

# SLEEF: A Portable Vectorized Library of C Standard Mathematical Functions

Naoki Shibata, *Member, IEEE*, and Francesco Petrogalli

**Abstract**—In this paper, we present techniques used to implement our portable vectorized library of C standard mathematical functions written entirely in C language. In order to make the library portable while maintaining good performance, intrinsic functions of vector extensions are abstracted by inline functions or preprocessor macros. We implemented the functions so that they can use sub-features of vector extensions such as fused multiply-add, mask registers and extraction of mantissa. In order to make computation with SIMD instructions efficient, the library only uses a small number of conditional branches, and all the computation paths are vectorized. We devised a variation of the Payne-Hanek argument reduction for trigonometric functions and a floating point remainder, both of which are suitable for vector computation. We compare the performance of our library to Intel SVML.

**Index Terms**—Parallel and vector implementations, SIMD processors, elementary functions, floating-point arithmetic

## 1 INTRODUCTION

THE instruction set architecture of most modern processors provides *Single Instruction Multiple Data (SIMD)* instructions that process multiple instances of data concurrently [1]. The programming model that utilizes these instructions is a key technique for many computing systems to reach their peak performance. Most software SIMD optimizations are introduced manually by programmers. However, this approach introduces a portability problem because the code needs to be re-written when targeting a new vector extension. In order to improve portability of codes with SIMD optimizations, recent compilers have introduced auto-vectorizing capability [2]. To fully exploit the SIMD capabilities of a system, the transformation for auto-vectorization of a compiler must be able to invoke a version of functions that operates on concurrent iterations, or on a *vector function*. This applies particularly to C mathematical functions defined in `math.h` that are frequently called in hot-loops.

In this paper, we describe our implementation of a vectorized library of C standard math functions, called SLEEF library. SLEEF stands for *SIMD Library for Evaluating Elementary Functions*, and implements a vectorized version of all C99 real floating-point math functions. Our library provides 1-ULP accuracy version and 3.5-ULP accuracy version for most of the functions. We confirmed that our library satisfies such accuracy requirements on an empirical basis. Our library achieves both good performance and portability by abstracting intrinsic functions. This abstraction enables sub-features of vector extensions such as mask registers to

be utilized while the source code of our library is shared among different vector extensions. We also implemented a version of functions that returns bit-wise consistent results across all platforms. Our library is designed to be used in conjunction with vectorizing compilers. In order to help development of vectorizing compilers, we collaborated with compiler developers in designing a Vector Function Application Binary Interface (ABI). The main difficulty in vectorizing math functions is that conditional branches are expensive. We implemented many of the functions in our library without conditional branches. We devised reduction methods and adjusted domains of polynomials so that a single polynomial covers the entire input domain. For an increased vector size, a value requiring a slow path is more likely to be contained in a vector. Therefore, we vectorized all the code paths in order to speed up the computation in such cases. We devised a variation of the Payne-Hanek range reduction and a remainder calculation method that are both suitable for vectorized implementation.

We compare the implementation of several selected functions in our library to those in other open-source libraries. We also compare the reciprocal throughput of functions in our library, Intel SVML [3], FDLIBM [4], and Vector-libm [5]. We show that the performance of our library is comparable to that of Intel SVML.

The rest of this paper is organized as follows. Section 2 introduces related work. Section 3 discusses how portability is improved by abstracting vector extensions. Section 4 explains the development of a Vector ABI and a vectorized mathematical library. Section 5 shows an overview of the implementation of SLEEF, while comparing our library with FDLIBM and Vector-libm. Section 6 explains how our library is tested. Section 7 compares our work with prior art. In Section 8, the conclusions are presented.

• Naoki Shibata is with Graduate School of Information Science, Nara Institute of Science and Technology, Nara, Japan.

• Francesco Petrogalli is with ARM, 110 Fulbourn Road, Cambridge, CB1 9NJ, United Kingdom.

Manuscript received 25 Apr. 2019; revised 11 Dec. 2019; accepted 13 Dec. 2019. Date of publication 18 Dec. 2019.

(Corresponding author: Naoki Shibata.)

Recommended for acceptance by D. S. Nikolopoulos.

Digital Object Identifier no. 10.1109/TPDS.2019.2960333

## 2 RELATED WORK

### 2.1 C Standard Math Library

The C standard library (`libc`) includes the standard mathematical library (`libm`) [6]. There have been many implementations of `libm`. Among them, `FDLIBM` [4] and the `libm` included in the GNU C Library [7] are the most widely used libraries. `FDLIBM` is a freely distributable `libm` developed by Sun Microsystems, Inc., and there are many derivations of this library. Gal et al. described the algorithms used in the elementary mathematical library of the IBM Israel Scientific Center [8]. Their algorithms are based on the accurate tables method developed by Gal. It achieves high performance and produces very accurate results. `Crlibm` is a project to build a correctly rounded mathematical library [9].

There are several existing vectorized implementations of `libm`. Intel Short Vector Math Library (SVML) is a highly regarded commercial library [3]. This library provides highly optimized subroutines for evaluating elementary functions which can use several kinds of vector extensions available in Intel's processors. However, this library is proprietary and only optimized for Intel's processors. There are also a few commercial and open-source implementations of vectorized `libm`. AMD is providing a vectorized `libm` called AMD Core Math Library (ACML) [10].

Some of the code from SVML is published under a free software license, and it is now published as `Libmvec` [11], which is a part of `Glibc`. This library provides functions with 4-ULP error bound. It is coded in assembly language, and therefore it does not have good portability. C. K. Anand et al. reported their C implementation of 32 single precision `libm` functions tuned for the Cell BE SPU compute engine [12]. They used an environment called `Coconut` that enables rapid prototyping of patterns, rapid unit testing of assembly language fragments and patterns to develop their library. M. Dukhan published an open-source and portable SIMD vector `libm` library named `Yeppp!` [13], [14]. Most of vectorized implementations of `libm` utilizes assembly coding or intrinsic functions to specify which vector instruction is used for each operator. On the other hand, there are also other implementations of vector versions of `libm` which are written in a scalar fashion but rely on a vectorizing compiler to generate vector instructions and generate a vectorized binary code. Christoph Lauter published an open-source Vector-`libm` library implemented with plain C [5]. VDT Mathematical Library [15], is a math library written for the compiler's auto-vectorization feature.

### 2.2 Translation of SIMD Instructions

Manilov et al. propose a C source code translator for substituting calls to platform-specific intrinsic functions in a source code with those available on the target machine [16]. This technique utilizes graph-based pattern matching to substitute intrinsics. It can translate SIMD intrinsics between extensions with different vector lengths. This rewriting is carried out through loop-unrolling.

N. Gross proposes specialized C++ templates for making the source code easily portable among different vector extensions without sacrificing performance [17]. With these templates, some part of the source code can be written in a way that resembles scalar code. In order to vectorize

algorithms that have a lot of control flow, this scheme requires the bucketing technique is applied, to compute all the paths and choose the relevant results at the end.

Clark et al. proposes a method for combining static analysis at compile time and binary translation with a JIT compiler in order to translate SIMD instructions into those that are available on the target machine [18]. In this method, SIMD instructions in the code are first converted into an equivalent scalar representation. Then, a dynamic translation phase turns the scalar representation back into architecture-specific SIMD equivalents.

Leißa et al. propose a C-like language for portable and efficient SIMD programming [19]. With their extension, writing vectorized code is almost as easy as writing traditional scalar code. There is no strict separation in host code and kernels, and scalar and vector programming can be mixed. Switching between them is triggered by the type system. The authors present a formal semantics of their extension and prove the soundness of the type system.

Most of the existing methods are aiming at translating SIMD intrinsics or instructions to those provided by a different vector extension in order to port a code. Intrinsics that are unique in a specific extension are not easy to handle, and translation works only if the source and the target architectures have equivalent SIMD instructions. Automatic vectorizers in compilers have a similar weakness. Whenever possible, we have specialized the implementation of the math functions to exploit the SIMD instructions that are specific to a target vector extension. We also want to make special handling of FMA, rounding and a few other kinds of instructions, because these are critical for both execution speed and accuracy. We want to implement a library that is statically optimized and usable with Link Time Optimization (LTO). The users of our library do not appreciate usage of a JIT compiler. In order to minimize dependency on external libraries, we want to write our library in C. In order to fulfill these requirements, we take a cross-layer approach. We have been developing our abstraction layer of intrinsics, the library implementation, and the algorithms in order to make our library run fast with any vector extensions.

## 3 ABSTRACTION OF VECTOR EXTENSIONS

Modern processors supporting SIMD instructions have SIMD registers that can contain multiple data [1]. For example, a 128-bit wide SIMD register may contain four 32-bit single-precision FP numbers. A SIMD add instruction might take two of these registers as operands, add the four pairs of numbers, and overwrite one of these registers with the resulting four numbers. We call an array of FP numbers contained in a SIMD register a vector.

SIMD registers and instruction can be exposed in a C program with intrinsic functions and types [20]. An intrinsic function is a kind of inline function that exposes the architectural features of an instruction set at C level. By calling an intrinsic function, a programmer can make a compiler generate a specific instruction without hand-coded assembly. Nevertheless, the compiler can reorder instructions and allocate registers, and therefore optimize the code. When intrinsic functions corresponding to SIMD instructions are

```

1 // Type definition
2 typedef svbool_t vopmask;
3 typedef svfloat64_t vdouble;
4
5 // Abstraction of intrinsic functions
6 #define ptrue svptrue_b8()
7 vdouble vcast_vd_d(double d) { return
    svdup_n_f64(d); }
8 vdouble vsub_vd_vd(vdouble x, vdouble y) {
9     return svsub_f64_x(ptrue, x, y); }
10 vopmask veq_vo_vd_vd(vdouble x, vdouble y) {
11     return svcmpeq_f64(ptrue, x, y); }
12 vopmask vlt_vo_vd_vd(vdouble x, vdouble y) {
13     return svcmplt_f64(ptrue, x, y); }
14 vopmask vor_vo_vo_vo(vopmask x, vopmask y) {
15     return svorr_b_z(ptrue, x, y); }
16 vdouble vsel_vd_vo_vd_vd(vopmask mask, vdouble x,
    vdouble y) {
17     return svsel_f64(o, x, y); }

```

Fig. 1: A part of definitions in VEAL for SVE

```

1 vdouble xfdim(vdouble x, vdouble y) {
2     vdouble ret = vsub_vd_vd(x, y);
3     ret =
        vsel_vd_vo_vd_vd(vor_vo_vo_vo(vlt_vo_vd_vd(ret,
        vcast_vd_d(0)), veq_vo_vd_vd(x, y)),
        vcast_vd_d(0),
        ret);
4     return ret;
5 }
6 }
7 }

```

Fig. 2: Implementation of vectorized `fdim` (positive difference) function with VEAL

defined inside a compiler, C data types for representing vectors are also defined.

In SLEEF, we use intrinsic functions to specify which assembly instruction to use for each operator. We abstract intrinsic functions for each vector extension by a set of inline functions or preprocessor macros. We implement the functions exported from the library to call abstract intrinsic functions instead of directly calling intrinsic functions. In this way, it is easy to swap the vector extension to use. We call our set of inline functions for abstracting architecture-specific intrinsics *Vector Extension Abstraction Layer (VEAL)*.

In some of the existing vector math libraries, functions are implemented with hand-coded assembly [11]. This approach improves the absolute performance because it is possible to provide the optimal implementation for each microarchitecture. However, processors with a new microarchitecture are released every few years, and the library needs revision accordingly in order to maintain the optimal performance.

In other vector math libraries, the source code is written in a scalar fashion that is easy for compilers to auto-vectorize [5], [15]. Although such libraries have good portability, it is not easy for compilers to generate a well-optimized code. In order for each transformation rule in an optimizer to kick in, the source code must satisfy many conditions to guarantee that the optimized code runs correctly and faster. In order to control the level of optimization, a programmer must specify special attributes and compiler options.

### 3.1 Using Sub-features of the Vector Extensions

There are differences in the features provided by different vector extensions, and we must change the function implementation according to the available features. Thanks to the level of abstraction provided by the VEALs, we implemented the functions so that all the different versions of functions can be built from the same source files with different macros enabled. For example, the availability of FMA instructions is important when implementing double-double (DD) operators [21]. We implemented DD operators both with and without FMA by manually specifying if the compiler can convert each combination of multiplication and addition instructions to an FMA instruction, utilizing VEALs.

Generally, bit masks are used in a vectorized code in order to conditionally choose elements from two vectors. In some vector extensions, a vector register with a width that matches a vector register for storing FP values, is used to store a bit mask. Some vector extensions provide narrower vector registers that are dedicated to this purpose, which is SLEEF makes use of these opmask registers by providing a dedicated data type in VEALs. If a vector extension does not support an opmask, the usual bit mask is used instead of an opmask. It is also better to have an opmask as an argument of a whole math function and make that function only compute the elements specified by the opmask. By utilizing a VEAL, it is also easy to implement such a functionality.

### 3.2 Details of VEALs

Fig. 1 shows some definitions in the VEAL for SVE [22]. We abstract vector data types and intrinsic functions with typedef statements and inline functions, respectively.

The `vdouble` data type is for storing vectors of double precision FP numbers. `vopmask` is the data type for the opmask described in 3.1.

The function `vcast_vd_d` is a function that returns a vector in which the given scalar value is copied to all elements in the vector. `vsub_vd_vd_vd` is a function for vector subtraction between two `vdouble` data. `veq_vo_vd_vd` compares elements of two vectors of `vdouble` type. The results of the comparison can be used, for example, by `vsel_vd_vo_vd_vd` to choose a value for each element between two vector registers. Fig. 2 shows an implementation of a vectorized positive difference function using a VEAL. This function is a vectorized implementation of the `fdim` function in the C standard math library.

### 3.3 Making Results Bit-wise Consistent across All Platforms

The method of implementing math functions described so far, can deliver computation results that slightly differ depending on architectures and other conditions, although they all satisfy the accuracy requirements, and other specifications. However, in some applications, bit-wise consistent results are required.

To this extent, the SLEEF project has been working closely with Unity Technologies,<sup>1</sup> which specializes in developing frameworks for video gaming, and we discovered

1. <https://unity3d.com/>.

that they have unique requirements for the functionalities of math libraries. Networked video games run on many gaming consoles with different architectures and they share the same virtual environment. Consistent results of simulation at each terminal and server are required to ensure fairness among all players. For this purpose, fast computation is more important than accurate computation, while the results of computation have to perfectly agree between many computing nodes, which are not guaranteed to rely on the same architecture. Usually, fixed-point arithmetic is used for a purpose like this, however there is a demand for modifying existing codes with FP computation to support networking.

There are also other kinds of simulation in which bit-wise identical reproducibility is important. In [23], the authors show that modeled mean climate states, variability and trends at different scales may be significantly changed or even lead to opposing results due to the round-off errors in climate system simulations. Since reproducibility is a fundamental principle of scientific research, they propose to promote bit-wise identical reproducibility as a worldwide standard.

One way to obtain bit-wise consistent values from math functions is to compute correctly rounded values. However, for applications like networked video games, this might be too expensive. SLEEF provides vectorized math functions that return bit-wise consistent results across all platforms and other settings, and this is also achieved by utilizing VEALs. The basic idea is to always apply the same sequence of operations to the arguments. The IEEE 754 standard guarantees that the basic arithmetic operators give correctly rounded results [24], and therefore the results from these operators are bit-wise consistent. Because most of the functions except trigonometric functions do not have a conditional branch in our library, producing bit-wise consistent results is fairly straightforward with VEALs. Availability of FMA instructions is another key for making results bit-wise consistent. Since FMA instructions are critical for performance, we cannot just give up using FMA instructions. In SLEEF, the bit-wise consistent versions of functions have two versions both with and without FMA instructions. We provide a non-FMA version of the functions to guarantee bit-wise consistency among extensions such as Intel SSE2 that do not have FMA instructions. Another issue is that the compiler might introduce inconsistency by FP contraction, which is the result of combining a pair of multiplication and addition operations into an FMA. By disabling FP contraction, the compiler strictly preserves the order and the type of FP operations during optimization. It is also important to make the returned values from scalar functions bit-wise consistent with the vector functions. In order to achieve this, we also made a VEAL that only uses scalar operators and data types. The bit-wise consistent and non-consistent versions of vector and scalar functions are all built from the same source files, with different VEALs and macros enabled. As described in Section 5, trigonometric functions in SLEEF chooses a reduction method according to the maximum argument of all elements in the argument vector. In order to make the returned value bit-wise consistent, the bit-wise consistent version of the functions first applies the reduction method for small arguments to the elements covered by

this method. Then it applies the second method only to the elements with larger arguments which the first method does not cover.

## 4 THE DEVELOPMENT OF A VECTOR FUNCTION ABI AND SLEEF

Recent compilers are developing new optimization techniques to automatically vectorize a code written in standard programming languages that do not support parallelization [25], [26]. Although the first SIMD and vector computing systems [27] appeared a few decades ago, compilers with auto-vectorization capability have not been widely used until recently, because of several difficulties in implementing such functionality for modern SIMD architectures. Such difficulties include verifying whether the compiler can vectorize a loop or not, by determining data access patterns of the operations in the loop [2], [28]. For languages like C and C++, it is also difficult to determine the data dependencies through the iteration space of the loop, because it is hard to determine aliasing conditions of the arrays processed in the loop.

### 4.1 Vector Function Application Binary Interface

Vectorizing compilers convert calls to scalar versions of math functions such as sine and exponential to the SIMD version of the math functions. The most recent versions of Intel Compiler [29], GNU Compiler [30], and Arm Compiler for HPC [31], which is based on Clang/LLVM [32], [33], are capable of this transformation, and rely on the availability of vector math libraries such as SVML [3], Libmvec [11] and SLEEF respectively to provide an implementation of the vector function calls that they generate. In order to develop this kind of transformations, a target-dependent *Application Binary Interface (ABI)* for calling vectorized functions had to be designed.

The Vector Function ABI for AArch64 architecture [34] was designed in close relationship with the development of SLEEF. This type of ABI must standardize the mapping between scalar functions and vector functions. The existence of a standard enables interoperability across different compilers, linkers and libraries, thanks to the use of standard names defined by the specification.

The ABI includes a name mangling function, a map that converts the scalar signature to the vector one, and the calling conventions that the vector functions must obey. In particular, the name mangling function that takes the name of the scalar function to the vector function must encode all the information that is necessary to reverse the transformation back to the original scalar function. A linker can use this reverse mapping to enable more optimizations (Link Time Optimizations) that operate on object files, and does not have access to the scalar and vector function prototypes. There is a demand by users for using a different vector math library according to the usage. Reverse mapping is also handy for this purpose. A vector math library implements a function for each combination of a vector extension, a vector length and a math function to evaluate. As a result, the library exports a large number of functions. Some vector math libraries can only implement

part of all the combinations. By using the reverse mapping mechanism, the compiler can check the availability of the functions by scanning the symbols exported by a library.

The Vector Function ABI is also used with OpenMP [35]. From version 4.0 onwards, OpenMP provides the directive `declare simd`. A user can decorate a function with this directive to inform the compiler that the function can be safely invoked concurrently on multiple instances of its arguments [36]. This means that the compiler can vectorize the function safely. This is particularly useful when the function is provided via a separate module, or an external library, for example in situations where the compiler is not able to examine the behavior of the function in the call site. The scalar-to-vector function mapping rules stipulated in the Vector Function ABI are based on the classification of vector functions associated with the `declare simd` directive of OpenMP. Currently, work for implementing these OpenMP directives on LLVM is ongoing.

The Vector Function ABI specifications are provided for the Intel x86 and the Armv8 (AArch64) families of vector extensions [34], [37]. The compiler generates SIMD function calls according to the compiler flags. For example, when targeting AArch64 SVE auto-vectorization, the compiler will transform a call to the standard `sin` function to a call to the symbol `_ZGVsMxv_sin`. When targeting Intel AVX-512 [38] auto-vectorization, the compiler would generate a call to the symbol `_ZGVeNe8v_sin`.

## 4.2 SLEEF and the Vector Function ABI

SLEEF is provided as two separate libraries. The first library exposes the functions of SLEEF to programmers for inclusion in their C/C++ code. The second library exposes the functions with names mangled according to the Vector Function ABI. This makes SLEEF a viable alternative to `libm` and its SIMD counterpart `libmvec`, in `glibc`. This also enables a user work-flow that relies on the auto-vectorization capabilities of a compiler. The compatibility with `libmvec` enables users to swap from `libmvec` to `libsleef` by simply changing compiler options, without changing the code that generated the vector call. The two SLEEF libraries are built from the same source code, which are configured to target the different versions via auto-generative programs that transparently rename the functions according to the rules of the target library.

## 5 OVERVIEW OF LIBRARY IMPLEMENTATION

One of the objectives of the SLEEF project is to provide a library of vectorized math functions that can be used in conjunction with vectorizing compilers. When a non-vectorized code is automatically vectorized, the compiler converts calls to scalar math functions to calls to a SIMD version of the math functions. In order to make this conversion safe and applicable to wide variety of codes, we need functions with 1-ULP error bound that conforms to ANSI C standard. On the other hand, there are users who need better performance. Our library provides 1-ULP accuracy version and 3.5-ULP accuracy version for most of the functions. We confirmed that our library satisfies the accuracy requirements on an empirical basis. For non-finite inputs

and outputs, we implemented the functions to return the same results as `libm`, as specified in the ANSI C standard. They do not set `errno` nor raise an exception.

In order to optimize a program with SIMD instructions, it is important to eliminate conditional branches as much as possible, and execute the same sequence of instructions regardless of the argument. If the algorithm requires conditional branches according to the argument, it must prepare for the case where the elements in the input vector contain both values that would make a branch happen and not happen. Recent processors have a long pipeline and therefore branch misprediction penalty can reach more than 10 cycles [39]. Making a decision for a conditional branch also requires non-negligible computation, within the scope of our tests. A conditional move is an operator for choosing one value from two given values according to a condition. This is equivalent to a ternary operator and can be used in a vectorized code to replace a conditional branch. Some other operations are also expensive in vectorized implementation. A table-lookup is expensive. Although in-register table lookup is reported fast on Cell BE SPU [12], it is substantially slower than polynomial evaluation without any table lookup, within the scope of our tests. Most vector extensions do not provide 64-bit integer multiplication or a vector shift operator with which each element of a vector can be specified a different number of bits to shift. On the other hand, FMA and round-to-integer instructions are supported by most vector extensions. Due to the nature of the evaluation methods, dependency between operations cannot be completely eliminated. Latencies of operations become an issue when a series of dependent operations are executed. FP division and square root are not too expensive from this aspect.<sup>2</sup>

The actual structure of the pipeline in a processor is complex, and such level of details are not well-documented for most CPUs. Therefore, it is not easy to optimize the code according to such hardware implementation. In this paper, we define the latency and throughput of an instruction or a subroutine as follows [41]. The latency of an instruction or a subroutine is the delay that it generates in a dependency chain. The throughput is the maximum number of instructions or subroutines of the same kind that can be executed per unit time when the inputs are independent of the preceding instructions or subroutines. Several tools and methods are proposed for automatically constructing models of latency, throughput, and port usage of instructions [42], [43]. Within the scope of our tests, most of the instruction latency in the critical path of evaluating a vector math function tends to be dominated by FMA operations. In many processors, FMA units are implemented in a pipeline manner. Some powerful processors have multiple FMA units with out-of-order execution, and thus the throughput of FMA instruction is large, while the latency is long. In SLEEF, we try to maximize the throughput of computation in a versatile way by only taking account of dependencies among FMA operations. We regard each FMA operation as a job that can be executed in parallel and try to reduce the length of the critical path.

2. The latencies of 256-bit DP add, divide and sqrt instructions are 4, 14 and 18 cycles, respectively on Intel Skylake processors [40].

In order to evaluate a double-precision (DP) function to 1-ULP accuracy, the internal computation with accuracy better than 1 ULP is sometimes required. *Double-double* (DD) arithmetic, in which a single value is expressed by a sum of two double-precision FP values [44], [45], is used for this purpose. All the basic operators for DD arithmetic can be implemented without a conditional branch, and therefore it is suitable for vectorized implementation. Because we only need 1-ULP overall accuracy for DP functions, we use simplified DD operators with less than the full DD accuracy. In SLEEF, we omit re-normalization of DD values by default, allowing overlap between the two numbers. We carry out re-normalization only when necessary.

Evaluation of an elementary function often consists of three steps: range reduction, approximation, and reconstruction [21]. An approximation step computes the elementary function using a polynomial. Since this approximation is only valid for a small domain, a number within that range is computed from the argument in a range reduction step. The reconstruction step combines the results of the first two steps to obtain the resulting number.

An argument reduction method that finds an FP remainder of dividing the argument  $x$  by  $\pi$  is used in evaluation of trigonometric functions. The range reduction method suggested by Cody and Waite [46], [47] is used for small arguments. The Payne and Hanek's method [48] provides an accurate range-reduction for a large argument of trigonometric function, but it is expensive in terms of operations.

There are tools available for generating the coefficients of the polynomials, such as Maple [49] and Sollya [50]. In order to fine-tune the generated coefficients, we created a tool for generating coefficients that minimizes the maximum relative error. When a SLEEF function evaluates a polynomial, it evaluates a few lowest degree terms in DD precision while other terms are computed in double-precision, in order to achieve 1-ULP overall accuracy. Accordingly, coefficients in DD precision or coefficients that can be represented by FP numbers with a few most significant bits in mantissa are used in the last few terms. We designed our tool to generate such coefficients. We use Estrin's scheme [51] to evaluate a polynomial to reduce dependency between FMA operations. This scheme reduces bubbles in the pipeline, and allows more FMA operations to be executed in parallel. Reducing latency can improve the throughput of evaluating a function because the latency and the reciprocal throughput of the entire function are close to each other.

Below, we describe and compare the implementations of selected functions in SLEEF, FDLIBM [4] and Christoph Lauter's Vector-libm [5]. We describe 1-ULP accuracy version of functions in SLEEF. The error bound specification of FDLIBM is 1 ULP.

### 5.1 Implementation of $\sin$ and $\cos$

FDLIBM uses Cody-Waite range reduction if the argument is under  $2^{18}\pi$ . Otherwise, it uses the Payne-Hanek range reduction. Then, it switches between polynomial approximations of the sine and cosine functions on  $[-\pi/4, \pi/4]$ . Each polynomial has 6 non-zero terms.

$\sin$  and  $\cos$  in Vector-libm have 4-ULP error bound. They use a vectorized path if all arguments are greater

than  $3.05e-151$ , and less than 5.147 for sine and 2.574 for cosine. In the vectorized paths, a polynomial with 8 and 9 non-zero terms is used to approximate the sine function on  $[-\pi/2, \pi/2]$ , following Cody-Waite range reduction. In the scalar paths, Vector-libm uses a polynomial with 10 non-zero terms.

SLEEF switches among two Cody-Waite range reduction methods with approximation with different sets of constants, and the Payne-Hanek reduction. The first version of the algorithm operates for arguments within  $[-15, 15]$ , and the second version for arguments that are within  $[-10^{14}, 10^{14}]$ . Otherwise, SLEEF uses a vectorized Payne-Hanek reduction, which is described in 5.8. SLEEF only uses conditional branches for choosing a reduction method from Cody-Waite and Payne-Hanek. SLEEF uses a polynomial approximation of the sine function on  $[-\pi/2, \pi/2]$ , which has 9 non-zero terms. The sign is set in the reconstruction step.

### 5.2 Implementation of $\tan$

After Cody-Waite or Payne-Hanek reduction, FDLIBM reduces the argument to  $[0, 0.67434]$ , and uses a polynomial approximation with 13 non-zero terms. It has 10 `if` statements after Cody-Waite reduction.

$\tan$  in Vector-libm has 8-ULP error bound. A vectorized path is used if all arguments are less than 2.574 and greater than  $3.05e-151$ . After Cody-Waite range reduction, a polynomial with 9 non-zero terms for approximating sine function on  $[-\pi/2, \pi/2]$  is used twice to approximate sine and cosine of the reduced argument. The result is obtained by dividing these values. In the scalar path, Vector-libm evaluates a polynomial with 10 non-zero terms twice.

In SLEEF, the argument is reduced in 3 levels. It first reduces the argument to  $[-\pi/2, \pi/2]$  with Cody-Waite or Payne-Hanek range reduction. Then, it reduces the argument to  $[-\pi/4, \pi/4]$  with  $\tan a_1 = 1/\tan(\pi/2 - a_0)$ . At the third level, it reduces the argument to  $[-\pi/8, \pi/8]$  with the double-angle formula. Let  $a_0$  be the reduced argument with Cody-Waite or Payne-Hanek.  $a_1 = \pi/2 - a_0$  if  $|a_0| > \pi/4$ . Otherwise,  $a_1 = a_0$ . Then, SLEEF uses a polynomial approximation of the tangent function on  $[-\pi/8, \pi/8]$ , which has 9 non-zero terms, to approximate  $\tan(a_1/2)$ . Let  $t$  be the obtained value with this approximation. Then,  $\tan a_0 \approx 2t/(1-t^2)$  if  $|a_0| \leq \pi/4$ . Otherwise,  $\tan a_0 \approx (1-t^2)/(2t)$ . SLEEF only uses conditional branches for choosing a reduction method from Cody-Waite and Payne-Hanek. Annotated source code of  $\tan$  is shown in Appendix A

### 5.3 Implementation of $\text{asin}$ and $\text{acos}$

FDLIBM and SLEEF first reduces the argument to  $[0, 0.5]$  using  $\arcsin x = \pi/2 - 2 \arcsin \sqrt{(1-x)/2}$  and  $\arccos x = 2 \arcsin \sqrt{(1-x)/2}$ .

Then, SLEEF uses a polynomial approximation of arcsine on  $[0, 0.5]$  with 12 non-zero terms.

FDLIBM uses a rational approximation with 11 terms (plus one division). For computing arcsine, FDLIBM switches the approximation method if the original argument is over 0.975. For computing arccosine, it has three paths that are taken when  $|x| < 0.5$ ,  $x \leq -0.5$  and  $x \geq 0.5$ ,

respectively. It has 7 and 6 `if` statements in `asin` and `acos`, respectively.

`asin` and `acos` in Vector-libm have 6-ULP error bound. `asin` and `acos` in Vector-libm use vectorized paths if arguments are all greater than  $3.05e-151$  and  $2.77e-17$ , respectively. Vector-libm evaluates polynomials with 3, 8, 8, and 5 terms to compute arcsine. It evaluates a polynomial with 21 terms for arccosine.

## 5.4 Implementation of `atan`

FDLIBM reduces the argument to  $[0, 7/16]$ . It uses a polynomial approximation of the arctangent function with 11 non-zero terms. It has 9 `if` statements.

`atan` in Vector-libm have 6-ULP error bound. Vector-libm uses vectorized paths if arguments are all greater than  $1.86e-151$  and less than 2853. It evaluates four polynomials with 7, 9, 9 and 4 terms in the vectorized path.

SLEEF reduces argument  $a$  to  $[0, 1]$  using  $\arctan x = \pi/2 - \arctan(1/x)$ . Let  $a' = 1/a$  if  $|a| \geq 1$ . Otherwise,  $a' = a$ . It then uses a polynomial approximation of arctangent function with 20 non-zero terms to approximate  $r \approx \arctan a'$ . As a reconstruction, it computes  $\arctan a \approx \pi/2 - r$  if  $|a| \geq 1$ . Otherwise,  $\arctan a \approx r$ .

## 5.5 Implementation of `log`

FDLIBM reduces the argument to  $[\sqrt{2}/2, \sqrt{2}]$ . It then approximates the reduced argument with a polynomial that contains 7 non-zero terms in a similar way to SLEEF. It has 9 `if` statements.

`log` in Vector-libm has 4-ULP error bound. It uses a vectorized path if the input is a normalized number. It uses a polynomial with 20 non-zero terms to approximate the logarithm function on  $[0.75, 1.5]$ . It does not use division.

SLEEF multiplies the argument  $a$  by  $2^{64}$ , if the argument is a denormal number. Let  $a'$  be the resulting argument,  $e = \lfloor \log_2(4a'/3) \rfloor$  and  $m = a' \cdot 2^{-e}$ . If  $a$  is a denormal number,  $e$  is subtracted 64. SLEEF uses a polynomial with 7 non-zero terms to evaluate  $\log m \approx \sum_{n=0}^6 C_n \left(\frac{m-1}{m+1}\right)^{2n+1}$ , where  $C_0 \dots C_6$  are constants. As a reconstruction, it computes  $\log a = e \log 2 + \log m$ .

## 5.6 Implementation of `exp`

All libraries reduce the argument range to  $[-(\log 2)/2, (\log 2)/2]$  by finding  $r$  and integer  $k$  such that  $x = k \log 2 + r$ ,  $|r| \leq (\log 2)/2$ .

SLEEF then uses a polynomial approximation with 13 non-zero terms to directly approximate the exponential function of this domain. It achieves 1-ULP error bound without using a DD operation.

FDLIBM uses a polynomial with 5 non-zero terms to approximate  $f(r) = r(e^r + 1)/(e^r - 1)$ . It then computes  $\exp(r) = 1 + 2r/(f(r) - r)$ . It has 11 `if` statements.

The reconstruction step is to add integer  $k$  to the exponent of the resulting FP number of the above computation.

`exp` in Vector-libm has 4-ULP error bound. A vectorized path covers almost all input domains. It uses a polynomial with 11 terms to approximate the exponential function.

## 5.7 Implementation of `pow`

FDLIBM computes  $y \log_2 x$  in DD precision. Then, it computes  $\text{pow}(x, y) = e^{\log 2 \cdot y \log_2 x}$ . It has 44 `if` statements.

Vector-libm does not implement `pow`.

SLEEF computes  $e^{y \log x}$ . The internal computation is carried out in DD precision. In order to compute logarithm internally, it uses a polynomial with 11 non-zero terms. The accuracy of the internal logarithm function is around 0.008 ULP. The internal exponential function in `pow` uses a polynomial with 13 non-zero terms.

## 5.8 The Payne-Hanek Range Reduction

Our method computes  $\text{rfrac}(2x/\pi) \cdot \pi/2$ , where  $\text{rfrac}(a) := a - \text{round}(a)$ . The argument  $x$  is an FP number, and therefore it can be represented as  $M \cdot 2^E$ , where  $M$  is an integer mantissa and  $E$  is an integer exponent  $E$ . We now denote the integral part and the fractional part of  $2^E \cdot 2/\pi$  as  $I(E)$  and  $F(E)$ , respectively. Then,

$$\begin{aligned} \text{rfrac}(2x/\pi) &= \text{rfrac}(M \cdot 2^E \cdot 2/\pi) \\ &= \text{rfrac}(M \cdot (I(E) + F(E))) \\ &= \text{rfrac}(M \cdot F(E)). \end{aligned}$$

The value  $F(E)$  only depends on the exponent of the argument, and therefore, can be calculated and stored in a table, in advance. In order to compute  $\text{rfrac}(M \cdot F(E))$  in DD precision,  $F(E)$  must be in quad-double-precision. We now denote  $F(E) = F_0(E) + F_1(E) + F_2(E) + F_3(E)$ , where  $F_0(E) \dots F_3(E)$  are DP numbers and  $|F_0(E)| \geq |F_1(E)| \geq |F_2(E)| \geq |F_3(E)|$ . Then,

$$\begin{aligned} &\text{rfrac}(M \cdot F(E)) \\ &= \text{rfrac}(M \cdot F_0(E) + M \cdot F_1(E) \\ &\quad + M \cdot F_2(E) + M \cdot F_3(E)) \\ &= \text{rfrac}(\text{rfrac}(\text{rfrac}(\text{rfrac}(M \cdot F_0(E)) + M \cdot F_1(E)) \\ &\quad + M \cdot F_2(E)) + M \cdot F_3(E)), \end{aligned} \quad (1)$$

because  $\text{rfrac}(a + b) = \text{rfrac}(\text{rfrac}(a) + b)$ . In the method, we compute (1) in DD precision in order to avoid overflow. The size of the table retaining  $F_0(E) \dots F_3(E)$  is 32K bytes. Our method is included in the source code of `tan` shown in Appendix A.

FDLIBM seems to implement the original Payne-Hanek algorithm with more than 100 lines of C code, which includes 13 `if` statements, 18 `for` loops, 1 `switch` statement and 1 `goto` statement. The numbers of iterations of most of the `for` loops depend on the argument.

Vector-libm implements a non-vectorized variation of the Payne-Hanek algorithm which has some similarity with our method. In order to reduce argument  $x$ , it first decomposes  $|x|$  into  $E$  and  $n$  such that  $2^E \cdot n = |x|$ . A triple-double (TD) approximation to  $t(E) = 2^E/\pi - 2 \cdot \lfloor 2^{E-1}/\pi \rfloor$  is looked-up from a table. It then calculates  $m = n \cdot t(E)$  in TD. The reduced argument is obtained as a product of  $\pi$  and the fractional part of  $m$ . In Table 1, we compare the numbers of FP operators in the implementations. Note that the method in Vector-libm is used for trigonometric functions with 4-ULP error bound, while our method is used for functions with 1-ULP error bound.

$$\begin{array}{r}
0x1.397714bd79317p+133(\approx 1.333e+40) \quad 0x1.555555555553p+82(\approx 6.448e+24) \quad 0x1.aaaaaaaaaaaa8p+31(\approx 3.579e+10) \\
0.75 \quad 0x1.d6329f1c35ca5p+132(\approx 1.000e+40) \\
\hline
0x1.d6329f1c35ca1p+132(\approx 1.000e+40) \\
\hline
0x1p+82(\approx 4.836e+24) \\
0x1.fffffffffffffbp+81(\approx 4.836e+24) \\
\hline
0x1.4p+31(=2684354560) \\
0x1.3fffffffff8p+31(=2684354559.75) \\
0.25
\end{array}$$

Fig. 3: Example computation of FP remainder

TABLE 1: Number of FP operators in the Payne-Hanek implementations

Operator	SLEEF (1-ULP)	Vector-libm (4-ULP)
add/sub	36	71
mul	5	18
FMA	11	0
round	8	0

**ALGORITHM 1:** Exact remainder calculation**Input:** Finite positive numbers  $n$  and  $d$ **Output:** Returns  $n - d \lfloor n/d \rfloor$ 

- 1:  $r_0 := n, k := 0$
- 2: **while**  $d \leq r_k$  **do**
- 3:    $q_k$  is an arbitrary integer satisfying  $(r_k/d)/2 \leq q_k \leq r_k/d$
- 4:    $r_{k+1} := r_k - q_k d$
- 5:    $k := k + 1$
- 6: **end while**
- 7: **return**  $r_k$

**5.9 FP Remainder**

We devised an exact remainder calculation method suitable for vectorized implementation. The method is based on the long division method, where an FP number is regarded as a digit. Fig. 3 shows an example process for calculating the FP remainder of  $1e+40 / 0.75$ . Like a typical long division, we first find integer quotient  $1.333e+40$  so that  $1.333e+40 \cdot 0.75$  does not exceed  $1e+40$ . We multiply the found quotient with  $0.75$ , and then subtract it from  $1e+40$  to find the dividend  $4.836e+24$  for the second iteration.

Our basic algorithm is shown in Algorithm 1. If  $n$ ,  $d$  and  $q_k$  are FP numbers of the same precision  $p$ , then  $r_k$  is representable with an FP number of precision  $2p$ . In this case, the number of iterations can be minimized by substituting  $q_k$  with the largest FP number of precision  $p$  within the range specified at line 3. However, the algorithm still works if  $q_k$  is any FP number of precision  $p$  within the range. By utilizing this property, an implementer can use a division operator that does not return a correctly rounded result. The source code of an implementation of this algorithm is shown in Fig. 7 in Appendix B. A part of the proof of correctness is shown in Appendix C.

FDLIBM uses a method of shift and subtract. It first converts the mantissa of two given arguments into 64-bit integers, and calculates a remainder in a bit-by-bit basis. The main loop iterates  $ix - iy$  times, where  $ix$  and  $iy$  are the exponents of the arguments of `fmod`. This loop includes 10 integer additions and 3 `if` statements. The number of iterations of the main loop can reach more than 1000.

Vector-libm does not implement FP remainder.

**5.10 Handling of Special Numbers, Exception and Flags**

Our implementation gives a value within the specified error bound without special handling of denormal numbers, unless otherwise noted.

When a function has to return a specific value for a specific value of an argument (such as a NaN or a negative zero) is given, such a condition is checked at the end of each function. The return value is substituted with the special value if the condition is met. This process is complicated in functions like `pow`, because they have many conditions for returning special values.

SLEEF functions do not give correct results if the computation mode is different from round-to-nearest. They do not set `errno` nor raise an exception. This is a common behavior among vectorized math libraries including Libmvec [11] and SVML [3]. Because of SIMD processing, functions can raise spurious exceptions if they try to raise an exception.

**5.11 Summary**

FDLIBM extensively uses conditional branches in order to switch the polynomial according to the argument (`sin`, `cos`, `tan`, `log`, etc), to return a special value if the arguments are special values (`pow`, etc.), and to control the number of iterations (the Payne-Hanek reduction).

Vector-libm switches between a few polynomials in most of the functions. It does not provide functions with 1-ULP error bound, nevertheless, the numbers of non-zero terms in the polynomials are larger than other two libraries in some of the functions. A vectorized path is used only if the argument is smaller than  $2.574$  in `cos` and `tan`, although these functions are frequently evaluated with an argument up to  $2\pi$ . In most of the functions, Vector-libm uses a non-vectorized path if the argument is very small or a non-finite number. For example, it processes  $0$  with non-vectorized paths in many functions, although  $0$  is a frequently evaluated argument in normal situations. If non-finite numbers are once contained in data being processed, the whole processing can become significantly slower afterward. Variation in execution time can be exploited for a side-channel attack in cryptographic applications.

SLEEF uses the fastest paths if all the arguments are under  $15$  for trigonometric functions, and the same vectorized path is used regardless of the argument in most of the non-trigonometric functions. SLEEF always uses the same polynomial regardless of the argument in all functions.

Although reducing the number of conditional branches has a few advantages in implementing vector math libraries, it seems to be not given a high priority in other libraries.

## 6 TESTING

SLEEF includes three kinds of testers. The first two kinds of testers test the accuracy of all functions against high-precision evaluation using the MPFR library. In these tests, the computation error in ULP is calculated by comparing the values output by each SLEEF function and the values output by the corresponding function in the MPFR library, and it is checked if the error is within the specified bounds.

### 6.1 Perfunctory Test

The first kind of tester carries out a perfunctory set of tests to check if the build is correct. These tests include standards compliance tests, accuracy tests and regression tests.

In the standards compliance tests, we test if the functions return the correct values when values that require special handling are given as the argument. These argument values include  $\pm\text{Inf}$ , NaN and  $\pm 0$ . `Atan2` and `pow` are binary functions and have many combinations of these special argument values. These are also all tested.

In the accuracy test, we test if the error of the returned values from the functions is within the specified range, when a predefined set of argument values are given. These argument values are basically chosen between a few combinations of two values at regular intervals. The trigonometric functions are also tested against argument values close to integral multiples of  $\pi/2$ . Each function is tested against tens of thousands of argument values in total.

In the regression test, the functions are tested with argument values that triggered bugs in the previous library release, in order to prevent re-emergence of the same bug.

The executables are separated into a tester and IUTs (Implementation Under Test). The tests are carried out by making these two executables communicate via an input/output pipeline, in order to enable testing of libraries for architectures which the MPFR library does not support.

### 6.2 Randomized Test

The second kind of tester is designed to run continuously. This tester generates random arguments and compare the output from each function to the output calculated with the corresponding function in the MPFR library. This tester is expected to find bugs if it is run for a sufficiently long time.

In order to randomly generate an argument, the tester generates random bits of the size of an FP value, and reinterprets the bits as an FP value. The tester executes the randomized test for all the functions in the library at several thousand arguments per second for each function on a computer with a Core i7-6700 CPU.

In the SLEEF project, we use randomized testing in order to check the correctness of functions, rather than formal verification. It is indeed true that proving correctness of implementation contributes to the reliability of implementation. However, there is a performance overhead because the way of implementation is limited in a form that is easy to prove the correctness. There would be an increased cost of maintaining the library because of the need for updating the proof each time the implementation is modified.

### 6.3 Bit-Identity Test

The third kind of tester is for testing if bit-identical results are returned from the functions that are supposed to return such results. This test is designed to compare the results among the binaries compiled with different vector extensions. For each predetermined list of arguments, we calculate an MD5 hash value of all the outputs from each function. Then, we check if the hash values match among functions for different architectures.

## 7 PERFORMANCE COMPARISON

In this section, we present results of a performance comparison between FDLIBM Version 5.3 [4], Vector-libm [5], SLEEF 3.4, and Intel SVML [3] included in Intel C Compiler 19.

We measured the reciprocal throughput of each function by measuring the execution time of a tight loop that repeatedly calls the function in a single-threaded process. In order to obtain useful results, we turned off optimization flags when compiling the source code of this tight loop,<sup>3</sup> while the libraries are compiled with their default optimization options. We did not use LTO. We confirmed that the calls to the function are not compiled out or inlined by checking the assembly output from the compiler. The number of function calls by each loop is  $10^{10}$ , and the execution time of this loop is measured with the `clock_gettime` function.

We compiled SLEEF and FDLIBM using `gcc-7.3.0` with `"-O3 -mavx2 -mfma"` optimization options. We compiled Vector-libm using `gcc-7.3.0` with the default `"-O3 -march=native -ftree-vectorize -ftree-vectorizer-verbose=1 -fno-math-errno"` options. We changed `VECTOR_LENGTH` in `vector.h` to 4 and compiled the source code on a computer with an Intel Core i7-6700 CPU.<sup>4</sup> The accuracy of functions in SVML can be chosen by compiler options. We specified an `"-fimf-max-error=1.0"` and an `"-fimf-max-error=4.0"` options for `icc` to obtain the 1-ULP and 4-ULP accuracy results, respectively.

We carried out all the measurements on a physical PC with Intel Core i7-6700 CPU @ 3.40GHz without any virtual machine. In order to make sure that the CPU is always running at the same 3.4GHz clock speed during the measurements, we turned off Turbo Boost. With this setting, 10 nano sec. corresponds to 34 clock cycles.

The following results compare the the reciprocal throughput of each function. If the implementation is vectorized and each vector has  $N$  elements of FP numbers, then a single execution evaluates the corresponding mathematical function  $N$  times. We generated arguments in advance and stored in arrays. Each time a function is executed, we set a

3. If we turn on the optimizer, there is concern that the compiler optimizes away the call to a function. In order to prevent this, we have to introduce extra operations, but this also introduces overhead. After trying several configurations of the loop and optimizer settings, we decided to turn off the optimizer in favor of reproducibility, simplicity and fairness. We checked the assembly output from the compiler and confirmed that the unoptimized loop simply calls the target function and increments a counter, and therefore that the operations inside a loop are minimal.

4. Vector-libm evaluates functions with 512 bits of vector length by default. Because SLEEF and SVML are 256-bit wide, the setting is changed.

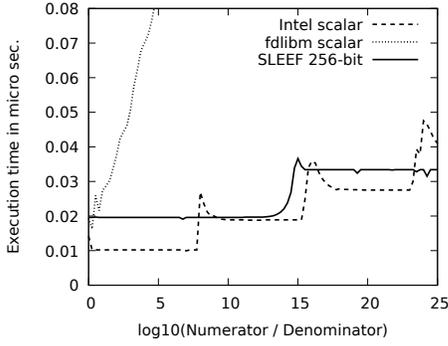


Fig. 4: Reciprocal throughput of double-precision fmod functions

randomly generated argument to each element of the argument vector (each element is set with a different value). The measurement results do not include the delay for generating random numbers.

### 7.1 Execution Time of Floating Point Remainder

We compared the reciprocal throughput of double-precision `fmod` functions in the `libm` included in Intel C Compiler 19, FDLIBM and SLEEF. All the FP remainder functions always return a correctly-rounded result. We generated a random denominator  $d$  and a numerator uniformly distributed within  $[1, 100]$  and  $[0.95r \cdot d, 1.05r \cdot d]$ , respectively, where  $r$  is varied from 1 to  $10^{25}$ . Fig. 4 shows the reciprocal throughput of the `fmod` function in each library. Please note that SVML does not contain a vectorized `fmod` function.

The graph of reciprocal throughput looks like a step function, because the number of iterations increases in this way.

### 7.2 Comparison of Overall Execution Time

We compared the reciprocal throughput of 256-bit wide vectorized double-precision functions in Vector-libm, SLEEF and SVML, and scalar functions in FDLIBM. We generated random arguments that were uniformly distributed within the indicated intervals for each function. In order to check execution speed of fast paths in trigonometric functions, we measured the reciprocal throughput with arguments within  $[0.4, 0.5]$ . The result is shown in Table 2.

The reciprocal throughput of functions in SLEEF is comparable to that of SVML in all cases. This is because the latency of FP operations is generally dominant in the execution time of math functions. Because there are two levels of scheduling mechanisms, which includes the optimizer in a compiler and the out-of-order execution hardware, there is small room for making a difference to the throughput or latency.

Execution speed of FDLIBM is not very slow despite many conditional branches. This seems to be because of a smaller number of FP operations, and faster execution speed of scalar instructions compared to equivalent SIMD instructions.

Vector-libm is slow even if only the vectorized path is used. This seems to be because Vector-libm evaluates polynomials with a large number of terms. Auto-vectorizers

TABLE 2: Reciprocal throughput in nano sec.

Func, error bound, domain	Vector-libm	FDLIBM	SLEEF	SVML
sin, 1 ulp, [0.4, 0.5]		4.927	11.43	13.68
sin, 4 ulps, [0.4, 0.5]	9.601		7.504	6.679
sin, 1 ulp, [0, 6.28]		18.96	11.41	13.86
sin, 4 ulps, [0, 6.28]	12.48		7.507	6.723
sin, 1 ulp, [0, 1e + 100]		162.3	48.79	41.72
sin, 4 ulps, [0, 1e + 100]	288.6		43.82	34.96
cos, 1 ulp, [0.4, 0.5]		11.42	13.75	12.99
cos, 4 ulps, [0.4, 0.5]	9.557		7.850	7.917
cos, 1 ulp, [0, 6.28]		18.45	13.74	13.18
cos, 4 ulps, [0, 6.28]	13.97		7.850	7.838
cos, 1 ulp, [0, 1e + 100]		162.1	50.38	38.38
cos, 4 ulps, [0, 1e + 100]	360.3		46.09	36.57
tan, 1 ulp, [0.4, 0.5]		7.819	17.30	15.71
tan, 4+ ulps, [0.4, 0.5]	15.58		9.367	7.570
tan, 1 ulp, [0, 6.28]		22.24	17.28	15.78
tan, 4+ ulps, [0, 6.28]	20.16		9.367	7.595
tan, 1 ulp, [0, 1e + 100]		177.0	48.82	43.31
tan, 4+ ulps, [0, 1e + 100]	399.4		36.54	40.50
asin, 1 ulp, [-1, 1]		14.87	12.99	12.10
asin, 4+ ulps, [-1, 1]	20.75		5.552	9.627
acos, 1 ulp, [-1, 1]		12.07	16.09	12.11
acos, 4+ ulps, [-1, 1]	23.62		7.572	10.23
atan, 1 ulp, [-700, 700]		10.16	22.12	19.97
atan, 4+ ulps, [-700, 700]	35.54		9.251	12.09
log, 1 ulp, [0, 1e + 300]		31.66	15.46	12.05
log, 4 ulps, [0, 1e + 300]	39.64		9.636	8.842
exp, 1 ulp, [-700, 700]		12.19	7.663	7.968
exp, 4 ulps, [-700, 700]	17.35			6.756
pow, 1 ulp, [-30, 30][−30, 30]		69.40	55.53	75.18

are still developing, and the compiled binary code might not be well optimized. When a slow path has to be used, Vector-libm is even slower since a scalar evaluation has to be carried out for each of the elements in the vector.

Vector-libm uses Horner's method to evaluate polynomials, which involves long latency of chained FP operations. In FDLIBM, this latency is reduced by splitting polynomials into even and odd terms, which can be evaluated in parallel. SLEEF uses Estrin's scheme. In our experiments, there was

only a small difference between Estrin’s scheme and splitting polynomials into even and odd terms with respect to execution speed.

## 8 CONCLUSION

In this paper, we showed that our SLEEF library shows performance comparable to commercial libraries while maintaining good portability. We have been continuously developing SLEEF since 2010.<sup>5</sup> [52] We distribute SLEEF under the Boost Software License [53], which is a permissive open source license. We actively communicate with developers of compilers and members of other projects in order to understand the needs of real-world users. The Vector Function ABI is important in developing vectorizing compilers. The functions that return bit-identical results are added to our library to reflect requests from our multiple partners. We thoroughly tested these functionalities, and SLEEF is already adopted in multiple commercial products.

## APPENDIX A

### ANNOTATED SOURCE CODE OF `tan`

Fig. 5 and 6 shows a C source code of our implementation of the tangent function with 1-ULP error bound. We omitted the second Cody-Waite reduction for the sake of simplicity. The definitions of DD operators are provided in Table 3. In the implementation of our library, we wrote all the operators with VEAL, as described in Sec. 3. The only conditional branch in this source code is the `if` statement at line 6. We implemented the other `if` statements at line 14, 26, 57, and 62 with conditional move operators. The `for` loop at line 16 is unrolled. We implemented `round` functions at line 8, 19 and 20 with a single instruction for most of the vector extensions. Macro `ESTRIN` at line 35 evaluates a polynomial in double-precision with Estrin’s scheme.

In the loop from line 16 to 23 in the Payne-Hanek reduction, Eq. (1) is computed. The result is multiplied by  $\pi/2$  at line 24. The numbers of FP operators shown in Table 1 are the numbers of operators from line 12 to line 26. The path with the Payne-Hanek reduction is also taken if the argument is non-finite. In this case, variable  $x$  is set to NaN at line 26, and this will propagate to the final result.

## APPENDIX B

### ANNOTATED SOURCE CODE OF THE FP REMAINDER

Fig. 7 shows a sample C source code of our FP remainder. In this implementation, both dividend and divisor are DP numbers. It correctly handles signs and denormal numbers. It executes a division instruction only once throughout the computation of the remainder. The implementation also supports the case where a division operator does not give a correctly rounded result. In this case, the `nextafter` function must be applied multiple times at line 12 according to the maximum error. We implemented `if` statements at line 3, 21 and 23 with conditional move operators. In the main loop, the algorithm finds the remainder of  $n/d$ , and at line 22, the correct sign is assigned to the resulting remainder.

5. <https://sleef.org/>

TABLE 3: DD Functions

Function name	Output
<code>dd(x, y)</code>	DD number $x + y$
<code>ddadd2_d2_d2_d2(x, y)</code>	Sum of DD numbers $x$ and $y$
<code>ddadd_d2_d2_d2</code>	Addition of two DD numbers $x$ and $y$ , where $ x  \geq  y $
<code>ddadd_d2_d_d</code>	Addition of two DP numbers $x$ and $y$ , where $ x  \geq  y $
<code>ddadd_d2_d_d2</code>	Addition of DP number $x$ and DD number $y$ , where $ x  \geq  y $
<code>ddmul_d2_d2_d2(x, y)</code>	Product of DD numbers $x$ and $y$
<code>ddmul_d2_d2_d(x, y)</code>	Product of DD number $x$ and DP number $y$
<code>ddmul_d2_d_d(x, y)</code>	Product of DP numbers $x$ and $y$
<code>ddddiv_d2_d2_d2(x, y)</code>	Returns $x/y$ , where $x$ and $y$ are DD numbers
<code>ddsqu_d2_d2(x)</code>	Returns $x^2$ , where $x$ is a DD number
<code>ddsqu_d2_d2_d(x, y)</code>	Product of DD number $x$ and DP number $y$ , where $y = 2^N$
<code>ddnormalize_d2_d2(x)</code>	Re-normalize DD number $x$

The stopping condition  $r \cdot x \geq d$  of the `for` loop at line 11 is not strictly required, and this loop can be terminated after all the values in vectors satisfy the stopping condition in a vectorized implementation. Since we assume that all operators return a round-to-nearest FP number, we use the `nextafter` function at line 9 and 12 to find a value close to  $r/d$  but not exceeding it. This method is applicable only if  $x/y$  is smaller than `DBL_MAX`. In order to find the correct  $q$  when  $r/d$  is between 1 and 3, we use a few comparisons to detect this case at line 13 and 14.

## APPENDIX C

### CORRECTNESS OF FP REMAINDER

It is obvious that Algorithm 1 returns a correct result. We show partial proof that  $r_k$  is representable as a radix-2 FP number of precision  $2p$  if  $n$ ,  $d$  and  $q_k$  are radix-2 FP numbers of precision  $p$ .

We now define property  $FP_p$  and function  $RD_p$ .  $FP_p(x)$  holds iff there are integer  $0 \leq m < 2^p$  and integer  $e$  that satisfy  $x = m \cdot 2^e$ . If  $x$  and  $y$  are finite DP numbers,  $FP_{53}(x)$  and  $FP_{106}(xy)$  holds.  $RD_p(x)$  denotes the maximum number that does not exceed  $x$  and  $FP_p(x)$  holds.

We now show that  $FP_{2p}(r_k)$  and  $FP_{2p}(q_k d)$  hold if  $FP_p(n)$  and  $FP_p(d)$  hold. By Lemma 1, integer  $q_k$  exists that satisfies  $FP_p(q_k)$ . Then,  $FP_{2p}(r_0)$  and  $FP_{2p}(q_k d)$  hold.  $FP_{2p}(r_k - q_k d)$  holds because  $r_k/2 \leq q_k d \leq 2r_k$ . Then,  $FP_{2p}(r_{k+1})$  holds.

**Lemma 1** Given  $s \geq 1$  and integer  $p \geq 2$ . There exists an integer  $q$  that satisfies  $FP_p(q)$  and  $s/2 \leq q \leq s$ .

**Proof:** If  $s < 2^p$ ,  $FP_p(\lfloor s \rfloor)$  holds. In this case,  $q$  can be  $\lfloor s \rfloor$  since  $\lfloor s \rfloor \geq 1$ . Otherwise,  $RD_p(s)$  is an integer.  $(s - RD_p(s))/s < 2^{1-p}$ , and therefore  $(1 - 2^{1-p})s < RD_p(s)$ . Because  $p \geq 2$ ,  $s/2 < RD_p(s)$ , and therefore  $q$  can be  $RD_p(s)$ .

We now discuss the number of iterations. We suppose  $q_k$  is set to  $\min(RD_p(r_k/d), \lfloor r_k/d \rfloor)$ . If  $r_k/d < 2^p$ ,  $\lfloor r_k/d \rfloor = q_k$  and the loop terminates after  $k$ -th iteration. Otherwise,  $q_k = RD_p(r_k/d) < \lfloor r_k/d \rfloor$ , and  $(r_k/d - q_k)/(r_k/d) < 2^{1-p}$ .

Then,  $r_{k+1} = r_k - q_k d < 2^{1-p} r_k$ . Therefore the number of iterations is  $\lceil (\log_2(n/d))/(p-1) \rceil$ .

In order to show that the number that substitutes  $q$  at line 12 in Fig. 7 satisfies the condition at line 3 in Algorithm 1, we prove  $(r/d)/2 < \lfloor RN(\text{nextafter}(RN(r), 0) \cdot \text{nextafter}(RN(1.0/d), 0)) \rfloor$  assuming  $p = 53$ ,  $r/d \geq 3$  and  $d > 0$ .  $RN(x)$  is the floating-point number that is the closest to  $x$ . We assume no overflow or underflow.

It is obvious that  $q$  is substituted with a number that is smaller or equal to  $r/d$ .  $(r - \text{nextafter}(RN(r), 0))/r < 2^{2-p}$ , where  $p$  is precision. Therefore,  $(1 - 2^{2-p}) \cdot r < \text{nextafter}(RN(r), 0)$ . Similarly,  $(1 - 2^{2-p}) \cdot 1.0/d < \text{nextafter}(RN(1.0/d), 0)$ . Therefore,  $(1 - 2^{2-p})^2 (1 - 2^{1-p}) r/d < RN(\text{nextafter}(RN(r), 0) \cdot \text{nextafter}(RN(1.0/d), 0))$ . Let  $u = (1 - 2^{2-p})^2 (1 - 2^{1-p})$ .  $ur/d - 1 \leq \lfloor ur/d \rfloor$ .  $(r/d)/2 < ur/d - 1 \Leftrightarrow 1/(u - 1/2) < r/d$ .  $1/(u - 1/2) < r/d$  holds since  $p = 53$  and  $r/d \geq 3$ . Therefore,  $(r/d)/2 < \lfloor RN(\text{nextafter}(RN(r), 0) \cdot \text{nextafter}(RN(1.0/d), 0)) \rfloor$ .

## ACKNOWLEDGMENTS

We wish to acknowledge Will Lovett and Srinath Vadlamani for their valuable input and suggestions. The authors are particularly grateful to Robert Werner, who reviewed the final manuscript. The authors would like to thank Prof. Leigh McDowell for his suggestions. The authors would like to thank all the developers that contributed to the SLEEP project, in particular Diana Bite and Alexandre Mutel.

## REFERENCES

- [1] G. Conte, S. Tommesani, and F. Zanichelli, "The long and winding road to high-performance image processing with mmx/sse," in *Proceedings Fifth IEEE International Workshop on Computer Architectures for Machine Perception*, Sep. 2000, pp. 302–310.
- [2] D. Naishlos, "Autovectorization in GCC," in *Proceedings of the 2004 GCC Developers Summit*, 2004, pp. 105–118. [Online]. Available: <http://people.redhat.com/lockhart/.gcc2004/MasterGCC-2side.pdf>
- [3] Intel Corporation. (2019) Intel short vector math library. [Online]. Available: <https://software.intel.com/en-us/node/523613>
- [4] Sun Microsystems, Inc. (2010) Sun freely distributable libm version 5.3. [Online]. Available: <http://www.netlib.org/fdlibm/>
- [5] C. Lauter, "A new open-source SIMD vector libm fully implemented with high-level scalar C," in *2016 50th Asilomar Conference on Signals, Systems and Computers*, November 2016, pp. 407–411.
- [6] ISO, *ISO/IEC 9899:2011 Information technology — Programming languages — C*. International Organization for Standardization, December 2011. [Online]. Available: <https://www.iso.org/standard/57853.html>
- [7] GNU Project. (2018) The GNU C library (glibc). [Online]. Available: <https://www.gnu.org/software/libc/>
- [8] S. Gal, "An accurate elementary mathematical library for the IEEE floating point standard," *ACM Trans. Math. Softw.*, vol. 17, no. 1, pp. 26–45, March 1991.
- [9] C. Daramy-Loirat, D. Defour, F. de Dinechin, M. Gallet, N. Gast, C. Lauter, and J.-M. Muller, "CR-LIBM: a correctly rounded elementary function library," in *Proc. SPIE 5205, Advanced Signal Processing Algorithms, Architectures, and Implementations XIII*, vol. 5205, December 2003, pp. 458–464.
- [10] Advanced Micro Devices, Inc. (2013) AMD core math library. [Online]. Available: <http://developer.amd.com/tools-and-sdks/archive/acml-product-features>
- [11] GNU Project. (2015) Libmvec in glibc. [Online]. Available: <https://sourceware.org/glibc/wiki/libmvec>
- [12] C. K. Anand and W. Kahl, "An optimized cell be special function library generated by coconut," *IEEE Transactions on Computers*, vol. 58, no. 8, pp. 1126–1138, Aug 2009.
- [13] M. Dukhan et al. (2013) Yeppp! library. [Online]. Available: <https://bitbucket.org/MDukhan/yeppp>
- [14] M. Dukhan and R. Vuduc, "Methods for high-throughput computation of elementary functions," in *Parallel Processing and Applied Mathematics: 10th International Conference*, May 2014, pp. 86–95.
- [15] D. Piparo, V. Innocente, and T. Hauth, "Speeding up HEP experiment software with a library of fast and auto-vectorisable mathematical functions," *Journal of Physics: Conference Series*, vol. 513, no. 5, p. 052027, 2014.
- [16] S. Manilov, B. Franke, A. Magrath, and C. Andrieu, "Free rider: A source-level transformation tool for retargeting platform-specific intrinsic functions," *ACM Trans. Embed. Comput. Syst.*, vol. 16, no. 2, pp. 38:1–38:24, Dec. 2016. [Online]. Available: <http://doi.acm.org/10.1145/2990194>
- [17] M. Gross, "Neat SIMD: Elegant vectorization in C++ by using specialized templates," in *2016 International Conference on High Performance Computing Simulation (HPCS)*, July 2016, pp. 848–857.
- [18] N. Clark, A. Hormati, S. Yehia, S. Mahlke, and K. Flautner, "Liquid SIMD: Abstracting SIMD Hardware using Lightweight Dynamic Mapping," in *2007 IEEE 13th International Symposium on High Performance Computer Architecture*, Feb 2007, pp. 216–227.
- [19] R. Leiða, S. Hack, and I. Wald, "Extending a c-like language for portable SIMD programming," *SIGPLAN Not.*, vol. 47, no. 8, pp. 65–74, Feb. 2012.
- [20] Arm Limited, *ARM NEON Intrinsics Reference*, September 2014, no. IHI0073A. [Online]. Available: [http://infocenter.arm.com/help/topic/com.arm.doc.ih0073a/IHI0073A\\_arm\\_](http://infocenter.arm.com/help/topic/com.arm.doc.ih0073a/IHI0073A_arm_)
- [21] J.-M. Muller, *Elementary Functions: Algorithms and Implementation*. Birkhauser Boston, Inc., 1997.
- [22] Arm Limited. (2017) ARM C language extensions for SVE documentation. [Online]. Available: <https://developer.arm.com/docs/100987/0000>
- [23] L. Liu, S. Peng, C. Zhang, R. Li, B. Wang, C. Sun, Q. Liu, L. Dong, L. Li, S. Yanyan, Y. He, W. Zhao, and G. Yang, "Importance of bitwise identical reproducibility in earth system modeling and status report," *Geoscientific Model Development Discussions*, vol. 8, pp. 4375–4400, June 2015.
- [24] J.-M. Muller, N. Brisebarre, F. de Dinechin, C.-P. Jeannerod, V. Lefèvre, G. Melquiond, N. Revol, D. Stehlé, and S. Torres, *Handbook of Floating-Point Arithmetic*, 1st ed. Birkhauser Boston, Inc., 2009.
- [25] O. Krzikalla, K. Feldhoff, R. Müller-Pfefferkorn, and W. E. Nagel, "Auto-Vectorization Techniques for Modern SIMD Architectures," in *Proc. of the 16th Workshop on Compilers for Parallel Computing*, January 2012.
- [26] D. Nuzman, S. Dyshel, E. Rohou, I. Rosen, K. Williams, D. Yuste, A. Cohen, and A. Zaks, "Vapor SIMD: Auto-vectorize once, run everywhere," in *International Symposium on Code Generation and Optimization*, ser. CGO 2011, April 2011, pp. 151–160.
- [27] R. M. Russell, "The CRAY-1 computer system," *Commun. ACM*, vol. 21, no. 1, pp. 63–72, Jan. 1978.
- [28] D. Nuzman, I. Rosen, and A. Zaks, "Auto-vectorization of interleaved data for SIMD," in *Proceedings of the 27th ACM SIGPLAN Conference on Programming Language Design and Implementation*, ser. PLDI '06, 2006, pp. 132–143.
- [29] X. Tian, A. Bik, M. Girkar, P. Grey, H. Saito, and E. Su, "Intel OpenMP C++/Fortran Compiler for Hyper-Threading Technology: Implementation and Performance," in *Intel Technology Journal*, vol. 6, no. 1, Feb 2002, pp. 36–46.
- [30] GNU Project. (1984) GCC, the GNU compiler collection. [Online]. Available: <https://www.gnu.org/software/gcc/>
- [31] Arm Limited. Arm Compiler for HPC. [Online]. Available: <https://developer.arm.com/products/software-development-tools/hpc/arm-c>
- [32] The LLVM compiler infrastructure project. (2004) The LLVM Compiler Infrastructure. [Online]. Available: <https://llvm.org/>
- [33] C. Lattner and V. Adve, "LLVM: A Compilation Framework for Lifelong Program Analysis & Transformation," in *Proceedings of the 2004 International Symposium on Code Generation and Optimization (CGO'04)*, Palo Alto, California, Mar 2004.
- [34] Arm Limited. (2018) Vector function application binary interface specification for AArch64. [Online]. Available: <https://developer.arm.com/docs/101129/latest>
- [35] OpenMP Architecture Review Board. (2013) OpenMP application program interface. [Online]. Available: <https://www.openmp.org/wp-content/uploads/OpenMP4.0.0.pdf>
- [36] J. Lee, F. Petrogalli, G. Hunter, and M. Sato, "Extending OpenMP SIMD support for target specific code and application to ARM SVE," in *Scaling OpenMP for Exascale Performance and Portability*, August 2017, pp. 62–74.

- [37] X. Tian, H. Saito, S. Kozhukhov, K. B. Smith, R. Geva, M. Girkar, and S. V. Preis. (2015, November) Vector function application binary interface. Intel Corporation. [Online]. Available: <https://software.intel.com/en-us/articles/vector-simd-function-abi>
- [38] A. Sodani, R. Gramunt, J. Corbal, H. S. Kim, K. Vinod, S. Chinthamani, S. Hutsell, R. Agarwal, and Y. C. Liu, "Knights Landing: Second-Generation Intel Xeon Phi Product," *IEEE Micro*, vol. 36, no. 2, pp. 34–46, March 2016.
- [39] S. Eyerhan, J. E. Smith, and L. Eckhout, "Characterizing the branch misprediction penalty," in *2006 IEEE International Symposium on Performance Analysis of Systems and Software*, March 2006, pp. 48–58.
- [40] Intel Corporation, *Intel 64 and IA-32 Architectures Optimization Reference Manual*, June 2016, no. 248966-033. [Online]. Available: <https://software.intel.com/en-us/download/intel-64-and-ia-32-architectures-optimization-reference-manual>
- [41] A. Fog. Instruction tables: Lists of instruction latencies, throughputs and micro-operation breakdowns for Intel, AMD and VIA CPUs. [Online]. Available: [https://www.agner.org/optimize/instruction\\_tables.pdf](https://www.agner.org/optimize/instruction_tables.pdf)
- [42] A. Abel and J. Reineke, "Uops.info: Characterizing latency, throughput, and port usage of instructions on intel microarchitectures," in *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*, ser. ASPLOS '19. New York, NY, USA: ACM, 2019, pp. 673–686.
- [43] Google's Compiler Research Team and The LLVM compiler infrastructure project. (2018) *llvm-exegesis: Automatic measurement of instruction latency/uops*. [Online]. Available: <https://github.com/google/EXEgesis>
- [44] T. J. Dekker, "A floating-point technique for extending the available precision," *Numer. Math.*, vol. 18, no. 3, pp. 224–242, June 1971.
- [45] J. Richard Shewchuk, "Adaptive precision floating-point arithmetic and fast robust geometric predicates," *Discrete & Computational Geometry*, vol. 18, no. 3, pp. 305–363, Oct 1997.
- [46] W. Cody and W. Waite. (1980) *Software manual for the elementary functions*. Englewood Cliffs, N.J.
- [47] W. J. Cody, "Implementation and testing of function software," in *Problems and Methodologies in Mathematical Software Production*, November 1980, pp. 24–47.
- [48] M. H. Payne and R. N. Hanek, "Radian reduction for trigonometric functions," *SIGNUM Newsl.*, vol. 18, no. 1, pp. 19–24, January 1983.
- [49] B. W. Char, K. O. Geddes, W. M. Gentleman, and G. H. Gonnet, "The design of maple: A compact, portable, and powerful computer algebra system," in *Computer Algebra*, J. A. van Hulzen, Ed. Berlin, Heidelberg: Springer Berlin Heidelberg, 1983, pp. 101–115.
- [50] S. Chevillard, M. Joldes, and C. Lauter, "Sollya: An environment for the development of numerical codes," in *Mathematical Software - ICMS 2010*, ser. Lecture Notes in Computer Science, K. Fukuda, J. van der Hoeven, M. Joswig, and N. Takayama, Eds., vol. 6327. Heidelberg, Germany: Springer, September 2010, pp. 28–31.
- [51] G. Estrin, "Organization of computer systems: The fixed plus variable structure computer," in *Papers Presented at the May 3-5, 1960, Western Joint IRE-AIEE-ACM Computer Conference*, ser. IRE-AIEE-ACM '60 (Western). New York, NY, USA: ACM, 1960, pp. 33–40. [Online]. Available: <http://doi.acm.org/10.1145/1460361.1460365>
- [52] N. Shibata, "Efficient evaluation methods of elementary functions suitable for SIMD computation," *Computer Science - Research and Development*, vol. 25, no. 1, pp. 25–32, May 2010.
- [53] D. Abrahams, D. Cabell, D. Smith, and E. Chan. (2003, August) Boost Software License 1.0 (BSL-1.0). [Online]. Available: [https://www.boost.org/LICENSE\\_1\\_0.txt](https://www.boost.org/LICENSE_1_0.txt)

PLACE  
PHOTO  
HERE

**Naoki Shibata** is an associate professor at Nara Institute of Science and Technology. He received the Ph.D. degree in computer science from Osaka University, Japan, in 2001. He was an assistant professor at Nara Institute of Science and Technology 2001-2003 and an associate professor at Shiga University 2004-2012. His research areas include distributed and parallel systems, and intelligent transportation systems. He is a member of IPSJ, ACM and IEEE.

PLACE  
PHOTO  
HERE

**Francesco Petrogalli** is a software engineer working on the development of Arm Compiler for HPC. He contributed to the implementation of the Vector Length Agnostic (VLA) vectorizer for the Scalable Vector Extension (SVE) of Arm, to the ABI specifications for the vector functions on AArch64, and to the open source library SLEEF. Francesco has also worked on optimizing a variety of computational kernels that are core to Machine Learning algorithms.

```

1 double xtan(double d) {
2   double u;
3   double2 x = dd(0, 0), y;
4   int q = 0;
5
6   if (fabs(d) < 15) {
7     // Cody-Waite
8     q = round(d * M_2_PI);
9     u = fma(q, -0x1.921fb54442d18p0, d);
10    x = ddadd_d2_d_d(u, q * -0x1.1a62633145c07p-54);
11  } else {
12    // Payne-Hanek
13    int ex = ilogb(d), M = ex > 700 ? -130 : -2;
14    if (ex < 0) ex = 0;
15    u = ldexp(d, M);
16    for(int i=0;i<4;i++) {
17      y = ddmul_d2_d_d(u, tab[ex*4+i]);
18      x = ddadd2_d2_d2_d2(x, y);
19      double r = round(4*x.x);
20      q += (int32_t)((r - 4*round(x.x)));
21      x.x -= r * 0.25;
22      x = ddnormalize_d2_d2(x);
23    }
24    x = ddmul_d2_d2_d2(x,
25      dd(0x1.921fb54442d18p2, 0x1.1a62633145c07p-52));
26    if (!isfinite(d)) x.x = NAN; // NAN handling
27  }
28
29  // Reduction with double-angle formula
30  y = ddscale_d2_d2_d(x, 0.5);
31
32  // Polynomial evaluation with Estrin's scheme
33  // Domain : |y| <= PI/8
34  x = ddsqu_d2_d2(y);
35  u = ESTRIN(x.x,
36    +0.3245098826639276316e-3,
37    +0.5619219738114323735e-3,
38    +0.1460781502402784494e-2,
39    +0.3591611540792499519e-2,
40    +0.8863268409563113126e-2,
41    +0.2186948728185535498e-1,
42    +0.5396825399517272970e-1,
43    +0.1333333333330500581e+0);
44
45  // Last two terms are evaluated with Horner's
46  // method
47  u = fma(u, x.x, +0.33333333333333343695e+0);
48
49  // Last term is evaluated in DD precision
50  x = ddadd_d2_d2_d2(y,
51    ddmul_d2_d2_d(ddmul_d2_d2_d2(x, y), u));
52
53  // Reconstruction with double-angle formula
54  y = ddadd_d2_d_d2(-1, ddsqu_d2_d2(x));
55  x = ddscale_d2_d2_d(x, -2);
56
57  // Reconstruction with tan(PI/2 - x) = cot(x)
58  if (q & 1) {
59    double2 t = x; x = y; y.x = -t.x; y.y = -t.y;
60  }
61  x = dddiv_d2_d2_d2(x, y);
62
63  if (d == 0) return d; // Negative-zero handling
64  return x.x + x.y;
65 }

```

Fig. 5: C source code of tan

```

1 #include <mpfr.h>
2
3 double tab[4096];
4
5 void init() {
6   mpfr_set_default_prec(1280);
7
8   mpfr_t pi, twoopi, m;
9   mpfr_inits(pi, twoopi, m, NULL);
10  mpfr_const_pi(pi, GMP_RNDN);
11  mpfr_d_div(twoopi, 2, pi, GMP_RNDN);
12
13  for(int ex=0;ex<1024;ex++) {
14    int M = ex > 700 ? -128 : 0;
15    mpfr_set(m, twoopi, GMP_RNDN);
16    mpfr_set_exp(m, mpfr_get_exp(m) + (ex - 53));
17    mpfr_frac(m, m, GMP_RNDN);
18    mpfr_set_exp(m, mpfr_get_exp(m) - (ex - 53 + M));
19
20    for(int i=0;i<4;i++) {
21      union { double d; int64_t i; } tmp = { .d =
22        mpfr_get_d(m, GMP_RNDN) };
23      tmp.i &= 0xfffffffffffffeLL;
24      tab[ex*4+i] = tmp.d;
25      mpfr_sub_d(m, m, tab[ex*4+i], GMP_RNDN);
26    }
27  }
28  mpfr_clears(pi, twoopi, m, NULL);
29 }

```

```

1 double xfmod(double x, double y) {
2   double n = fabs(x), d = fabs(y), s = 1, q;
3   if (d < DBL_MIN) {
4     n *= 1ULL << 54;
5     d *= 1ULL << 54;
6     s = 1.0 / (1ULL << 54);
7   }
8   double2 r = dd(n);
9   double rd = nextafter(1.0 / d, 0);
10
11  for(int i=0;i < 21 && r.x >= d;i++) {
12    q = trunc(nextafter(r.x, 0) * rd);
13    q = (3*d > r.x && r.x > d) ? 2 : q;
14    q = (2*d > r.x && r.x > d) ? 1 : q;
15    q = r.x == d ? (r.y >= 0 ? 1 : 0) : q;
16    r = ddadd2_d2_d2_d2(r, ddmul_d2_d_d(q, -d));
17    r = ddnormalize_d2_d2(r);
18  }
19
20  double ret = r.x * s;
21  if (r.x + r.y == d) ret = 0;
22  ret = copysign(ret, x);
23  if (n < d) ret = x;
24  return d == 0 ? NAN : ret;
25 }

```

Fig. 6: C source code of Payne-Hanek table generation

Fig. 7: C source code of the FP remainder function