

Elementary Functions from Kernels

W. Kahan Oct. 24, 1985

Given binary floating-point subprograms to calculate the "Kernels"
 $\ln(x)$ for $x \geq 0$ and $\lnp(x) := \ln(1+x)$ for $x \geq -1$,
 $\exp(x)$ and $\expm1(x) := \exp(x)-1$ for all x , and
 $\tan(x)$ for $|x| < \pi/8$ and $\arctan(x)$ for $|x| \leq \sqrt{2}-1$,
to nearly full working accuracy, we may calculate all the other
elementary transcendental functions almost as accurately, and with
no violation of (weak) monotonicity, as follows. Rounding must
conform to IEEE 754 or p854. We will need a threshold t
chosen about as large as possible subject to the constraint that
 $1 - t^2$ round to 1 to working precision; and we shall use
 $z := |x|$ and $s := \text{copysign}(1, x) = \pm 1$. We also abbreviate
 $\expm1$ to E and \lnp to L .

$\sinh(x) := x$ if $z < t$, else (provided $E(z)$ doesn't overflow)
 $:= s * (E(z) + E(z)/(1+E(z))) / 2$... certainly monotonic.

$\cosh(x) := 0.5 * \exp(z) + 0.25 / (0.5 * \exp(z))$... " " .

$\tanh(x) := x$ if $z < t$, else
 $:= -s * E(-2*z) / (2 + E(-2*z))$.

$\operatorname{asinh}(x) := x$ if $z < t$, else, unless $2z$ overflows,
 $:= s * L(z + z / (1/z + \sqrt{(1+(1/z)^2)}))$ ignoring underflow.

For slightly better accuracy when $z > 4/3$, use

$\operatorname{asinh}(x) := s * \ln(2z + 1 / (z + \sqrt{(1+z^2)}))$ if $z < 1/t$, else
 $:= s * (\ln(z) + \ln(2))$.

$\operatorname{acosh}(x) := +L(\sqrt{(x-1)} * (\sqrt{(x-1)} + \sqrt{(x+1)}))$ unless $2x$ overflows.

For slightly better accuracy,

$\operatorname{acosh}(x) := \ln(x) + \ln(2)$ if $x > 1/t$, else
 $:= \ln(2x - 1 / (x + \sqrt{(x^2-1)}))$ if $5/4 < x \leq 1/t$, else
 $:= L((x-1) + \sqrt{(2(x-1) + (x-1)^2)})$.

$\operatorname{atanh}(x) := x$ if $z < t$, else
 $:= s * L(2*z / (1-z)) / 2$.

$\arctan(x) := s * \pi/2 - \arctan(1/x)$ if $z > 1$, or (monotonically)
 $:= s * \pi/4 + \arctan((x-s) / (x+s))$ if $\sqrt{2}-1 < z < \sqrt{2}+1$.

$\arcsin(x) := x$ if $z < t$, else
 $:= \arctan(x / \sqrt{(1-z^2)})$ if $t \leq z \leq 1/2$, else
 $:= \arctan(x / \sqrt{(2(1-z) - (1-z)^2)})$ ignoring divide-by-zero.

$\arccos(x) := 2 * \arctan(\sqrt{((1-x)/(1+x))})$ ignoring divide-by-zero.

For $z \leq \pi/4$ let $T(x) := 2 \tan(x/2)$; then

$T(x) := \tan(x) := \sin(x) := x$ and $\cos(x) := 1$ if $z < t$.

Otherwise compute $\tan(x)$, $\sin(x)$ and $\cos(x)$ thus for $z \leq \pi/2$:

$\tan(x) :=$ if $z < \pi/8$ then $T(2*x)/2$
else if $3\pi/8 < z$ then $2s/T(\pi-2*z)$
else $s * (2 + T(2*z-\pi/2)) / (2 - T(2*z-\pi/2))$.

(Check monotonicity as z passes through $\pi/8$ and $3\pi/8$.)

If $\pi/4 \leq z \leq \pi/2$ then the formulas $\sin(x) = s \cdot \cos(\pi/2 - z)$ and $\cos(x) = \sin(\pi/2 - z)$ reduce the argument x to y satisfying $|y| \leq \pi/4$, wherein we compute $T := T(y)$, $q := T^2$, and then

$$\sin(y) := y - y/(1+4/q);$$

$$\cos(y) := \text{if } q < 4/15 \text{ then } 1 - 2/(1+4/q)$$

$$\text{else } 3/4 + ((1-2*q) + q/4)/(4+q).$$

Monotonicity is preserved except possibly as x passes through multiples of $\pi/4$, where the accuracy of $T(x)$ matters.

Some implementations of $\tan(x/2)$ actually deliver two functions $A(x)$ and $B(x)$ satisfying $A(x)/B(x) = \tan(x/2)$ for $|x| \leq \pi/4$, on which range $|A(x)/B(x)| < \sqrt{2} - 1 \approx 0.414\dots$. These can be used to deliver \sin , \cos and \tan more economically than above, and monotonically too provided $A(x)/B(x)$ is monotonic. For $t < z \leq \pi/4$ let $r := B(x)/A(x) > \sqrt{2} + 1$; and then

$$\sin(x) := 2/(r+1/r) \quad \text{and} \quad \cos(x) := 1 - 2/(1+r^2).$$

If both of $\sin(x)$ and $\cos(x)$ are wanted simultaneously, a more economical pair of formulas is

$$\sin(x) := 2/(r+1/r) \quad \text{and} \quad \cos(x) := 1 - (1/r) \sin(x).$$

To ensure monotonicity as x passes through multiples of $\pi/4$, check that computed $\sin(\pi/4) \leq$ computed $\cos(\pi/4)$; else use a better formula for \cos (see above). Computing $\tan(x)$ for $|x| \leq \pi/2$ from $A(x)$ and $B(x)$ is much like before:

$$\tan(x) := \text{if } z < \pi/8 \text{ then } A(2*x)/B(2*x)$$

$$\text{else if } 3\pi/8 < z \text{ then } B(s*\pi-2*x)/A(s*\pi-2*x)$$

$$\text{else } s*(B(y)+A(y))/(B(y)-A(y)) \quad \text{where } y := 2*z-\pi/2.$$

Monotonicity must be checked as z passes through $\pi/8$ and $3\pi/8$.

Other topics to be added later:

y^x
 $\text{atan2}(y,x) = \text{Arg}(x + iy)$, especially with ± 0 and $\pm \infty$
 $\text{cabs}(x + iy) = \sqrt{x^2 + y^2}$
 other complex elementary functions

approximating $\tan(z)$ for $0 < z < \pi/8$
 $\arctan(z)$ for $0 < z \leq \sqrt{2} - 1$
 $\lnlp(x)$ and $\ln(x)$ and $\text{expml}(x)$ and $\exp(x)$
 argument reduction

Given $A(x)$ and $B(x)$ above, which is better:

$$r := B(x)/A(x) \quad \text{and then compute } 1/r, \quad \text{or}$$

$$r := B(x)/A(x); \quad (1/r) := A(x)/B(x); \quad ?$$

What is wrong with

$$v := 2A/(A^2+B^2); \quad \sin(x) := vB; \quad \cos(x) := 1 - vA; \quad ?$$

Elementary Inequalities among Elementary Functions

W. Kahan Aug. 19, 1985

Programmers, like other people, frequently take familiar properties of elementary functions for granted. If $x \leq y$, for instance, they expect $\exp(x) \leq \exp(y)$; the possibility that computed $\exp(x) > \text{computed } \exp(y)$ might occur because of rounding errors is unlikely to be considered until after it has caused a disagreeable surprise. Such a violation of expected monotonicity is potentially more troublesome than an error of several ulps in the computed value of $\exp(x)$. Fortunately, library programs that compute $\exp(x)$ can easily be made monotonic even when, for very large $|x|$, they cannot easily be kept accurate within an ulp. For some other functions, like \cos and \log , the preservation of monotonicity can challenge the implementor. And if that challenge is overcome, inequalities among different but related elementary functions can pose problems of a still higher order of difficulty. How far is an implementor obliged to go to protect inequalities among elementary functions from roundoff?

To appreciate better the limits upon an implementor's powers, let us consider the following examples of elementary inequalities:

- L: $x/(1+x) \leq \lnp(x) := \ln(1+x) \leq x$ for all $x > -1$.
- E: $x \leq \text{expml}(x) := \exp(x) - 1$ for all x ; and
- EL: $\text{expml}(x) \leq -\lnp(-x) \leq x/(1-x)$ for all $x < 1$.

The inequalities $\lnp(x) \leq x$ and $x \leq \text{expml}(x)$ can be enforced by keeping the errors in the implementations of \lnp and expml below one ulp when $|x|$ is tiny; this is not hard to do. But no amount of care in the implementation of \lnp can enforce the inequality $x/(1+x) \leq \lnp(x)$ despite roundoff in $x/(1+x)$. For instance take $x = 0.00499$ and perform arithmetic rounded to 3 significant decimals. Then $1+x = 1.00499$ rounds to $[1+x] = 1$, and then $x/[1+x]$ rounds to x . But $\lnp(x) = 0.0049775912\dots$ rounds to $0.00498 < x$, violating the inequality in question. A similar example disposes of $\text{expml}(x) \leq x/(1-x)$. The inequality $\text{expml}(x) \leq -\lnp(-x)$ is more subtle; now try $x = 0.00000\ 99999$ in a context where arithmetic is performed to 5 sig. dec. Since

$$\begin{aligned} \text{expml}(x) &= 0.00000\ 99999\ 49999\ 1667\dots \\ &< 0.00000\ 99999\ 49999\ 3333\dots = -\lnp(-x), \end{aligned}$$

an implementor could not round these to $0.00000\ 99999$, that is to 5 sig. dec., without first knowing them to at least 10 sig. dec., twice as many. If each value were computed independently in error by as much as $\pm 0.00000\ 00000\ 00001$, rounding them subsequently to 5 sig. dec. could yield $0.00001\ 0000$ for $\text{expml}(x)$ and $0.00000\ 99999$ for $\lnp(x)$, violating the inequality in question.

It seems extravagant to carry more than twice as many figures as will be returned; and doing so would not by itself guarantee no argument x exists for which far more precision than that is needed to round well enough to preserve an inequality. Another unsatisfactory strategy for preserving inequalities is to use only algorithms designed for the purpose; the strategy is unattractive because the only such algorithms known at this time involve the use of Taylor series to the exclusion of economized polynomials or continued fractions or other more interesting schemes. Therefore the thoughtful programmer must acquiesce to the occasional violation of some familiar inequalities by roundoff.

What relations among elementary functions deserve to be taken for granted? One of them, monotonicity, is a subject too delicate to be discussed here; my report on the subject appears elsewhere. A second relation concerns "Cardinal Values"; these are exact values taken by transcendental functions. A collection of them is displayed in Table 1. A third relation concerns "Functional Identities"; the best-known examples are the odd functions like $\sin(-x) = -\sin(x)$, $\arctan(-x) = -\arctan(x)$, ... and the even ones like $\cos(-x) = \cos(x)$, Less well-known, perhaps because they are wrongly taken for granted, are identities like $\sqrt{x^2} = |x|$, which is satisfied, for all floating-point numbers x for which x^2 does not over/underflow, by correctly rounded square and square root operations in binary and quaternary floating-point arithmetic. The identity fails for some x when the arithmetic's radix exceeds 4. The complementary identity $(\sqrt{x})^2 = x$, on the other hand, cannot survive roundoff for all positive x , regardless of radix or rounding correctness. The most general discussion so far of Functional Identities was published in *Math. of Computation* in 1971 by Harry Diamond.

A fourth relation among elementary functions includes inequalities of the forms $f(x) \leq \text{Constant}$ and $f(x) \leq x$ or $f(x) \geq x$. Such inequalities can be preserved in implementations of $f(x)$ by keeping its error below one ulp, so they deserve to be taken for granted. Table 2 contains a collection of inequalities involving a representable Constant. Inequalities E and L above are instances of inequalities involving x , and some more follow:

The following string of inequalities involves only odd functions of x , and is therefore stated only for all sufficiently small positive values of x . Reversing the sign of x reverses the sense of all the inequalities in the string.

$$x \cos x < \tanh x < \arctan x < \sin x < \operatorname{arcsinh} x < x \dots$$

$$x < \sinh x < \arcsin x < \tan x < \operatorname{arctanh} x < x \cosh x .$$

Some of these inequalities remain valid as x increases from 0 only so long as x remains below some threshold. The thresholds are tabulated below:

At	$x = 0.74461\ 14991\ 45\dots$,	$\operatorname{arctanh} x = x \cosh x .$
At	$x = 0.97743\ 48912\ 2\dots$,	$\tan x = x \cosh x .$
At	$x = 0.99990\ 60124\ 1267\dots$,	$\arcsin x = \tan x .$
For	$x > 1$	remove	$\arcsin x$ and $\operatorname{arctanh} x$ from the string.
At	$x = 1.55708\ 58155\ \dots$,	$\arctan x = \sin x .$
For	$x \geq \pi/2 = 1.57079\ 63268\ \dots$	remove	$\tan x .$
At	$x = 1.87510\ 40687\ \dots$,	$\tanh x = \sin x .$
At	$x = 4.49340\ 94579\ \dots$,	$x \cos x = \sin x .$
At	$x = 4.91716\ 45703\ \dots$,	$x \cos x = \tanh x .$
At	$x = 4.99108\ 47512\ \dots$,	$x \cos x = \arctan x .$
At	$x = 5.18250\ 39692\ \dots$,	$x \cos x = \operatorname{arcsinh} x .$

Much as we might wish that the whole string of inequalities would persist as long as x remains between 0 and whatever threshold is pertinent, any of those inequalities demanding more than a comparison with x can succumb to roundoff when x is tiny.

Table 1: EXACT CARDINAL VALUES

~~~~~

**Positive zeros:**  $\ln(1) = \operatorname{arccosh}(1) = \arccos(1) = \exp(-\infty) =$   
 $= (+0) \text{ (even } > 0) = (+\infty) \text{ (even } < 0) = (+0) \text{ (noninteger } > 0) =$   
 $= (+\infty) \text{ (noninteger } < 0) = (\text{fraction})^{+\infty} = (+(>1))^{-\infty} = +0 .$

**Signed zeros:**  $\sin(+0) = \arcsin(+0) = \sinh(+0) = \operatorname{arcsinh}(+0) =$   
 $= \operatorname{lnip}(+0) = \tan(+0) = \arctan(+0) = \tanh(+0) = \operatorname{arctanh}(+0) =$   
 $= \operatorname{expml}(+0) = \sqrt{(+0)} = (+0) \text{ (odd } > 0) = (+\infty) \text{ (odd } < 0) = \pm 0 \text{ resp.}$

Whether  $\sin(n\pi)$ ,  $\tan(n\pi)$  or  $\cos((n+1/2)\pi)$  can vanish and, if so, what sign to assign to 0, depend upon how trigonometric argument reduction is performed.

**Ones:**  $\cos(0) = \cosh(0) = \tanh(+\infty) = \exp(0) = (\text{anything})^0 = 0! =$   
 $= 1! = 1^{\text{finite}} = (+1)^{\text{even}} = 1 ; \quad (-1)^{\text{odd}} = \tanh(-\infty) = -1 .$

Whether  $\cos(2n\pi) = \sin((2n+1/2)\pi) = \tan((n+1/4)\pi) = 1$  exactly depends upon how trigonometric argument reduction is performed.

**Integers:**  $\sqrt[n^2]{10^n} = \log_{10}(10^n) = n$  for all sufficiently small nonnegative integers  $n$ ;  $m**n = m^n$  is an integer too if  $|m|$  is an integer.

**Silent Infinities:**  $\sinh(+\infty) = \operatorname{arcsinh}(+\infty) = (+\infty) \text{ (odd } > 0) = +\infty \text{ resp.}$   
 $\cosh(+\infty) = \operatorname{arccosh}(+\infty) = \sqrt{(+\infty)} = \ln(+\infty) = \exp(+\infty) = (+(>1))^{+\infty} =$   
 $= (+\infty) \text{ (noninteger } > 0) = (+\infty) \text{ (even } > 0) = (\text{fraction})^{-\infty} = +\infty .$

**Signaled Infinities:**  $\operatorname{arctanh}(+1) = (+0) \text{ (odd } < 0) = +\infty \text{ resp.}$   
 $-\ln(0) = 0 \text{ (even } < 0) = 0 \text{ (noninteger } < 0) = +\infty .$

Whether  $\tan((n+1/2)\pi)$  is infinite and, if so, its sign depend upon how trigonometric argument reduction is performed. None the less, the identity  $\tan(-x) = -\tan(x)$  should still hold.

**Arg(x + iy) = ATAN2(y,x)** has values some of which are determined by consistency with complex arithmetic; to describe these special values we let  $\omega$  and  $\Omega$  stand for arbitrary real variables subject only to the constraints  $0 \leq \omega < \Omega \leq +\infty$ :

$\operatorname{ATAN2}(+0, +0) = \operatorname{ATAN2}(+0, +\Omega) = \operatorname{ATAN2}(+\omega, +\infty) = \pm 0 \text{ resp. ;}$   
 $\operatorname{ATAN2}(+0, -0) = \operatorname{ATAN2}(+0, -\Omega) = \operatorname{ATAN2}(+\omega, -\infty) = \pm\pi \text{ resp. ;}$   
 $\operatorname{ATAN2}(+\Omega, +\Omega) = \pm\pi/4 \text{ resp. ;} \quad \operatorname{ATAN2}(+\Omega, -\Omega) = \pm 3\pi/4 \text{ resp. ;}$   
 $\operatorname{ATAN2}(+\infty, +\infty) = \operatorname{ATAN2}(+\infty, -\infty) = \operatorname{ATAN2}(+\Omega, 0) = \pm\pi/2 \text{ resp.}$

Table 2: CONSTANT BOUNDS

~~~~~

$|\sin| \leq 1 ; |\cos| \leq 1 ; |\tanh| \leq 1 \leq \cosh ; 0 \leq \exp ; 0 \leq \sqrt{ ;}$
 $0 \leq \operatorname{arccosh} ; 0 \leq \arccos \leq \pi ; |\arcsin| \leq \pi/2 ; |\arctan| \leq \pi/2 .$