dragonbox	Dragonbox32	Dragonbox64	author:Junekey Jeon, <a href="https://github.com/jk-jeon/Dragonbox">https://github.com/jk-jeon/Dragonbox</a> .
fmt [22]	fmt32	fmt64	author:Victor Zverovich, <a href="https://github.com/fmtlib/fmt">https://github.com/fmtlib/fmt</a> version:12.1.0
yy_double	-	yy_double	author:Guo YaoYuan, <a href="https://github.com/ibireme/yyjson">link:yy_double</a> .
yy_json [23]	yy_json32	yy_json64	author:Guo YaoYuan, <a href="https://github.com/ibireme/yyjson">https://github.com/ibireme/yyjson</a> version:0.12.0
teju_jagua [24]	teju32	teju64	author:Cassio Neri, <a href="https://github.com/cassioneri/teju_jagua">https://github.com/cassioneri/teju_jagua</a> .
xjb	xjb32	xjb64	this paper, <a href="https://github.com/xjb714/xjb">https://github.com/xjb714/xjb</a> .
zmij	zmij32	zmij64	author:Victor Zverovich, <a href="https://github.com/vitaut/zmij">https://github.com/vitaut/zmij</a> .
jnum [25]	jnum32	jnum64	author:Jing Leng, <a href="https://github.com/lengjingzju/json/jnum.c">https://github.com/lengjingzju/json/jnum.c</a> .
uscalec	-	uscalec	author:Russ Cox, <a href="https://github.com/rsc/fpfmt">src:https://github.com/rsc/fpfmt</a> commit 6255750 (19 Jan 2026).

## 4. Experimental Evaluation

This section presents a comprehensive experimental evaluation of the xjb algorithm, comparing its performance against state-of-the-art floating-point to string conversion algorithms across multiple hardware platforms and compilers.

### 4.1. Correctness Verification

Before conducting performance benchmarks, we verified the correctness of the xjb algorithm through extensive testing:

- **Single-precision (float):** We tested all possible input values ( $2^{32}$  values) and verified that the output matches the Schubfach algorithm, ensuring complete correctness for the binary32 format.
- **Double-precision (double):** Due to the impracticality of exhaustively testing all  $2^{64}$  possible values, we conducted comprehensive random testing with a large sample size to verify correctness for the binary64 format.

All test results confirmed that xjb produces output identical to Schubfach while satisfying the SW principle.

### 4.2. Experimental Setup

#### 4.2.1. Hardware Platforms

We conducted benchmarks on two representative hardware platforms:

- **AMD Ryzen 7 7840H:** A modern x86-64 processor with support for AVX2 and AVX-512 instruction sets, running Ubuntu 26.04.
- **Apple M1:** A representative ARM-based processor with NEON SIMD support, running macOS 26.4.

#### 4.2.2. Compilers and Compilation Flags

- **AMD R7-7840H:** Intel C++ Compiler (icpx) version 2025.0.4
- **Apple M1:** Apple Clang version 21.0.0

All benchmarks were compiled with the optimization flags `-O3 -march=native` to enable maximum compiler optimizations and architecture-specific instruction generation.

#### 4.2.3. Benchmark Methodology

The benchmark methodology was designed to provide fair and reproducible performance measurements:

1. Generate  $2^{24}$  random floating-point numbers, excluding special values (0, NaN, and Inf).
2. Measure the total time required to convert all numbers to decimal representation.
3. Calculate the average conversion time per floating-point number.

This methodology ensures statistical significance while avoiding the overhead of special-case handling that would skew results.

### 4.3. Algorithms Compared

We compared xjb against the following state-of-the-art algorithms:

Tables 7 and (8) present the benchmark results for float-to-decimal and double-to-decimal conversions on AMD R7-7840H and Apple M1, respectively. Figure 1a–f presents the benchmark results for the conversion from float to string and from double to string. Figures 1a–d present the benchmark results when the input values are completely random. Since the significant digit length range of the printing results is from 1 to 17, comparative tests were also conducted for different effective lengths, as shown in Figures 1e,f.

- **Schubfach:** The baseline algorithm from which xjb is derived.

- **Schubfach\_xjb**: A variant of Schubfach with xjb optimizations.
- **Ryu**: A widely-used high-performance algorithm.
- **Dragonbox**: Another state-of-the-art algorithm with competitive performance.
- **fmt**: A modern formatting library.
- **yy\_json / yy\_double**: Fast algorithms optimized for JSON serialization.
- **teju\_jagua**: A recent algorithm (note: only supports float/double to decimal conversion).
- **zmij**: A modern algorithm with competitive performance.
- **uscalec**: An algorithm supporting double-precision only.

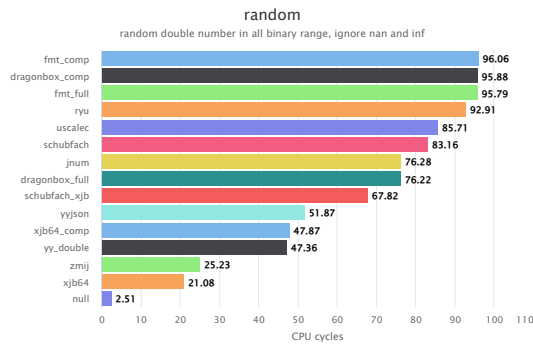
#### 4.4. Performance Results

**Table 7.** Float/double to decimal benchmark results on AMD R7-7840H and ubuntu 26.04. The unit is nanosecond(ns).

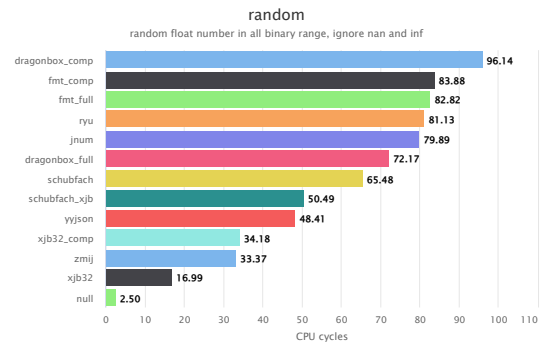
algorithm	float	double
	icpx 2025.0.4	icpx 2025.0.4
Schubfach	12.2	11.51
Schubfach_xjb	4.44	6.33
Ryu	14.02	13.08
Dragonbox	10.19	10.05
yy_json	4.67	5.72
yy_double	-	5.24
teju_jagua	14.99	14.37
zmij	4.76	4.78
uscalec	-	11.27
xjb	2.24	3.76

**Table 8.** Float/double to decimal benchmark results on Apple M1 and MacOS 26.4. The unit is nanosecond(ns).

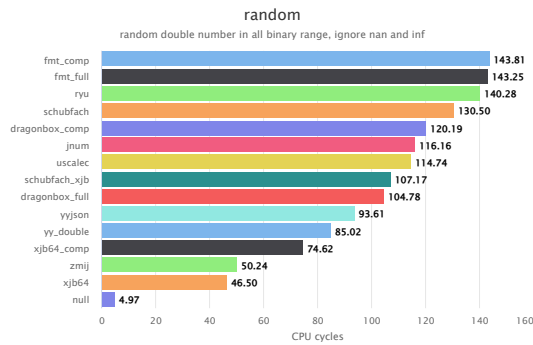
algorithm	float	double
	apple clang 21.0.0	apple clang 21.0.0
Schubfach	11.64	13.12
Schubfach_xjb	5.16	6.58
Ryu	15.75	14.16
Dragonbox	11.78	12.03
yy_json	3.97	4.46
yy_double	-	4.08
teju_jagua	20.25	18.66
zmij	4.11	3.83
uscalec	-	15.26
xjb	2.15	2.58



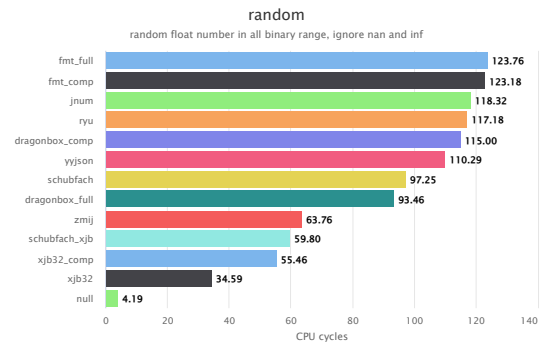
(a) Apple M1 random double benchmark result



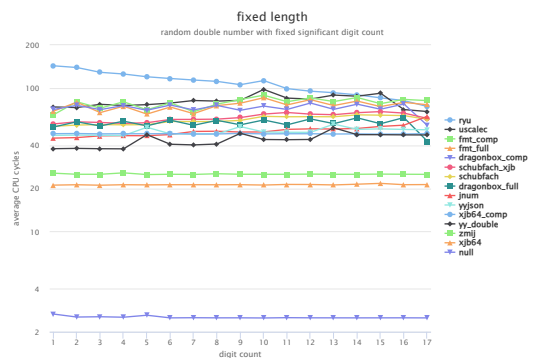
(b) Apple M1 random float benchmark result



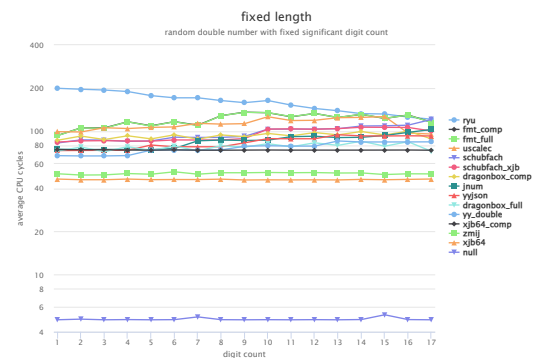
(c) AMD R7-7840H random double benchmark result



(d) AMD R7-7840H random float benchmark result



(e) Apple M1 fixed-length double benchmark result



(f) AMD R7-7840H fixed-length double benchmark result

Figure 1. float/double to string benchmark results.

## 4.5. Analysis and Discussion

### 4.5.1. Performance Comparison

The benchmark results demonstrate that xjb achieves superior performance compared to all other algorithms across both hardware platforms:

#### On AMD R7-7840H:

- For float-to-decimal conversion, xjb achieves 2.24 ns, which is  $1.98\times$  faster than the second-fastest algorithm (Schubfach\_xjb at 4.44 ns) and  $5.44\times$  faster than the baseline Schubfach (12.2 ns).
- For double-to-decimal conversion, xjb achieves 3.76 ns, which is  $1.27\times$  faster than the second-fastest algorithm (yy\_double at 5.24 ns) and  $3.06\times$  faster than the baseline Schubfach (11.51 ns).
- For float-to-string conversion, xjb is about 60% faster than zmij.
- For double-to-string conversion, xjb is about 8% faster than zmij.

#### On Apple M1:

- For float-to-decimal conversion, xjb achieves 2.15 ns, which is  $1.84\times$  faster than the second-fastest algorithm (zmij at 4.11 ns) and  $5.41\times$  faster than the baseline Schubfach (11.64 ns).

- For double-to-decimal conversion, xjb achieves 2.58 ns, which is  $1.49\times$  faster than the second-fastest algorithm (yy\_double at 4.08 ns) and  $5.08\times$  faster than the baseline Schubfach (13.12 ns).
- For float-to-string conversion, xjb is about 96% faster than zmij.
- For double-to-string conversion, xjb is about 20% faster than zmij.

#### 4.5.2. Performance Consistency

A key advantage of the xjb algorithm is its consistent performance across different input distributions. Since all branches in the xjb implementation are designed to be unlikely branches, the algorithm maintains a high branch prediction success rate regardless of input patterns. This property is particularly valuable in real-world applications where input distributions may vary unpredictably.

#### 4.5.3. Cross-Platform Performance

The benchmark results demonstrate that xjb performs well across different processor architectures:

- On x86-64 (AMD R7-7840H), xjb benefits from the compiler's ability to generate efficient code for the algorithm's arithmetic operations.
- On ARM64 (Apple M1), xjb maintains its performance advantage, demonstrating the algorithm's portability and effectiveness across different instruction set architectures.

#### 4.5.4. Comparison with Related Algorithms

- **vs. Schubfach:** xjb achieves  $3\text{--}5\times$  speedup over the baseline Schubfach algorithm, validating the effectiveness of our optimizations in reducing multiplication operations and minimizing branch prediction penalties.
- **vs. yy\_json/yy\_double:** While these algorithms are highly optimized for JSON serialization workloads, xjb outperforms them by  $1.3\text{--}1.8\times$ , demonstrating broader applicability beyond specific use cases.
- **vs. zmij:** xjb achieves  $1.5\text{--}1.9\times$  speedup over zmij, another recent high-performance algorithm, showing the benefits of our approach to instruction dependency reduction.
- **vs. Ryu and Dragonbox:** xjb significantly outperforms these widely-used algorithms ( $2.5\text{--}6\times$  faster), demonstrating that the optimizations in xjb provide substantial benefits over established approaches.

#### 4.6. Summary

The experimental evaluation demonstrates that xjb achieves state-of-the-art performance for floating-point to decimal conversion across multiple hardware platforms. The algorithm's key strengths include:

1. Consistently superior performance across both x86-64 and ARM64 architectures.
2. Significant speedups over the baseline Schubfach algorithm ( $3\text{--}5\times$ ).
3. Competitive or superior performance compared to all other state-of-the-art algorithms.
4. Stable performance regardless of input distribution due to effective branch prediction optimization.

These results validate the effectiveness of our optimization strategies: reducing multiplication operations, minimizing instruction dependencies, and designing branch patterns that work well with modern processor pipelines.

## 5. Conclusions

This paper presented xjb, a novel high-performance algorithm for converting IEEE 754 floating-point numbers to decimal string representations. Building upon the Schubfach algorithm, xjb introduces several key optimizations that significantly improve conversion speed while maintaining full compliance with the Steele-White principle for accurate and minimal-length output.

### 5.1. Summary of Contributions

The main contributions of this work are:

1. **Reduced computational complexity:** By restructuring the calculation process and minimizing the number of high-precision multiplication operations, xjb achieves substantial performance improvements over the baseline Schubfach algorithm.
2. **Optimized branch structure:** The algorithm employs unlikely branch patterns that enable efficient branch prediction on modern processors, resulting in consistent performance across diverse input distributions.
3. **Cross-platform portability:** xjb demonstrates excellent performance across both x86-64 and ARM64 architectures, making it suitable for deployment in heterogeneous computing environments.
4. **Comprehensive validation:** The algorithm has been thoroughly tested and verified to produce correct output for all IEEE 754 binary32 and binary64 floating-point values.

### 5.2. Experimental Findings

Our experimental evaluation on AMD R7-7840H and Apple M1 processors demonstrates that xjb achieves state-of-the-art performance:

- 3–5× speedup over the baseline Schubfach algorithm
- Faster than other high-performance algorithms such as yy\_json, yy\_double, and zmij
- Consistent performance across different input patterns and hardware platforms

### 5.3. Practical Implications

The xjb algorithm has immediate practical applications in numerous domains:

- **Data serialization:** JSON and other text-based data formats require efficient floating-point to string conversion for serialization operations.
- **Scientific computing:** Applications that output numerical results in human-readable format benefit from faster conversion without sacrificing accuracy.
- **Database systems:** Export operations and query result formatting can leverage xjb for improved throughput.
- **Web services:** RESTful APIs and web applications that return numerical data can achieve lower latency with efficient conversion.

### 5.4. Limitations and Future Work

While xjb demonstrates strong performance, several areas warrant further investigation:

1. **Extended precision support:** Future work could extend xjb to support extended precision formats (e.g., 80-bit and 128-bit floating-point numbers) for applications requiring higher precision.
2. **SIMD vectorization:** Although xjb is designed to be SIMD-friendly, explicit vectorization using AVX-512 or NEON could yield additional performance gains for batch conversion workloads.
3. **Compiler compatibility:** Further optimization for different compilers (particularly MSVC) would improve portability across development environments.
4. **Memory-constrained environments:** Investigating memory-efficient variants of xjb could benefit embedded systems and other resource-constrained platforms.

### 5.5. Availability

The complete implementation of the xjb algorithm, along with benchmark tools and test suites, is publicly available at <https://github.com/xjb714/xjb/releases/tag/v1.0.0>. We encourage the community to integrate, test, and contribute to the ongoing development of this work.

## 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. G. L. Steel Jr. and J. L. White. How to Print Floating-Point Numbers Accurately (Retrospective). *ACM SIGPLAN Notices* 39(4), April 2004 (Best of PLDI, 1979-1999). <https://dl.acm.org/doi/10.1145/989393.989431>
3. Robert G. Burger and R. Kent Dybvig. 1996. Printing Floating-point Numbers Quickly and Accurately. In *Proceedings of the ACM SIGPLAN1996 Conference on Programming Language Design and Implementation (PLDI '96)*. ACM, New York, NY, USA, 108&116. <https://doi.org/10.1145/231379.231397>
4. F. Loitsch. Printing Floating-Point Numbers Quickly and Accurately with Integers. In *Proceedings of the ACM SIGPLAN 2010 Conference on Programming Language Design and Implementation, PLDI 2010*. ACM, New York, NY, USA, 233-243. <https://doi.org/10.1145/1806596.1806623>
5. M. Andryscio, R. Jhala, and S. Lerner. Printing Floating-Point Numbers: a Faster, Always Correct Method. In *Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2016*. ACM, New York, NY, USA, 555-567. <https://doi.org/10.1145/2837614.2837654>
6. 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>
7. 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>
8. R. Giulietti. The Schubfach Way to Render Doubles. 2020. [https://drive.google.com/file/d/1KLtG\\_LalBk9ETXI290zqCvBW94dj058/view](https://drive.google.com/file/d/1KLtG_LalBk9ETXI290zqCvBW94dj058/view) (Sep. 2020)
9. J. Jeon. Grisu-Exact: A Fast and Exact Floating-Point Printing Algorithm. 2020. [https://github.com/jk-jeon/Grisu-Exact/blob/master/other\\_files/Grisu-Exact.pdf](https://github.com/jk-jeon/Grisu-Exact/blob/master/other_files/Grisu-Exact.pdf). (Sep. 2020)
10. Junekey Jeon. 2024. Dragonbox: A New Floating-Point Binary-to-Decimal Conversion Algorithm. <https://github.com/jk-jeon/Dragonbox>
11. Guo YaoYuan. [https://github.com/ibireme/c\\_numconv\\_benchmark/blob/master/vendor/yy\\_double/yy\\_double.c](https://github.com/ibireme/c_numconv_benchmark/blob/master/vendor/yy_double/yy_double.c) (Nov. 2024)
12. Russ Cox. <https://github.com/rsc/fpfmt> (Jan. 2026)
13. Russ Cox. Floating-Point Printing and Parsing Can Be Simple And Fast. published online, <https://research.swtch.com/fp>(Jan. 2026)
14. Russ Cox. Fast Unrounded Scaling: Proof by Ivy. published online, <https://research.swtch.com/fp-proof>(Jan. 2026)
15. Victor Zverovich. <https://github.com/vitaut/zmij> (Mar. 2026)
16. IEEE Standard for Binary Floating-Point Arithmetic, in ANSI/IEEE Std 754-1985 , vol., no., pp.1-20, 12 Oct. 1985. <https://doi.org/10.1109/IEEESTD.1985.82928>
17. IEEE Standard for Floating-Point Arithmetic, in IEEE Std 754-2019 (Revision of IEEE 754-2008) , vol., no., pp.1-84, 22 July 2019, <https://doi.org/10.1109/IEEESTD.2019.8766229>
18. Paul Khuong. How to print integers really fast (with Open Source AppNexus code!).published online, <https://pvk.ca/Blog/2017/12/22/appnexus-common-framework-its-out-also-how-to-print-integers-faster/>.(Dec. 2017)
19. Dougall Johnson. Converting integers to fixed-width strings faster with Neon SIMD on the Apple M1, published online, <https://dougallj.wordpress.com/2022/04/01/converting-integers-to-fixed-width-strings-faster-with-neon-simd-on-the-apple-m1/> (Apr. 2022)
20. Wojciech Muła. SSE: conversion integers to decimal representation. published online, <http://0x80.pl/notesen/2011-10-21-sse-itoa.html>. (Oct. 2011)
21. Daniel Lemire. Converting integers to decimal strings faster with AVX-512, published online, <https://lemire.me/blog/2022/03/28/converting-integers-to-decimal-strings-faster-with-avx-512/> (Mar. 2022)
22. Victor Zverovich. <https://github.com/fmtlib/fmt> (Oct. 2025)
23. Guo YaoYuan. <https://github.com/ibireme/yyjson> (Aug. 2025)
24. Cassio Neri. [https://github.com/cassioneri/teju\\_jagua](https://github.com/cassioneri/teju_jagua) (Nov. 2025)
25. Jing Leng. <https://github.com/lengjingzju/json/jnum.c> (Nov. 2025)

**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.