

Формальные подходы к тестированию математических функций

В. В. Кулямин

Институт системного программирования РАН (ИСП РАН),
Б. Коммунистическая, 25, Москва, Россия
E-mail: kuliamin@ispras.ru

Аннотация

Данная статья рассматривает вопросы проверки корректности вычисления математических функций на числах с плавающей точкой, формат которых определяется стандартом IEEE 754. В ней описывается метод разработки тестов для реализаций таких функций, основанный на формальных спецификациях их поведения. Предлагаемый метод основан на технологии разработки тестов UniTESK и двух дополнительных методиках: методике формирования строгих требований к реализации математической функции и методике построения набора тестовых данных для ее тестирования. Описанные методики опираются на специфические свойства представления чисел с плавающей точкой и особенности поведения самой тестируемой функции.

1. Введение

На современном уровне развития общества большое значение в его технологической инфраструктуре играют информационные системы. На них перекладывается все больше функций по хранению, обработке и адекватному представлению информации, являющейся одним из базовых элементов сегодняшней экономической и культурной жизни человечества.

Удивительно, но до сих пор нет эффективных технологий разработки информационных систем, работающих *правильно*, т. е. так, как того ожидают их пользователи. Это обусловлено двумя причинами. Во-первых, само понятие правильности такой системы неуловимо и неформально, поскольку оно тесно связано с набором нечетких, противоречивых и меняющихся со временем ожиданий и потребностей людей, общающихся с ней. Во-вторых, известные недостатки человеческой природы не позволяют разработчикам программ писать их без ошибок, даже когда они четко понимают поставленную задачу и имеют математически выверенное ее решение.

Итеративная разработка с использованием тестирования является единственным практически работающим способом создания достаточно качественного программного обеспечения (ПО), хотя и не корректного формально, но достаточно хорошо удовлетворяющего потребности пользователей. В ходе такой разработки проектирование и написание кода чередуется с проверками работоспособности и правильности результатов работы проектируемых и разрабатываемых компонентов.

Правильность их проверяется как соответствие этих результатов ожиданиям пользователей в соответствующей ситуации.

В то же время, область использования сложного ПО в человеческой деятельности все расширяется, и во многих случаях решаемые им задачи таковы, что даже по косвенным признакам уже нелегко судить о правильности полученных с его помощью результатов. Можно привести следующие примеры таких задач.

- Моделирование космогонических процессов: развития звезд, планетных систем и галактик. Моделирование космических катастроф, последствия которых могут угрожать существованию человечества.
- Моделирование физических процессов в экстремальных условиях: движения со сверхвысокими скоростями в вязких средах, поведения элементарных частиц, поведения плазмы в термоядерном реакторе, поведения материалов в условиях сверхнизких или сверхвысоких температур и давлений, в сверхсильных магнитных полях и т.п. Все эти задачи требуются решать для обеспечения технологического прогресса, в частности, для создания и доведения до практической применимости двигателей и энергетических установок, работающих на новых принципах; создания новых материалов с заданными свойствами и пр.
- Моделирование построения и работы наносистем, необходимое для успешного развития нанотехнологий и создания механизмов, решающих важные для человека задачи при помощи манипулирования отдельными молекулами и атомами.
- Моделирование биохимических процессов, связанных с функционированием различных белков, обменом веществ между живой клеткой и окружающей средой, модификацией генетической информации, активизацией и дезактивацией генов. Понимание их механизмов поможет создать более эффективные и безопасные лекарства против многих болезней, в том числе, плохо поддающихся лечению современными средствами, получить новые, более стойкие и полезные сорта и виды используемых в сельском хозяйстве растений и животных, а в перспективе — продлить жизнь каждого человека до ее биологических пределов.
- Моделирование сложных экологических, климатических, экономических и социальных систем. Оно позволит более глубоко понять законы их развития, выработать механизмы разрешения и предотвращения кризисов и катастроф, приносящих серьезный ущерб отдельным людям и обществу в целом, а также поможет человечеству более ответственно относиться к среде его существования и заложить основы для гармоничного и стабильного развития в дальнейшем.

Все эти задачи решаются при помощи использования тех или иных методов математического моделирования и дают результаты, проверка правильности которых связана с большими затратами ресурсов и усилий многих людей, а иногда и вообще невозможна или создает очень серьезные угрозы для отдельных людей или общества в целом.

Тем не менее, есть возможность повысить надежность и правильность работы программных систем, выполняющих моделирование такого рода. Это может быть сделано за счет формальной проверки корректности работы библиотечных математических функций, на которые во многом это ПО опирается. Уверенность в надежности фундамента, на котором построены такие системы, позволит разрабатывать их более качественно и с меньшими усилиями, сосредоточившись на поиске и исправлении ошибок в других компонентах.

В данной работе изучаются проблемы проверки корректности реализаций математических функций, работающих с числами с плавающей точкой, анализируются имеющиеся достижения в этой области и предлагается метод разработки тестов для таких функций на основе формальных спецификаций их поведения. Полученный метод базируется на технологии UniTESK [1-3], дополняя ее методиками формирования точных требований к реализациям математических функций и выбора тестовых данных для их тестирования.

Сперва рассмотрим проблемы, с которыми сталкиваются исследования в этой области.

2. Проблемы корректного вычисления функций

Основные проблемы корректного вычисления математических функций связаны с дискретностью представления действительных чисел в компьютерах. Чтобы иметь возможность эффективно выполнять операции с действительными числами, они представлены в виде так называемых *чисел с плавающей точкой*, формат и правила действий над которыми определены в стандартах IEEE 754 [4] (он же — IEC 60559 [5]) и IEEE 854 [6].

IEEE 754 определяет представление двоичных чисел с плавающей точкой, IEEE 854 обобщает его, определяя и десятичные числа с плавающей точкой. Однако, поскольку в большинстве случаев на практике используется двоичное представление чисел, мы будем рассматривать только его.

2.1. Числа с плавающей точкой

Двоичное число с плавающей точкой имеет следующую структуру [4,6,7].

- Число представлено в виде набора из n бит, из которых первый бит является *знаковым битом числа*, следующие k бит отданы под представление его *экспоненты*, а оставшиеся $(n-k-1)$ бит представляют его *мантиссу*.

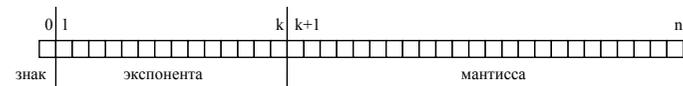


Рисунок 1. Битовое представление чисел с плавающей точкой.

- Знаковый бит S , экспонента E и мантисса M числа x определяют его значение по следующим правилам.

$$x = (-1)^S \cdot 2^e \cdot m$$
, где
 - S — знаковый бит, 0 для положительных чисел, и 1 для отрицательных;
 - если $E > 0$, то $e = E - 2^{(k-1)} + 1$;
иначе, если $E = 0$, $e = -2^{(k-1)} + 2$;
число $(2^{(k-1)} - 1)$ называется *смещением экспоненты (bias)*;
 - если $0 < E < 2^k - 1$, то m имеет двоичное представление $1.M$, т.е. целая часть m равна 1, а последовательность цифр дробной части совпадает с последовательностью бит M ;
если же $E = 0$, то m имеет двоичное представление $0.M$,
такие числа (с нулевой экспонентой) называются *денормализованными*.
- Максимальное возможное значение экспоненты $E = 2^k - 1$ зарезервировано для представления положительной $+\infty$ и отрицательной $-\infty$ бесконечностей и специального значения NaN (not-a-number), которое возникает, если результат выполняемых действий нельзя корректно представить ни обычным числом, ни бесконечностью, например, $0/0$ или $(-\infty) + (+\infty)$.
 $+\infty$ имеет нулевой знаковый бит, максимальную экспоненту и нулевую мантиссу; $-\infty$ отличается только единичным знаковым битом.
 Любое число, имеющее максимальную экспоненту и ненулевую мантиссу, считается представлением NaN.
 Все такие числа называются *исключительными*.
- Стандарты IEEE 754 и IEEE 854 определяют несколько возможных типов чисел с плавающей точкой, из которых чаще всего используются *числа однократной точности* (single precision), *числа двойной точности* (double precision) и *числа расширенной двойной точности* (double-extended precision).
 Для чисел однократной точности $n = 32$ и $k = 8$. Соответственно, для мантиссы используется 23 бита и смещение экспоненты равно 127.
 Для чисел двойной точности $n = 64$ и $k = 11$. Для мантиссы используется 52 бита и смещение экспоненты равно 1023.
 Для чисел расширенной двойной точности определенные значения k и n

в результате -0 , за исключением квадратного корня из -0 . При сложении, вычитании, умножении, делении, вычислении остатка от деления и преобразованиях типов всегда в случае нулевого результата возвращается 0 .

2.2 Требования стандартов к реализациям математических функций

Поскольку не все действительные числа представимы, возникает проблема представления результатов математических функций в тех случаях, когда такой результат не представим.

Казалось бы, естественно потребовать, чтобы реализация функции возвращала в качестве результата число с плавающей точкой, являющееся ближайшим к результату точного вычисления функции для тех же самых значений аргументов.

Однако, ни стандарты IEEE 754 и IEEE 854, ни стандарт языка C ISO/IEC 9899 [8], ни стандарт переносимого интерфейса операционной системы IEEE 1003.1 [9] (известный как POSIX), описывающие библиотеку математических функций языка C не фиксируют такого требования для большинства функций.

Стандарты IEEE 754 и IEEE 854 описывают работу только сложения, умножения, вычитания и деления чисел с плавающей точкой, а также вычисления остатка от деления, извлечения квадратного корня и преобразований между типами с плавающей точкой и между ними и целочисленными типами. Соответственно, только для этих операций требуется возвращать результат, полученный из точного приведением к ближайшему представимому числу согласно действующему режиму округления. Возможны 4 режима округления: просто к ближайшему, к 0 , к $+\infty$ и к $-\infty$. Кроме того, эти стандарты требуют аккуратного выставления флагов переполнения, слишком маленького результата или неточного результата при работе этих операций.

Стандарт C ссылается на требования IEEE 754, добавляя только ограничения на значения результатов ряда функций для некоторых значений параметров (например, $\exp(0) = \cos(0) = 1$, а $\sin(0) = \operatorname{tg}(0) = 0$). Стандарт POSIX, в свою очередь, ссылается на требования стандарта языка C, добавляя описание поведения реализаций математических функций в случае возникновения переполнения или слишком маленьких результатов и для тех значений параметров, где соответствующая функция не определена.

Отсутствие ограничений на точность вычисления математических функций может привести к накоплению погрешностей и серьезным ошибкам при многократном использовании этих функций в приложениях для математического моделирования, и, соответственно, неверным результатам работы таких приложений.

В последние 5-10 лет появились предложения стандартизировать необходимые для аккуратного моделирования требования к реализациям математических функций [10,11]. Многие инициаторы этой деятельности работают в проекте

Arenaire [12], совместно проводимом INRIA, CNRS и Высшей Нормальной школой Лиона, Франция. В результате активно разрабатывается набор стандартов ISO/IEC 10967 [13-15], формулирующий естественные ограничения на работу реализаций большинства математических функций. Эти ограничения касаются нескольких аспектов.

- Возможные погрешности вычисления функций выражены в терминах *единиц последнего разряда (unit in the last place, ulp)* [7]. Наилучшее приближение дало бы точность в 0.5 ulp , т.е. вычисленный результат отличался бы от точного значения функции не более чем на половину единицы последнего разряда мантиссы результата. Однако природа математических функций такова, что такая точность не является практически обоснованной во многих случаях, хотя ее достижение потребовало бы значительных усилий от разработчиков библиотек. Это связано с тем, что число с плавающей точкой является приближенным представлением любого действительного числа, к которому оно является ближайшим. Таким образом, уже в значениях аргументов функции может иметься погрешность, которую ее вычисление не в силах исправить. Для широко используемых математических функций были проведены оценки возрастания погрешности при их вычислении, и на основании этих оценок были сформулированы более практичные требования к точности вычислений, ограничивающие погрешность результата величиной от 0.5 ulp до 2 ulp , в зависимости от функции [14].
- Стандарты серии ISO/IEC 10967 требуют от реализации математической функции сохранения знака ее точного значения для данного значения параметра. Также требуется, чтобы во всех интервалах монотонности функции ее реализации были монотонны таким же образом, т.е. там, где сама функция убывает, ее численная реализация должна убывать, а там, где функция возрастает, — возрастать. Исключением из этого правила являются тригонометрические функции в области больших значений аргумента, где интервал смены знака и интервал монотонности становятся сравнимы с единицей последнего разряда аргумента.
- ISO/IEC 10967 требует соблюдения специфических требований при вычислении функций в окрестностях точек, где они имеют известные представимые значения. Например, реализация экспоненты для значений аргументов, достаточно близких к 0 , должна возвращать в точности 1 . Это требование связано с тем, что плотность чисел с плавающей точкой в окрестности 0 гораздо больше, чем их плотность в окрестности 1 — между 1 и ближайшим к нему числом с плавающей точкой умещается много чисел, близких к 0 . Для двойной точности ближайшее к 1 число равно $1 - 2^{-53}$, а ближайшее к 0 — это 2^{-1074} .

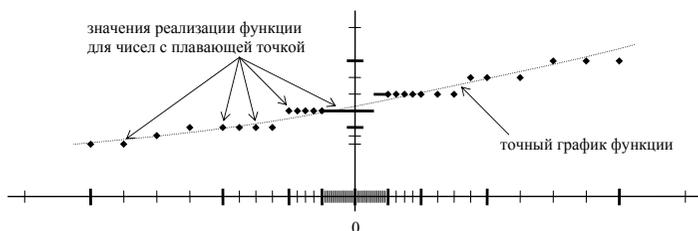


Рисунок 3. График функции, имеющей ненулевое значение в 0.

2.3. Дилемма составителя таблиц

Сформулированные требования к точности вычисления математических функций приводят к так называемой *дилемме составителя таблиц (Table Maker's Dilemma)* [16,17].

Эта проблема состоит в том, что для выбора правильно округленного ближайшего числа с плавающей точкой при приближенных вычислениях иногда нужно вычислить много дополнительных бит мантииссы результата, значительно больше, чем имеется в рассматриваемом типе чисел с плавающей точкой.

Для иллюстрации дилеммы составителя таблиц приведем следующий пример. Пусть вычисляется функция \sin для двоичных чисел с плавающей точкой, имеющих 6 битов мантииссы. Синус числа $11.1010_2 = 3.625_{10}$ (выделены биты мантииссы) равен $0.011101101111110..._2 = 0.063225984913..._{10}$ (снова выделены биты мантииссы). Приближенное вычисление 6-ти бит мантииссы результата может дать как 0.0111011_2 , так и 0.0111100_2 , поскольку точное значение очень близко к их среднему арифметическому. Только получив точный 14-й бит, мы сможем уверенно выбрать первое из них в качестве значения, ближайшего к точному результату.



Рисунок 4. Дилемма составителя таблиц.

Чтобы определить ближайшее число с плавающей точкой, погрешность вычислений иногда должна быть значительно меньше половины расстояния между соседними такими числами.

Приведем более реалистичный пример для чисел с двойной точностью, имеющих 52 бита в мантииссе. Вычисляя значение натурального логарифма для

$$1.011000010011100101010101110111001000000001011111000_2 \cdot 2^{-35}$$

получим

$$-10111.1111000000101111100110111010111011000000011010101^{60}0011..._2$$

Обозначение 1^{60} означает, что единица повторяется 60 раз. Таким образом, для получения корректно округленного значения нужно вычислять логарифм в этой точке с относительной погрешностью, не превосходящей 2^{-113} .

Описанные примеры показывают, что для выбора ближайшего числа с плавающей точкой иногда нужно использовать гораздо более точные вычисления, чем это позволяет сделать тип таких чисел. Аналогичные примеры существуют и для других режимов округления. Если используется режим округления к 0, к $+\infty$ или к $-\infty$, а точное значение функции лежит очень близко к представимому числу, необходимо добиться погрешности настолько маленькой, чтобы точно определить, превосходит оно это число или нет.

Например, натуральный логарифм для числа

$$110101100.01010000101101000000100111001000101011101110_2$$

равен

$$110.00001111010100101111001101111010111011001111100111^{61}0101..._2.$$

Для практически всех часто используемых функций их значения в «обычных» (т.е., не равных 0, 1 или 2) двоично-рациональных числах не являются рациональными, и поэтому не могут ни быть представимыми, ни лежать в точности посередине между двумя представимыми числами. Из этого в силу конечности множества представимых чисел следует, что для каждой функции есть такое число $\epsilon > 0$, что вычисляя значения этой функции с погрешностью, не превосходящей ϵ , можно всегда точно определить корректное округление, являющееся представимым числом. Однако вычислять функцию с такой точностью для всех значений аргумента может оказаться слишком неэффективно.

Например, для натурального логарифма на числах двойной точности при произвольном режиме округления такое ϵ можно взять равным 2^{-118} [17]. Однако реально такая точность нужна только для единственного значения аргумента, во всех остальных случаях можно использовать меньшую. Для подавляющего же большинства представимых чисел двойной точности корректное округление их логарифма можно получить, вычисляя его с погрешностью, не превосходящей 2^{-54} .

Дилемма составителя таблиц приводит к необходимости организации значительно более точных вычислений, чем это позволяют сделать стандартные типы чисел с плавающей точкой, как при построении правильных реализаций математических функций, так и проверке их корректности. В то же время, проводить настолько точные вычисления для всех значений аргументов очень неэффективно. Для повышения эффективности вычислений нужно уметь выбирать их точность в зависимости от аргументов функции.

3. Обзор работ по проверке корректности реализаций математических функций

Методам вычисления математических функций посвящено огромное количество работ. Одним из классических трудов на эту тему является сборник статей под редакцией Abramowitz и Stegun [18], хотя он был выпущен уже довольно давно и частично устарел. Более современное изложение методов вычисления элементарных функций (являющихся только подмножеством рассмотренных в [18]) можно найти в книге Muller [19].

Значительно реже встречаются исследования, в которых не только формулируются методы вычисления каких-либо функций, но и доказываются их корректность. Под корректностью имеется в виду достижение определенной точности результатов при определенных значениях параметров метода.

Только часть из этих работ посвящена вычислениям на числах с плавающей точкой, представление которых определяется в IEEE 754. Большинство из таких исследований связано с корректностью работы алгоритмов вычисления функций, реализованных в специализированном аппаратном обеспечении. Помимо точности вычислений здесь приходится учитывать разнообразные режимы округления, определяемые в IEEE 754, корректность выставляемых флагов, например, переполнения, а также правильность работы алгоритма на специальных значениях — бесконечностях и NaN (см. далее).

Гораздо реже встречаются исследования по строгому тестированию правильности реализаций математических функций, работающих с числами с плавающей точкой.

Таким образом, все имеющиеся работы по проверке правильности работы реализации математических функций можно разделить на следующие группы.

- Работы по формальной верификации определенных алгоритмов.
- Работы по тестированию реализаций математических функций.

3.1. Работы по формальной верификации определенных алгоритмов

Формальная верификация корректности алгоритмов вычислений математических функций чаще всего проводится при проектировании блоков вычислений с плавающей точкой универсальных процессоров или при разработке специализированных вычислительных процессоров [28,32,34-52,54].

Основой для такой верификации всегда является точное знание реализованного алгоритма и формализация основных требований стандарта IEEE 754, необходимая для их строгой проверки. Формализация этих требований проводилась несколько раз в различных формализмах: в языках формальных спецификаций Z [20], Barrett в 1989 [21], и VDM [22], Wichmann в том же году [23], в формализмах инструментов автоматизации доказательств Nurpl [24], PVS [25], HOL [26] и ACL2 [27] —

O'Leary с соавторами для Nurpl в 1994 [28], Miner для PVS в 1995 [29,30], Carreno для HOL в 1995 [30,31], Moore с соавторами для ACL2 в 1996 [32], Harrison для HOL в 1996 [33].

В работах Verkest с соавторами [34] и Cornea-Hagesan [35,36] все доказательства выполнялись вручную. Но в большей части исследований такого рода, ввиду значительной сложности верификации практически важных систем, применяются только формализации, которые могут использоваться инструментами для автоматизированного доказательства теорем.

Статьи [32,34,37-41] имеют дело только с алгоритмами умножения и/или деления чисел с плавающей точкой и представляют примеры верификации корректности таких алгоритмов по отношению к требованиям IEEE 754.

В работах [35,36,42-48] верифицируются также алгоритмы вычисления квадратного корня, иногда еще и вычисление остатка от деления для чисел с плавающей точкой и преобразования между различными такими типами и между ними и целыми числами.

В связи с тем, что в стандарте IEEE 754 определены только арифметические действия, вычисление остатка от деления и квадратного корня, проверке правильности вычисления других функций посвящено значительно меньше работ. В [49,50] описываются практические примеры формальной верификации алгоритмов, вычисляющих экспоненциальную функцию. В статье [51] верифицировалось вычисление синуса и косинуса, а в [52] — натурального логарифма. Все четыре работы были выполнены Harrison и его коллегами. В работе [53] нескольких исследователей из проекта Aenaire [12], представлен пример верификации вычисления натурального логарифма с помощью специализированного инструмента Garra, позволяющего автоматизировать доказательства, касающиеся свойств элементарных функций. Других примеров верификации математических функций в доступной литературе найти не удалось.

В перечисленных работах верифицировались отдельные элементы и блоки процессоров компаний Intel (Pentium II, Pentium III, Pentium 4 и Itanium [35,36,41-43,46,47,49-52]) AMD (K5, Athlon [32,44,45]), IBM (Power4 [48]). При этом использовались следующие инструменты автоматизации доказательств: HOL ([43,49-52]) и близкий к нему Nurpl ([29]), PVS ([38,40,54]), ACL2 ([32,44,45,48]) и его предыдущая версия Nqthm ([34]). Помимо этого применялись различные комбинации автоматизированного доказательства теорем с проверкой моделей (model checking) и символической проверкой эквивалентности моделей ([37,39,41,42,46,47]).

По-видимому, до сих пор не было удачных попыток полностью проверить соответствие блока вычислений с плавающей точкой стандарту IEEE 754. Во всех найденных работах либо проверяются не все операции, определяемые стандартом, либо не проверяется их работа для специальных чисел с плавающей точкой — денормализованных, -0, бесконечных и NaN. Работа [54], похоже, единственная, в которой систематически рассматривается поведение части описанных в IEEE 754 операций на таких исключительных значениях.

Можно также отметить, что в процессоре, верифицированном в проекте, описанном в [42] впоследствии была найдена ошибка в операции преобразования числа с плавающей точкой в целое (так называемый FIST bug [55]). Это показывает, что формальная верификация сложной системы, будучи тоже достаточно сложной деятельностью, из которой нельзя исключить участие людей, сама по себе подвержена ошибкам.

3.2. Работы по тестированию реализаций математических функций

В Интернет можно найти огромное количество различных программ для тестирования функций, работающих с числами с плавающей точкой, см. например, [56]. К сожалению, подавляющее большинство таких тестов крайне несистематично и проверяет какой-то один аспект вычислений, реже — два-три таких аспекта.

Как указывается в [56] (см. также иллюстрации из [57,58]), несмотря на то, что стандартизация вычислений с плавающей точкой началась около 20 лет назад, до сих пор многие поставщики библиотек и аппаратного обеспечения не придерживаются имеющихся стандартов достаточно строго, поэтому тесты на правильность поведения реализаций математических функций по-прежнему необходимы.

Среди наиболее систематичных работ по тестированию вычислений с плавающей точкой можно назвать следующие.

- **Работы по тестированию на соответствие стандарту IEEE 754.**

- В работе [59] описывается самый первый из известных систематических тестовых наборов для проверки корректности реализации операций над числами с плавающей точкой. Он появился еще до введения в действие стандарта IEEE 754, сделан в виде набора программ на Fortran и предназначен для тестирования только сложения, вычитания, умножения и деления.
- Специально для проверки на соответствие IEEE 754 был разработан тестовый набор, который описан в статьях [60,61] и может быть получен с сайта [62]. В этом наборе проверяются все требования стандарта к арифметическим операциям, вычислению квадратных корней и взятию остатков, а также преобразования между типами чисел с плавающей точкой и целыми.
- Программа PARANOIA [63,64] была создана одним из авторов стандарта IEEE 754 Кэханом (W. Kahan) и остается довольно популярным средством проверки на соответствие ему, хотя такая проверка менее тщательна, чем с помощью тестового набора, описанного выше. Она также проверяет только базовые арифметические операции.
- Другой подход к построению тестов для операций IEEE 754 используется в среде FPgen [65,66]. Здесь, помимо специальных значений, в качестве тестовых данных используются числа с плавающей точкой,

удовлетворяющие некоторым шаблонам — например, в которых нулевые и единичные биты мантиссы чередуются, или в которых мантисса содержит ровно 7 единиц.

- **Работы по тестированию широкого набора математических функций.**

- Тестовый набор UCBTTEST [67] предназначен для тестирования базовых арифметических действий и достаточно широкого набора математических функций. Он оформлен как набор программ на разных языках, включая Fortran и C, и наборов predetermined входных данных для разных функций. В каждом тесте проверяется, что для заданных значений параметров данная функция возвращает число с плавающей точкой, ближайшему к точному значению функции.
- Тестовый набор ELEFUNT [68], основанный на книге [69], также содержит тесты для многих математических функций в виде программ на C и Java и текстовых файлов с тестовыми данными и ожидаемыми результатами. Этот и предыдущий тестовый набор построены на основе проверки значений, возвращаемых реализациями функций для некоторых наборов аргументов. Методика выбора этих наборов аргументов, скорее всего, использовала несколько разных соображений.
 - Выделялись особые значения чисел с плавающей точкой: 0, -0, NaN +∞, -∞, минимальное положительное, максимальное положительное, числа, имеющие ровно один бит в мантиссе, и пр. (см. ниже).
 - Выделялись значения аргументов, значение функции для которых может быть точно представлено числом с плавающей точкой (например, $\exp(0) = 1$, $\cos(0) = 1$ и пр.).
 - Некоторые значения аргументов выбирались, по-видимому, случайно или из соображений, связанных со структурой известных алгоритмов для вычисления элементарных функций. Например, в ряде алгоритмов вычисления логарифма сначала значение аргумента при помощи умножения или деления на 2 приводится к интервалу (0.5, 1]. Соответственно, выбираются границы этого интервала и нескольких соседних с ним.
- Аналогичные подходы — использование ряда специальных значений, границ интервалов, определяемых часто используемыми алгоритмами вычисления данной функции, и случайных значений — применялись для построения более полных тестовых наборов, например, набора Беркли [70], а также в статье [71].

Отдельно стоит отметить работы в рамках проекта Arenalre [16,17,72], посвященные дилемме составителя таблиц и поиску чисел, для которых корректное вычисление функций с заданной точностью наиболее трудоемко. Эти числа можно использовать в

качестве «неудобных» тестовых значений для практически любой реализации соответствующей функции.

В рамках того же проекта был разработан инструмент MPCheck [73] для тестирования корректности реализации элементарных функций с точки зрения сохранения монотонности, симметрий, ограничений на область значений и корректности округления.

Доступная литература по исследованиям, посвященным проверке корректности реализации математических функций, показывает, что формальные спецификации не используются для разработки тестов для таких функций, хотя все предпосылки и технологические возможности для применения такого подхода имеются. Математические функции имеют четко определенное поведение, которое практически однозначно понимается разработчиками ПО и легко может быть зафиксировано в виде международных стандартов, поэтому разработка их спецификаций не является слишком трудоемкой задачей.

Использование формальных спецификаций также дает ряд преимуществ, связанных с повышением удобства сопровождения и снижением трудоемкости модификации больших тестовых наборов, по сравнению с традиционными методами разработки тестов. Все систематические тестовые наборы для тестирования библиотек математических функций достаточно объемны, поэтому построение тестов для них на основе формальных спецификаций способно реализовать эти преимущества.

4. Предлагаемый подход

Подход к тестированию реализаций математических функций, предлагаемый в данном исследовании, основан на технологии UniTESK, использующей формальные спецификации требований к программному обеспечению для автоматизированного построения тестов на соответствие им.

Основные элементы технологии UniTESK [1-3], применительно к реализациям математических функций, таковы.

- Требования к поведению тестируемой системы представляются в виде *формальных спецификаций*, которые состоят из следующих частей:
 - *предусловия*, описывающие области определения функций;
 - *постусловия*, описывающие условия корректности возвращаемых функциями результатов;
 - *инварианты* типов данных, описывающих условия целостности данных.
- Чтобы сделать спецификации независимыми от конкретной сигнатуры функции (и, возможно, от языка программирования, на котором она реализована), разрабатывается слой *адаптеров* или *медиаторов*, связывающих спецификации и реализации соответствующих функций друг с другом.

- Постусловие каждой функции анализируется с тем, чтобы выделить из него различные возможные варианты ее поведения. Они чаще всего соответствуют ветвлениям в теле постусловия и различным выражениям, описывающим ограничения на результат функции. Такие различные варианты поведения функции называются ее *функциональными ветвями*.

Набор функциональных ветвей определяет набор ситуаций, в которых реализация обязательно должна быть протестирована для того, чтобы проверить хотя бы один раз все выписанные ограничения.

- Помимо функциональных ветвей могут существовать другие ситуации, в которых тестирование данной реализации необходимо. Часть из этих ситуаций может быть получена более тонким анализом требований, а другая часть — на основе анализа возможных ошибок в конкретной реализации.

Такие ситуации также могут быть описаны в спецификациях с помощью специальных конструкций.

- На основе полученного набора ситуаций, в которых поведение реализации функции должно быть проверено, разрабатывается *тестовый сценарий*. Он определяет множество наборов аргументов функции, с которыми она будет вызываться во время тестирования. Это множество должно обеспечивать покрытие всех выделенных ситуаций, что контролируется при помощи автоматического построения отчетов о покрытии тестовых ситуаций, определенных в спецификации.

Возможности технологии UniTESK по автоматическому построению тестовой последовательности для тестирования математических функций не требуются. Единственное возможное исключение — наличие гипотез о зависимости работы реализаций математических функций от каких-то элементов внутреннего состояния тестируемой системы. В этом случае дополнительно должна быть построена обобщенная модель состояния, каким-то образом учитывающая те элементы, которые, как предполагается, могут влиять на работу функций.

Чтобы адекватно применять технологию UniTESK для построения тестов для математических функций, необходимо ответить на два вопроса: какие требования должны предъявляться к поведению их реализации, т.е. что именно должно быть написано в спецификациях; и на каких значениях параметров функций должно проводиться тестирование.

В рамках предлагаемого подхода были разработаны метод определения требований к реализации конкретной математической функции и метод выбора тестовых данных для тестирования конкретной функции.

4.1. Метод определения требований к математическим функциям

Данный метод определения требований к поведению реализаций математических функций заимствует часть идей из стандарта ISO 10967 [13-15] и работ [10,11], посвященных вопросу разработки серии стандартов с повышенными требованиями к корректности вычисления математических функций и удовлетворяющих им библиотек. Некоторые элементы предлагаемого метода являются новыми и не встречаются в доступной литературе.

Требования к реализации математической функции могут быть разделены на несколько аспектов, которые должны быть рассматриваться отдельно.

• Область определения функции и особые точки функции.

- Для всех значений аргументов, где математическая функция определена, ее реализация должна возвращать некоторый результат, который может быть равен $+\infty$ или $-\infty$, если значение самой функции находится за пределами интервала чисел с плавающей точкой, но не должен быть NaN.
- Для всех значений, для которых математическая функция не определена (в том числе, и для ее особых точек), но имеет однозначно определенный предел, может быть, равный $+\infty$ или $-\infty$, реализация должна возвращать значение этого предела.
- Если функция имеет особенность в точке 0, не имеет там предела, равного $+\infty$ или $-\infty$, нужно рассматривать односторонние пределы функции. Значение реализации функции в 0 нужно считать равным ее пределу при $x \rightarrow +0$, если он существует, а значение в -0 — пределу при $x \rightarrow -0$, если он есть. Примером такой функции служит котангенс.
- В остальных случаях должен возвращаться результат NaN. Особо нужно рассматривать такие значения аргументов, по поводу которых нет однозначного мнения о принадлежности их к области определения функции или о возможном продолжении ее в эту точку по непрерывности. Примером служит значение 0^0 , которое иногда интерпретируется как 1, а иногда как NaN.
- Для полюсов функции, где ее значение стремится к бесконечности, необходимо точно определить окрестности, в которых оно уже не является представимым. При наличии представимых чисел в такой окрестности, реализация функции для них должна возвращать значения $+\infty$ или $-\infty$ в соответствии со знаком точного значения.
- Для функций, стремящихся к бесконечности при $x \rightarrow +\infty$ или $x \rightarrow -\infty$, должны быть точно определены пределы представимости их значений. За этим пределами реализация также должна возвращать $+\infty$ или $-\infty$ в соответствии со знаком точного значения. Например, для значений $x > \ln(2^{104} \cdot (2^{24} - 1)) = 88,722839052\dots$ $\exp(x)$

не попадает в диапазон чисел однократной точности, поэтому реализация экспоненты для таких чисел должна возвращать $+\infty$.

• Специальные значения, значения в 0, касательные и асимптоты.

- Нужно наиболее естественным образом определить значения функции для особых значений аргумента: -0 , $+\infty$, $-\infty$. Обычно достаточно определять их как пределы, если те существуют, иначе как NaN.
- Значение функции для значения аргумента NaN должно быть равно NaN.
- Для некоторых значений аргумента значения функции известны точно. Если оба значения представимы, естественно потребовать от реализации этой функции возвращать именно точное ее значение в таких точках. Например, $\exp(0) = \cos(0) = \operatorname{ch}(0) = 1$, $\sin(0) = \tan(0) = \operatorname{arcsin}(0) = 0$ и т.п.
- Кроме этого, если в такой точке производная функции равна 0, то для любого аргумента из некоторой ее окрестности реализация должна возвращать то же самое значение.
- Так как около 0 плотность чисел с плавающей точкой больше, чем около любого другого значения, если значение функции в 0 не равно 0, даже если ее производная там ненулевая, должно быть выполнено то же самое правило: в некоторой окрестности 0 для всех чисел с плавающей точкой значение ее реализации должно быть одинаковым. Например, для экспоненты $e^0 = 1$, при этом $(e^x - 1) < 2^{-24}$ при $x > 1$ и $(x-1) < \ln(1 + 2^{-24}) = 5.9604642999\dots \cdot 10^{-8}$, поэтому для всех таких x , которых довольно много, реализация экспоненты с однократной точностью должна возвращать 1.
- В тех случаях, когда функция имеет горизонтальные асимптоты, необходимо аккуратно определить границы, после которых ее значение должно стать постоянным. Например, для значений $x < \ln(2^{-150}) = -103.97207708399\dots$ значение e^x становится ближе к 0, чем к какому либо еще числу, представимому с однократной точностью. Поэтому реализация экспоненты для таких аргументов должна возвращать 0.
- Казалось бы, естественно предъявить аналогичные требования к функциям, имеющим негоризонтальные асимптоты или асимптотически близких к другим функциям. Например, $\operatorname{ch}(x) = (e^x + e^{-x})/2 \sim e^x$ для достаточно больших значений аргумента, или $\sin(x) \sim x$ при $x \sim 0$. Однако во многих случаях такое требование не может быть сформулировано достаточно аккуратно с учетом различных режимов округления, зафиксированных в IEEE 754. Дело в том, что даже очень маленькая разность между двумя асимптотически близкими выражениями может дать отличие в

представимы,— это свойства $\Gamma(1+x) = x\Gamma(x)$ и $\Gamma(1-x) = -x\Gamma(-x)$ гамма-функции $\Gamma(x)$. Иначе может нарушиться важное соотношение $\Gamma(n) = (n-1)!$ для целых положительных n .

- Для свойства периодичности функции остаются верными все те же аргументы. Если период представим, то можно проверять это свойство только для чисел, представимых вместе со своим сдвигом на число, кратное периоду. Если период не представим, реализация может быть только приблизительно периодична.

Для значений аргумента, у которых единица последнего разряда больше, чем величина периода функции, проверка ее периодичности становится практически бесполезной.

- **Корректное округление.**

Помимо всех перечисленных ограничений, нужно потребовать, чтобы результат, возвращаемый реализацией, получался из точного результата функции для данного аргумента при помощи принятой в текущей конфигурации процедуры округления. Иногда практически бессмысленно требовать точности 0.5-1 ulp, но во всяком случае погрешность выше 1.5-2.0 ulp должна рассматриваться как неточность соответствующей реализации.

Стандарт IEEE 754 предписывает поддержку 4-х видов округления: к ближайшему представимому числу, к $+\infty$, к $-\infty$ и к 0.

Часто три последних вида округления противоречат требованиям, сформулированным по остальным аспектам. В этом случае иногда можно принимать решение в зависимости от функции, потому что некоторые ее важные свойства могут нарушиться. Но чаще удобнее считать, что поддержка режима округления имеет более высокий приоритет, поскольку пользователь, применяющий такой режим, осведомлен о его последствиях.

После определения требований они оформляются в виде формальных спецификаций. Для записи требований, касающихся корректного округления точного значения математической функции, необходимо уметь вычислять ее правильно округленные значения. Сделать это можно несколькими способами.

- С помощью систем математических вычислений, например, Maple [74], Mathematica [75], MATLAB [76].
- С помощью библиотек корректно округляемых функций, таких, как разработанная на основе работ Ziv [77] IBM Accurate Portable MathLib [78], GNU MPFR [79], разрабатываемая Sun libmcr [80] или библиотек SCSLib [81] и CRlibm [82], разрабатываемых в рамках проекта Arenaire [16,83-86].
- Можно также разработать собственную реализацию функции на основе методов,

изложенных в книгах [18,19] и многочисленных статьях, или на основе методов интервальных вычислений [87-90].

4.2. Метод выбора тестовых данных

Предлагаемая методика выбора тестовых данных для тестирования математических функций основана на особенностях представление чисел с плавающей точкой, результатах работ [16,17,72] по вычислению «неудобных» значений аргументов и приведенном выше методе определения требований к реализациям таких функций, а также на технике построения тестов при помощи разбиения интервалов входных данных.

Основные шаги этого метода состоят в следующем.

- При определении требований к реализации функции по описанной выше методике числа с плавающей точкой разбиваются на ряд интервалов, в рамках каждого из которых действуют свои собственные ограничения. Это, например, интервалы монотонности и сохранения знака, интервалы, на которых функция имеет постоянное значение, включая бесконечные, интервалы, на которых она не определена. Множество концов этих интервалов будем называть *исходным множеством*.
- К исходному множеству добавляем числа 0, -0, $+\infty$, $-\infty$, минимальные и максимальные по абсолютной величине представимые числа, минимальное и максимальное денормализованные числа.
- Полученное исходное множество разбивает числа с плавающей точкой на набор интервалов.

Для выбранных числовых параметров n и k на каждом из этих интервалов возможные тестовые значения выбираются следующим образом. Сначала интервал разбивается на n более мелких интервалов, равных по количеству содержащихся в них чисел с плавающей точкой. При этом возникает $(n+1)$ точка. Затем берутся все числа, лежащие в рассматриваемом интервале и отстоящие не более чем на k чисел с плавающей точкой от полученных точек. Получаемое множество точек назовем *пробным множеством*.

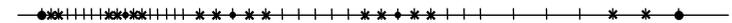


Рисунок 5. Схема выбора тестовых данных на одном интервале с параметрами $n=4$, $k=2$.

- К пробному множеству необходимо добавить ряд значений, которые приводят к необходимости гораздо более точного вычисления значения рассматриваемой функции для них (см. выше о дилемме составителя таблиц и [17]). Кроме того,

$29 \cdot \pi/2 =$ 101101.10001101100101111000100010111011100100011011101 $1^{13}010 \dots_2$ $tg(1.01101.10001101100101111000100010111011100100011011110_2) =$ $-1.0110011010111001111010111100010010000101000011000110$ $0110 \dots_2 \cdot 2^{60}$
$9206271 \cdot \pi/2 =$ 110111001010100011111000.10101011100101110101110100101 $0^{30}111 \dots_2$ $tg(1.10111001010100011111000.10101011100101110101110100101_2) =$ $1.0000010101110101100001001100010000101001101100111010$ $0110 \dots_2 \cdot 2^{59}$
$138049179777104367775 \cdot \pi/2 =$ 1.0111100000101011011110100010000011011111011011010100 $0^{72}101 \dots_2 \cdot 2^{67}$ $tg(1.011110000010101101111010001000001101111011011010100_2 \cdot 2^{67}) =$ $1.0101101001100011110010000010101101101101011100010$ $0100 \dots_2 \cdot 2^{57}$
$1104381301317041933816362171 \cdot \pi/2 =$ 1.0110011010111101010101000010010011100101011001010101 $0^{95}101 \dots_2 \cdot 2^{90}$ $tg(1.01100110101111010101010000100100111001010110010101_2 \cdot 2^{90}) =$ $1.01011010010111100010101110001000011010100010101010$ $1^3010 \dots_2 \cdot 2^{57}$
$2276647626870351330792862584889027286951 \cdot \pi/2 =$ 1.010100000100110010101100010100011111000111010101111 $0^{137}100 \dots_2 \cdot 2^{131}$ $tg(1.01010000010011001010110001010001111100011101010111_2 \cdot 2^{131}) =$ $1.111001111100011111010000111101000011111100000011101$ $10^511 \dots_2 \cdot 2^{58}$
$9617272285741355388977328677436746192449293 \cdot \pi/2 =$ 1.0101101011010101101001100010110010110001110011001000 $1^{147}011 \dots_2 \cdot 2^{143}$ $tg(1.010110101101010110100110001011001011000111001100100_2 \cdot 2^{143}) =$ $-1.110001111010101101110011100010011100001001001101011$ $0110 \dots_2 \cdot 2^{56}$
$3982124097056689019860931670115672371080901279195265295 \cdot \pi/2 =$ 1.0000010100111001101101001000110100010100110001010101 $0^{186}100 \dots_2 \cdot 2^{182}$ $tg(1.00000101001110011011010010001101000101001100010101_2 \cdot 2^{182}) =$ $1.1010011111001011011111001011010000101000000010110000$ $0^3110 \dots_2 \cdot 2^{56}$
$243709684824005521714504726240306266173066571317974432621719 \cdot \pi/2 =$ 1.1110011111100100010010100111100010101100000110001011 $1^{203}000 \dots_2 \cdot 2^{197}$ $tg(1.11100111110010001001010011110001010110000011000100_2 \cdot 2^{197}) =$ $-1.0000100110111010011111100000011000011111001100000000$ $1001 \dots_2 \cdot 2^{58}$
$10177187664528444025108553957230483888384395218254373185418290015777$ $4575619 \cdot \pi/2 =$ 1.0110100111101010101100001001100001010001011110011011 $0^{251}101 \dots_2 \cdot 2^{246}$ $tg(1.0110100111101010101100001001100001010001011110011011_2 \cdot 2^{246}) =$ $1.011100100101011001111100101110010000100011110101000$ $1^3001 \dots_2 \cdot 2^{57}$

15624927791716915161759873344092142005165621843671972171731789642962 $257024455 \cdot \pi/2 =$ 1.1011001000011001011000110110010011010111010100001010 $1^{258}011 \dots_2 \cdot 2^{253}$ $tg(1.1011001000011001011000110110010011010111010100001011_2 \cdot 2^{253}) =$ $-1.1101110100010110100010001100011010010110100110101100$ $1^3011 \dots_2 \cdot 2^{57}$
71610307166660314043397602176728270128783849452911391884618409783207 $463478073649776722681 \cdot \pi/2 =$ 1.110001000101110011010001000100010101010011011111101 $0^{300}100 \dots_2 \cdot 2^{295}$ $tg(1.110001000101110011010001000100010101010011011111101_2 \cdot 2^{295}) =$ $1.11111101101110100100011011110101010011101000001000$ $10^310 \dots_2 \cdot 2^{57}$
79531021962858779827512141402968504124268502489089591753854287169082 $2757973300060931533046037509210740152847560365069332549 \cdot \pi/2 =$ 1.1110001111001010100110110110110001100101010111001011 $0^{412}100 \dots_2 \cdot 2^{408}$ $tg(1.111000111100101010011011011011000110010101011100101_2 \cdot 2^{408}) =$ $1.1001111110011110010101010001001001100110110000111100$ $01^300 \dots_2 \cdot 2^{56}$
43776696795279815301055904915945700987301503939761311288501927384592 $34187554875937911988865352801115519693682101780670312164186360396953$ $23195218767 \cdot \pi/2 =$ 1.1011100010001100101110110100111000110010010101110101 $1^{491}011 \dots_2 \cdot 2^{487}$ $tg(1.1011100010001100101110110100111000110010010101110110_2 \cdot 2^{487}) =$ $-1.1010001011000010110011001011110000011010011000110010$ $0101 \dots_2 \cdot 2^{56}$
11589330310381988626449129598171351352584091074116125961273394956149 $19456683050991296076480489809564268564293673421768988424019510581475$ $74305293329168873163869152174973 \cdot \pi/2 =$ 1.1000101100101000011001110110110011011100110001011010 $1^{560}000 \dots_2 \cdot 2^{555}$ $tg(1.100010110010100001100111011011001101110011000101101_2 \cdot 2^{555}) =$ $-1.000010100001000101100110110100011001010011011100010$ $0101 \dots_2 \cdot 2^{57}$
40912440333305918481550440392603355393318151657874369033644916645778 $80843430257120834658307295698261362420794202530508444695188784113250$ $59263543780254873744888889331055895843 \cdot \pi/2 =$ 1.0100110010010110110000010001000100110100110100110101 $1^{583}010 \dots_2 \cdot 2^{577}$ $tg(1.0100110010010110110000010001000100110100110100110110_2 \cdot 2^{577}) =$ $-1.01100101010111001111001111000100011100010011010111$ $0^3101 \dots_2 \cdot 2^{58}$
65429133524241263178237370719527517029251219555784828550385052067786 $68733143053622851340353340621053795171428918994903139424630388205704$ $7519616430579116455387138549345777195516511603705 \cdot \pi/2 =$ 1.1000001100000000100111100010111010011110001011101010 $1^{619}011 \dots_2 \cdot 2^{614}$ $tg(1.100000110000000010011110001011101001111000101110101_2 \cdot 2^{614}) =$ $-1.111000010110001000011111101100010011001101010100100$ $1^2011 \dots_2 \cdot 2^{57}$

<p>59287932212270746098210389073106720676473902052684288108187468991207 27939955976870793398691724674236683942417720739021351535557283241575 82755209798363213599827165644437783394500726283264184124810624461674 $9\cdot\pi/2 =$ 1.110110110100001111100111100101101110001110101111011 $0^{686}111\dots_2^{680}$ $\operatorname{tg}(1.110110110100001111100111100101101110001110101111011_2^{-680}) =$ 1.000100001011101001101111001111101000110011111010110 $1^2010\dots_2^{58}$</p>
<p>30636572260931163506332425320948767147146434821312896452845246587739 4299960543512781934069734059379292127565945745701814493759190791305 15095712133324352639295031679514263477241946753392456572308669671533 6885$\cdot\pi/2 =$ 1.110111110101000110100011000111100101011001111101101 $1^{693}011\dots_2^{689}$ $\operatorname{tg}(1.11011111010100011010001100011110010101100111110110_2^{-689}) =$ -1.10111001111100000011001101011010001011010111100111 $10^511\dots_2^{56}$</p>
<p>94971440572386015595571166139691201230707691489980939180815478595748 85577578104219175989760881732308451273779191276410784862786933003137 63136525734828774533668148005758919334015282384829230967867833601055 56253079245416317281232917440947587$\cdot\pi/2 =$ 1.0110111010001101011101111000110010010100110101100110 $0^{798}100\dots_2^{794}$ $\operatorname{tg}(1.0110111010001101011101111000110010010100110101100110_2^{-794}) =$ 1.10011100111010001001101111100111010110100011111100 $0101\dots_2^{56}$</p>
<p>33864178045159811206438920823311565991202393932998380352421215184285 37554064774221620930267583474709602068045686026362989271814411863708 49986972132271594662263430201169763297290792255889271083061603403854 1342154669787134871905353772776431251615694251273653$\cdot\pi/2 =$ 1.011010101100010110110010011000101100101000011111110 $1^{857}011\dots_2^{849}$ $\operatorname{tg}(1.011010101100010110110010011000101100101000011111111_2^{-849}) =$ -1.1101100110111010100110100111100101110101011000110101 $1010\dots_2^{60}$</p>
<p>88685758294706272204189579187820280422409877285489361868254170927334 81411325935922524369633276190848778076190586339609015909756150954578 31282700320668157612019880948068064298959801704400266492369256967459 7153720925148676684943087665815389546838546388224593403$\cdot\pi/2 =$ 1.11001111110010010000010001010000101111100011101101 $0^{866}111\dots_2^{860}$ $\operatorname{tg}(1.11001111110010010000010001010000101111100011101101_2^{-860}) =$ 1.000110100010010001101011011011110010010001010110 $1^2000\dots_2^{58}$</p>

<p>37105172474923906782576160374346777164937755101941743463232913447611 26458820318029413443054532013548715976758198761455677008210344281793 85754829765785467791950462488227489950141733334740102736828752663860 04962920828212934138031374382514718179363737540556512625927642564615 90125145265$\cdot\pi/2 =$ 1.010000010001011101010111001100111001011110101000001 $1^{944}000\dots_2^{939}$ $\operatorname{tg}(1.010000010001011101010111001100111001011110101000010_2^{-939}) =$ -1.000111000111001010010100101011101110011101110100100 $01^401\dots_2^{57}$</p>
<p>22795463807341151743454041517545705884702298848407846038370271485088 07014891141753237567442204768461160038124562964557539085826195999224 02145326867889285827944189828230765840337018227362357217203013966072 70621803425566565987094025434998287231727147847848195901676458199728 364279528014629$\cdot\pi/2 =$ 1.111000011001100001110001001000101011011111000000110 $0^{955}101\dots_2^{951}$ $\operatorname{tg}(1.111000011001100001110001001000101011011111000000110_2^{-951}) =$ 1.011100010111010000000111101001000000001111001100110 $1^3011\dots_2^{56}$</p>
<p>24982829712749183769184949361370722806939675044978450598655040276224 2206932656928540059089298879071869587152103490382783870770007392787 73724242538703521095769835440603290627555914001656747466652577114915 40415860945893559496946376666124169735622128694086982311970782271594 626856299313449263532321973$\cdot\pi/2 =$ 1.11100000000100111000101001100010100100010111100001 $0^{997}100\dots_2^{991}$ $\operatorname{tg}(1.11100000000100111000101001100010100100010111100001_2^{-991}) =$ 1.10111000110010111110111110100111110101111111010 $0^2111\dots_2^{58}$</p>
<p>19761831728800425220636434033446152324993825185277180432073766513146 26916965672712973121611642165275685950124449367411073847019492715148 94825117746798310262428988560159244225487630579194650628623442695547 95169108098262423683067093402882407403932831289539148165440830110821 552783846121469368255050296717288515$\cdot\pi/2 =$ 1.01100001101000111101101110001100100011010001001001 $0^{1026}111\dots_2^{1021}$ $\operatorname{tg}(1.111000000000100111000101001100010100100010111100001_2^{-1021}) =$ 1.000100101011110000010000100101010100001110011010100 $10^311\dots_2^{57}$</p>

Таблица 3. Нечетные кратные $\pi/2$, ближайшие к числам с плавающей точкой и соответствующие значения тангенса.

В 20-м ряду Таблицы 3 приведено число с плавающей точкой, ближайшее к нечетному кратному $\pi/2$. Это же число упоминается в статье [93], там оно вычислено с помощью программы, созданной W. Kahan и S. McDonald [94].

Из этого следует, что приведенное в 20-м ряду значение тангенса является максимальным по абсолютной величине. Таким образом, для всех чисел с

<p>65429133524241263178237370719527517029251219555784828550385052067786 68733143053622851340353340621053795171428918994903139424630388205704 7519616430579116455387138549345777195516511603705$\pi =$ 1.100000110000000010011110001011101001110001011101010 $1^{619}011\dots_2^{615}$ $tg(1.100000110000000010011110001011101001110001011101011_2^{615}) =$ 1.00010000010010000011000001000000000010010001001001 $1011\dots_2^{-57}$</p>
<p>59287932212270746098210389073106720676473902052684288108187468991207 27939955976870793398691724674236683942417720739021351535557283241575 82755209798363213599827165644437783394500726283264184124810624461674 $9\pi =$ 1.1101101101000001111100111100101101110001110101111011 $0^{686}111\dots_2^{681}$ $tg(1.110110110100000111110011110010110111000111010111011_2^{681}) =$ <math>-1.111000001001100001111000101001000101110001010011101 $0^2101\dots_2^{-58}$</math></p>
<p>33864178045159811206438920823311565991202393932998380352421215184285 37554064774221620930267583474709602068045686026362989271814411863708 49986972132271594662263430201169763297290792255889271083061603403854 1342154669787134871905353772776431251615694251273653$\pi =$ 1.0110101011000101101100100110001011001010000111111110 $1^{857}011\dots_2^{850}$ $tg(1.011010101100010110110010011000101100101000011111111_2^{850}) =$ 1.0001010010101110011100101110011010111010001000101110 $1^4010\dots_2^{-60}$</p>
<p>88685758294706272204189579187820280422409877285489361868254170927334 81411325935922524369633276190848778076190586339609015909756150954578 3128270032066815761201988094806806429895980170440026649236256967459 7153720925148676684943087665815389546838546388224593403$\pi =$ 1.110011111100100100000100010100001011111100011101101 $0^{866}111\dots_2^{861}$ $tg(1.110011111100100100000100010100001011111100011101101_2^{861}) =$ <math>-1.1101000010001111010110011100001100101011111110100010 $1^2010\dots_2^{-58}$</math></p>
<p>24982829712749183769184949361370722806939675044978450598655040276224 22069326569285540059089298879071869587152103490382783870770007392787 73724242538703521095769835440603290627555914001656747466652577114915 4041586094589355949646376666124169735622128694086982311970782271594 626856299313449263532321973$\pi =$ 1.1110000000001001110001010011000101001000101111100001 $0^{997}100\dots_2^{992}$ $tg(1.1110000000001001110001010011000101001000101111100001_2^{992}) =$ <math>-1.0010100101011010001110110000101001100100101100011101 $0100\dots_2^{-58}$</math></p>

Таблица 6. Кратные π , ближайшие к числам с плавающей точкой и соответствующие значения тангенса.

Отметим, что ближайшее к числу с плавающей точкой кратное π число, находящееся в 11-м ряду Таблицы 6, соответствует точно такому же целому

числу, что и ближайшее к числу с плавающей точкой нечетное кратное $\pi/2$.

• **Симметрии и периодичность.**

Тангенс — нечетная функция, $tg(-x) = -tg(x)$. Это свойство должно учитываться при построении всех тестов с отрицательными значениями аргумента. Периодичность с трансцендентным периодом π , как и симметрии вида $tg(\pi/2-x) = 1/tg(x)$, выполняются на числах с плавающей точкой лишь приблизительно и, поэтому не могут быть использованы для формулировки строгих требований.

• **Корректное округление.**

Погрешность вычисления тангенса согласно стандарту ISO/IEC 10967-2 [14] не должна превосходить 0.5-2.0 ulp. В качестве дополнительной информации можно выдавать предупреждение специального вида каждый раз, когда фактическая ошибка превышает 0.5 ulp.

Для вычисления правильно округленного значения тангенса можно использовать известные библиотеки вычислений с заданной точностью [79,81,82].

• **Тестовые значения.**

Идеальный набор интервалов для дальнейшего построения тестовых данных для тестирования реализаций тангенса мог бы состоять из всех его периодов, в которые попадает хотя бы два числа с плавающей точкой. Однако, таких интервалов слишком много — около $5.7 \cdot 10^{15}$. Поэтому необходимо выбрать только некоторые из таких периодов, руководствуясь каким-то правилом.

Предлагается в качестве основы для тестов выбрать периоды тангенса, содержащие хотя бы одну из упоминавшихся выше точек, близких к кратным π , $\pi/2$ и $\pi/4$ (включая 0). Таких точек меньше, чем 3000, поэтому получающийся набор тестовых данных будет одновременно достаточно представительным и обозримым. Кроме того, к ним можно добавить несколько периодов, близких к содержащему 0, например, от -29π до 29π (раз уж число 29 дает одно из самых близких кратных π к числу с плавающей точкой). Для каждого из выбранных интервалов нужно вычислить его левый и правый концы — точки, ближайшие к нечетным кратным $\pi/2$ на этом интервале, его точку, ближайшую к кратному π , а также две точки на этом интервале, ближайшие к нечетным кратным $\pi/4$. Эти точки разбивают каждый такой период на четыре под-интервала. Все их занесем в исходное множество.

Кроме того, в исходное множество попадут точки $x_1, -x_1, x_2, -x_2, x_3, -x_3, x_4, -x_4$, вычисленные при рассмотрении асимптотики тангенса в 0, сам 0, $-0, +\infty$ и $-\infty$, а также минимальное и максимальное положительные денормализованные числа и противоположные им.

- [19] J.-M. Muller. *Elementary Functions: Algorithms and Implementation*. Second edition. Birkhauser, Boston, 2006.
- [20] <http://www.zuser.org/z/>
- [21] G. Barrett. *Formal Methods Applied to a Floating-Point Number System*. IEEE Transactions on Software Engineering, 15(5):611–621, May 1989.
- [22] <http://www.csr.ncl.ac.uk/vdm/>
- [23] B. A. Wichmann. *Towards a Formal Specification of Floating Point*. The Computer Journal, 32:432–436, 1989.
- [24] <http://www.cs.cornell.edu/Info/Projects/NuPrI/nuprl.html>
- [25] <http://pvs.csl.sri.com/>
- [26] <http://www.cl.cam.ac.uk/Research/HVG/HOL/>
- [27] <http://www.cs.utexas.edu/users/moore/acl2/acl2-doc.html>
- [28] J. O’Leary, M. Leaser, J. Hickey, and M. Aagaard. *Non-restoring Integer Square Root: A Case Study in Design by Principled Optimization*. In T. Kropf, R. Kumar, eds. Proc. of the 2-nd International Conference on Theorem Provers in Circuit Design (TPCD’94): Theory, Practice, and Experience. LNCS 901, pp. 52–71, Springer-Verlag, 1994.
- [29] P. S. Miner. *Defining the IEEE-854 Floating-Point Standard in PVS*. Technical report TM-110167, NASA Langley Research Center, 1995.
- [30] V. A. Carreno and P. S. Miner. *Specification of the IEEE-854 Floating-Point Standard in HOL and PVS*. Proc. of 10-th International Workshop on Higher Order Logic Theorem Proving and its Applications, Aspen Grove, Utah, USA, September 1995.
- [31] V. A. Carreno. *Interpretation of IEEE-854 Floating-Point Standard and Definition in the HOL System*. Technical report TM-110189, NASA Langley Research Center, 1995.
- [32] J. Moore, T. Lynch, and M. Kaufmann. *A Mechanically Checked Proof of the Correctness of the Kernel of the AMD5K86 Floating-Point Division Algorithm*. IEEE Transactions on Computers, 47(9):913–926, September 1998.
- [33] J. Harrison. *Theorem Proving with the Real Numbers*. Technical report UCAM-CL-TR-408, University of Cambridge Computer Laboratory. UK, November 1996.
- [34] D. Verkest, L. Claesen, and H. De Man. *A Proof on the Nonrestoring Division Algorithm and its Implementation on an ALU*. Formal Methods in System Design, vol. 4, 1994.
- [35] M. Cornea-Hasegan. *Proving the IEEE Correctness of Iterative Floating-Point Square Root, Divide, and Remainder Algorithms*. Intel Technology Journal, 1998.
- [36] M. Cornea-Hasegan. *IA-64 Floating Point Operations and the IEEE Standard for Binary Floating-Point Arithmetic*. Intel Technology Journal, Q4, 1999.
- [37] M. D. Aagaard and C.-J. H. Seger. *The Formal Verification of a Pipelined Double-precision IEEE Floating-Point Multiplier*. Proc. ICCAD, IEEE, Nov. 1995, pp. 7–10.
- [38] P. S. Miner, and J. F. Leathrum. *Verification of IEEE Compliant Subtractive Division Algorithms*. Proc. FMCAD’96, November 1996.
- [39] E. M. Clarke, S. M. German, and X. Zhao. *Verifying the SRT Division Algorithm Using Theorem Proving Techniques*. Proc. CAV’96, LNCS 1102, Springer-Verlag, 1996.
- [40] H. Ruess, N. Shankar, and M. K. Srivas. *Modular Verification of SRT Division*. Proc. CAV’96, LNCS 1102, Springer-Verlag, 1996.
- [41] R. Kaivola and N. Narasimhan. *Formal Verification of the Pentium 4 Floating-Point Multiplier*. Proc. of Design, Automation and Test in Europe Conference and Exposition (DATE). IEEE, 2002.
- [42] Y.-A. Chen, E. M. Clarke, P.-H. Ho, Y. Hoskote, T. Kam, M. Khaira, J. W. O’Leary, and X. Zhao. *Verification of All Circuits in a Floating-Point Unit Using Word-level Model Checking*. In Formal Methods in Computer-Aided Design, LNCS 1166, pp. 19-33, Springer-Verlag, 1996.
- [43] J. R. Harrison. *Verifying the Accuracy of Polynomial Approximations in HOL*. TPHOLS’97, August 1997.
- [44] D. M. Russinoff. *A Mechanically Checked Proof of IEEE Compliance of the Floating Point Multiplication, Division and Square Root Algorithms of the AMD-K7 Processor*. LMS Journal of Computation and Mathematics, 1:148–200, 1998.
- [45] D. M. Russinoff. *A Mechanically Checked Proof of IEEE Compliance of the AMD K5 Floating-Point Square Root Microcode*. Formal Methods in System Design, 14 (1), January 1999.
- [46] J. O’Leary, X. Zhao, R. Gerth, and C. H. Seger. *Formally Verifying IEEE Compliance of Floating-Point Hardware*. Intel Technology Journal, 1999.
- [47] M. D. Aagaard, R. B. Jones, and R. Kaivola. *Formal Verification of Iterative Algorithms in Microprocessors*. Proc. DAC 2000. ACM, 2000.
- [48] J. Sawada. *Formal Verification of Divide and Square Root Algorithms Using Series Calculations*. Proc. of ACL2 Workshop, Grenoble, France, April 2002.
- [49] J. R. Harrison. *Floating-Point Verification in HOL Light: the Exponential Function*. Technical Report UCAM-CL-TR-428, University of Cambridge Computer Laboratory. UK, June 1997.
- [50] A. T. Abdel-Hamid. *A Hierarchical Verification of the IEEE-754 Table-driven Floating-point Exponential Function Using HOL*. Master’s thesis, Dept. Electrical and Computer Engineering, Concordia University, Montreal, Quebec, Canada, 2001.
- [51] J. Harrison. *Formal Verification of Floating Point Trigonometric Functions*. Proc. of Formal Methods in Computer-Aided Design, FMCAD 2000, LNCS 1954, pp. 217-233, Springer-Verlag, 2000.
- [52] J. Harrison. *Floating Point Verification in HOL*. E. T. Schubert, P. J. Windley, and J. Alves-Foss, eds. Proc. of 8-th International Workshop on Higher Logic Theorem Proving and its Applications, Aspen Grove, UT, USA, September 1995. LNCS 971, pp. 186–199, Springer-Verlag, 1995.
- [53] F. de Dinechin, C. Lauter, and G. Melquiond. *Assisted Verification of Elementary Functions*. INRIA Research Report RR-5683, September 2005.
- [54] C. Jacob, C. Berg. *Formal Verification of the VAMP Floating Point Unit*. In Formal Methods in System Design, 26, pp. 227-266, Springer, 2005.
- [55] <http://support.intel.com/support/processors/flag/tech.htm>, *Discussion of Flag Erratum*, 2002.
- [56] <http://www.math.utah.edu/~beebe/software/ieee/>
- [57] <http://people.redhat.com/drepper/libm/>
- [58] W. Kahan. *What Can You Learn about Floating-Point Arithmetic in One Hour?* <http://http.cs.berkeley.edu/~wkahan/ieee754status>, 1996.
- [59] N. L. Schryer. *A Test of Computer’s Floating-Point Arithmetic Unit*. Computer Science Technical Report 89, AT&T Bell Labs, 1981.
- [60] B. Verdonk, A. Cuyt, and D. Verschaeren. *A Precision- and Range-Independent Tool for Testing Floating-Point Arithmetic I: Basic Operations, Square Root and Remainder*. ACM TOMS 27(1):92–118, 2001.
- [61] B. Verdonk, A. Cuyt, and D. Verschaeren. *A Precision- and Range-Independent Tool for Testing Floating-Point Arithmetic II: Conversions*. ACM TOMS 27(1):119–140, 2001.
- [62] <http://www.cant.ua.ac.be/ieecc754.html>
- [63] R. Karpinski. *PARANOIA: A Floating-Point Benchmark*. Byte Magazine 10, 2 (Feb.), pp. 223–235, 1985.
- [64] <http://www.netlib.org/paranoia/>

- [65] A. Ziv, M. Aharoni, and S. Asaf. *Solving Range Constraints for Binary Floating-Point Instructions*. Proc. of 16-th IEEE Symposium on Computer Arithmetic (ARITH-16'03), pp. 158–163, 2003.
- [66] M. Aharoni, S. Asaf, L. Fournier, A. Koifman, and R. Nagel. *FPgen — A Test Generation Framework for Datapath Floating-Point Verification*. Proc. IEEE International High Level Design Validation and Test Workshop (HLDVT'03), pp. 17–22, 2003.
- [67] <http://www.netlib.org/fp/ucbtest.tgz>
- [68] <http://www.math.utah.edu/pub/elefun/>
- [69] W. Cody and W. Waite. *Software Manual for the Elementary Functions*. Prentice-Hall, Englewood Cliffs, NJ, 1980.
- [70] Z. A. Liu. *Berkeley Elementary Function Test Suite*. M.S. thesis, Computer Science Division, Dept. of Electrical Engineering and Computer Science, University of California at Berkeley, December 1987.
- [71] P.-T. P. Tang. *Accurate and Efficient Testing of the Exponential and Logarithm Functions*. ACM Transactions on Mathematical Software, 16(3):185–200, September 1990.
- [72] D. Stehle, V. Lefevre, P. Zimmermann. *Searching Worst Cases of a One-Variable Function Using Lattice Reduction*. IEEE Transactions on Computers, 54(3):340–346, March 2005.
- [73] <http://www.ioria.fr/~zimmerma/mpcheck/>
- [74] <http://www.maplesoft.com/>
- [75] <http://www.wolfram.com/products/mathematica/index.html>
- [76] <http://www.mathworks.com/products/matlab/>
- [77] A. Ziv. *Fast Evaluation of Elementary Mathematical Functions with Correctly Rounded Last Bit*. ACM Transactions on Mathematical Software, 17(3):410–423, September 1991.
- [78] *IBM Accurate Portable MathLib*
<http://rpmfind.net/linux/rpm2html/search.php?query=libultim.so.2>
- [79] <http://www.mpfr.org/>
- [80] <http://www.sun.com/download/products.xml?id=41797765>
- [81] <http://www.ens-lyon.fr/LIP/Arenaire/Ware/SCSLib/>
- [82] <http://lipforge.ens-lyon.fr/projects/crlibm/>
- [83] F. de Dinechin, A. Ershov, and N. Gast. *Towards the post-ultimate libm*. Proc. of 17-th Symposium on Computer Arithmetic. IEEE Computer Society Press, June 2005.
- [84] D. Defour, F. de Dinechin, J.-M. Muller. *Correctly Rounded Exponential Function in Double Precision Arithmetic*. INRIA Research report RR-2001-26, July 2001.
- [85] F. de Dinechin, C. Lauter, J.-M. Muller. *Fast and Correctly Rounded Logarithms in Double-Precision*. INRIA Research report RR-2005-37, September 2005.
- [86] S. Chevillard, N. Revol. *Computation of the Error Functions erf and erfc in Arbitrary Precision with Correct Rounding*. Proc. of 17-th IMACS Conf. on Scientific Computation, Applied Math. and Simulation, Paris, France, July 2005.
- [87] W. Kramer. *Multiple-Precision Computations with Result Verification*. In E. Adams, U. Kulisch, eds. Scientific Computing with Automatic Result Verification, pp. 325–356, Academic Press, 1993.
- [88] M. J. Schulte, E. E. Swartzlander. *Software and Hardware Techniques for Accurate, Self-Validating Arithmetic*. Applications of Interval Computations, pp. 381–404, 1996.
- [89] N. Revol, F. Rouillier. *Motivations for an Arbitrary Precision Interval Arithmetic and the MPFI Library*. Reliable Computing, 11(4):275–290, 2005.
- [90] *MPFI Library* http://perso.ens-lyon.fr/nathalie.revol/mpfi_toc.html
- [91] F. Lindemann. *Über die Zahl π* . Mathematische Annalen, vol. 20, pp. 213–225, 1882.

- [92] А. Я. Хинчин. *Ценные дроби*. М: Наука, 1978.
- [93] K. C. Ng. *Arguments Reduction for Huge Arguments: Good to the Last Bit*. 1992. Доступна как <http://www.validlab.com/arg.pdf>.
- [94] W. Kahan. *Minimizing $q^*m - n$* . 1983. Неопубликованные заметки, доступны по <http://http.cs.berkeley.edu/~wkahan/testpi/nearpi.c>.