

# Floating point accuracy

Vincent LAFAGE 

IJCLab, Laboratoire de Physique des 2 Infinis Irène Joliot-Curie  
CNRS/IN2P3 & Université Paris-Saclay, Orsay, France



April 15th 2026



Revisiting "What Every Computer Scientist Should Know About Floating-point Arithmetic"



- Numbers: real, algebraic, constructibles, decimal, binary, floating point...
- «Primitive» types: float, double, long double, quad, half...
- When computations don't turn out as expected...(why, how)
  - ▶ rounding errors
  - ▶ conversion errors
  - ▶ propagating errors
  - ▶ composing errors
- Heuristics for accuracy:
  - how a rough estimate can save epsilons
- Nondimensionalisation and formula entropy reduction
- How to reconcile nondimensionalisation and performance
- How to reconcile abstraction and accuracy: functions of a complex variable
- Why are geometrical computations so hard
- The hidden side of functional programming: towards total functions



## Numbers that cost and kill

- Vancouver Stock Exchange, from January 1982 to November 1983 : truncating instead of rounding, 4 digits instead of 3, error cumulated over two years on the value of the all-share stock market index  
52 % error : 524.811 \$ instead of 1098.892 \$
- Patriot Missiles, first Gulf War, February 25, 1991:  
600 m error (due to rounding error on time) for interception : 28 killed, a hundred injured
- Sleipner A offshore platform hull collapse, August 23, 1991:  
inaccurate finite element analysis (NASTRAN). Stresses on the ballast chambers were underestimated by 47 %. Sunk: 700 millions USD.
- Pentium FDIV bug, Intel, 1994:  
recall of million of chips that computed division wrong (up to 1/1000 relative error): about 475 millions USD
- Knight Capital Group Hedge Fund, August 1, 2012:  
They lost 440 millions USD in 45 minutes. Because of a floating point rounding error in their trading algorithm (plus bad deployment).



<https://doi.org/10.1145/103162.103163>

<https://www.validlab.com/goldberg/paper.pdf> (avec annexe)

*“Floating-point arithmetic is considered an esoteric subject by many people”*

## What Every Computer Scientist Should Know About Floating-Point Arithmetic

DAVID GOLDBERG

*Xerox Palo Alto Research Center, 3333 Coyote Hill Road, Palo Alto, California 94304*

Floating-point arithmetic is considered an esoteric subject by many people. This is rather surprising, because floating-point is ubiquitous in computer systems: Almost every language has a floating-point datatype; computers from PCs to supercomputers have floating-point accelerators; most compilers will be called upon to compile floating-point algorithms from time to time; and virtually every operating system must respond to floating-point exceptions such as overflow. This paper presents a tutorial on the aspects of floating-point that have a direct impact on designers of computer systems. It begins with background on floating-point representation and rounding



# Formats IEEE754

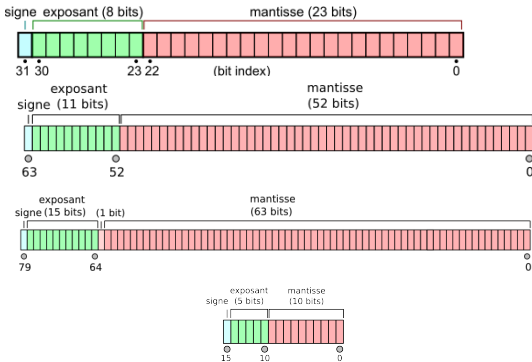
« computing is about representation »

Scientific notation:

significand  $\times$  base<sup>exponent</sup>      significand  $\in \mathbb{Z}$ , exponent  $\in \mathbb{Z}$

Standard form: mantissa, alias *normalized significand*

mantissa  $\times$  base<sup>exponent</sup>      Trick, for base 2: the most significant digit is always 1...



In the FPU registers, we widen mantissa with three bits: guard bit, round bit, "sticky" bit

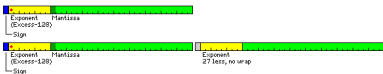


# Older Formats

<http://www.quadibloc.com/comp/cp0201.htm>

## Group 1

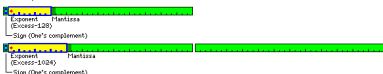
International Business Machines 704, 708, 7040, 7044, 7080, 7094, 7094 II



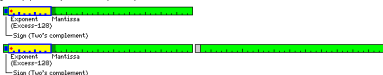
Hewlett-Packard 3000



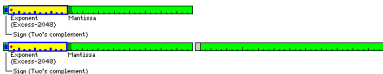
Univac 1107, 1108



Digital Equipment Corporation PDP-6, PDP-10, BECSYSTEM-20



Expanded range (RL-10 only)



KA-10 Double Precision

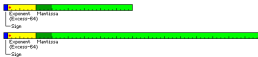


## Group 1 (continued)

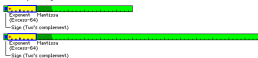
Digital Equipment Corporation PDP-8 Special (BK PDP7500)



International Business Machines System/360, System/370, ESA/390, z/Architecture



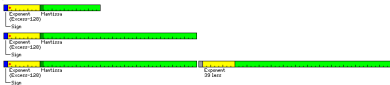
Series Data Systems Signo



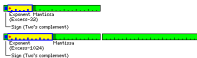
Control Data Corporation 3604, 3605



Digital Electric KDF9



Powerto PDP-1



Pulsar-Bell PB440

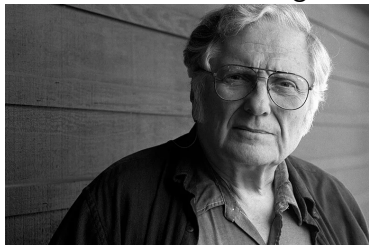


Univac 458





« The Father of Floating Point »



TURING Award 1989 for  
"his fundamental contributions to numerical analysis."

Consultant for Intel FPU  $\times 87$

<https://people.eecs.berkeley.edu/~wkahan/index.htm>



# Example float32

$$\text{float} = (-1)^S \times 2^{E-127} \times (1 + M), \quad M \in [0, 1[$$

31302928272625242322212019181716151413121110 9 8 7 6 5 4 3 2 1 0

S	E	M
0	0111 1111	000 0000 0000 0000 0000 0000
0	1000 0000	000 0000 0000 0000 0000 0000
0	1000 0000	100 0000 0000 0000 0000 0000
0	1000 0000	110 0000 0000 0000 0000 0000
0	1000 0000	111 0000 0000 0000 0000 0000
1	0111 1111	000 0000 0000 0000 0000 0000
0	0111 1110	000 0000 0000 0000 0000 0000
0	0111 1100	100 1100 1100 1100 1100 1101
0	0111 1101	010 1010 1010 1010 1010 1011
0	0111 1111	011 0101 0000 0100 1111 0011
0	1000 0000	100 1001 0000 1111 1101 1011
0	0000 0000	000 0000 0000 0000 0000 0000
1	0000 0000	000 0000 0000 0000 0000 0000
0	1111 1110	111 1111 1111 1111 1111 1111
0	1111 1111	000 0000 0000 0000 0000 0000
0	1111 1111	1xx xxxx xxxx xxxx xxxx xxxx
0	1111 1111	1xx xxxx xxxx xxxx xxxx xxxx
0	1111 1111	01x xxxx xxxx xxxx xxxx xxxx
0	0000 0001	000 0000 0000 0000 0000 0000
0	0000 0000	000 0000 0000 0000 0000 0001
0	0000 0000	111 1111 1111 1111 1111 1111

- }  $\text{float} = (-1)^S \times 2^{E-127} \times (1 + M)$
- }  $1 = 2^0 \times (1 + 0)$
- }  $2 = 2^1 \times (1 + 0)$
- }  $3 = 2^1 \times (1 + 1/2)$
- }  $3.5 = 2^1 \times (1 + 1/2 + 1/4)$
- }  $3.75 = 2^1 \times (1 + 1/2 + 1/4 + 1/8)$
- }  $-1 = -2^0 \times (1 + 0)$
- }  $1/2 = 2^{-1} \times (1 + 0)$
- }  $0.2 = 2^{-3} \times (1 + 1/2) \times \sum_n 1/16^n$
- }  $1/3 = 2^{-2} \times (1 + 1/4) \times \sum_n 1/16^n$
- }  $\sqrt{2}$
- }  $\pi \approx 2^1 \times (1 + 1/2 + 1/16 + 1/128 + \dots)$
- } 0 special representation
- } 0<sub>-</sub> special representation
- } largest float  $3.402823466 \times 10^{38}$
- }  $+\infty = \text{Inf}$  special representation
- } NaN special representation
- } qNaN *quiet* special representation
- } sNaN *signaling* special representation
- } smallest positive float  $1.17549435 \times 10^{-38}$
- } smallest **denormal** positive float  $1.401 \times 10^{-45}$
- } largest **denormal** positive float  $1.17549379 \times 10^{-38}$

⇒ using the link below, represent your favorite numbers:

<https://www.h-schmidt.net/FloatConverter/IEEE754.html>

<https://evanw.github.io/float-toy/>

<https://bartaz.github.io/ieee754-visualization/>

<https://float.exposed/0x40490fdb>





```
#include <stdio.h>

int main ()
{
    float x = 1.0f;
    printf ("%f = %a\n", x, x);
    x = 2.0f;
    printf ("%f = %a\n", x, x);
    x = 3.0f;
    printf ("%f = %a\n", x, x);
    x = 3.141592653589793f;
    printf ("%f = %a\n", x, x);
}
```

1.000000 = 0x1p+0

2.000000 = 0x1p+1

3.000000 = 0x1.8p+1

3.141593 = 0x1.921fb6p+1



# Get Hexadecimal displayed C++

```
#include <iostream>

int main ()
{
    float x = 1.0f;
    std::cout << x << " = " << std::hexfloat << x << std::defaultfloat << '\n';
    x = 2.0f;
    std::cout << x << " = " << std::hexfloat << x << std::defaultfloat << '\n';
    x = 3.0f;
    std::cout << x << " = " << std::hexfloat << x << std::defaultfloat << '\n';
    x = 3.141592653589793f;
    std::cout << x << " = " << std::hexfloat << x << std::defaultfloat << '\n';
}
```

1 = 0x1p+0

2 = 0x1p+1

3 = 0x1.8p+1

3.14159 = 0x1.921fb6p+1



```
program hexfloat
  use, intrinsic :: iso_fortran_env, only: real32
  implicit none
  real (real32) :: x

  x = 1
  write (*, '(F10.6,A,Z16)') x, ' = ', x
  x = 2
  write (*, '(F10.6,A,Z16)') x, ' = ', x
  x = 3
  write (*, '(F10.6,A,Z16)') x, ' = ', x
  x = acos (-1.0_real32)
  write (*, '(F10.6,A,Z16)') x, ' = ', x
end program hexfloat
```

```
1.000000 =      3F800000
2.000000 =      40000000
3.000000 =      40400000
3.141593 =      40490FDB
```



# Get Hexadecimal displayed Python

```
#!/usr/bin/python3
```

```
x = 1.0
print (x, " = ", float.hex(x))
x = 2.0
print (x, " = ", float.hex(x))
x = 3.0
print (x, " = ", float.hex(x))
x = 3.141592653589793
print (x, " = ", float.hex(x))
```

```
1.0 = 0x1.0000000000000p+0
2.0 = 0x1.0000000000000p+1
3.0 = 0x1.8000000000000p+1
3.141592653589793 = 0x1.921fb54442d18p+1
```



# Get Hexadecimal displayed

Ada

```
with Ada.Text_Io;

procedure Hexfloat is
  use Ada.Text_Io;
  X : Float := 1.0;
begin
  Put_Line (X'Image & " = 2^" & Float'Exponent (X)'Image & " x " & Float'Fractional (X)'Image);
  X := 2.0;
  Put_Line (X'Image & " = 2^" & Float'Exponent (X)'Image & " x " & Float'Fractional (X)'Image);
  X := 3.0;
  Put_Line (X'Image & " = 2^" & Float'Exponent (X)'Image & " x " & Float'Fractional (X)'Image);
  X := 3.141592653589793;
  Put_Line (X'Image & " = 2^" & Float'Exponent (X)'Image & " x " & Float'Fractional (X)'Image);
end Hexfloat;
```

1.00000E+00 = 2 <sup>1</sup> x 5.00000E-01

2.00000E+00 = 2 <sup>2</sup> x 5.00000E-01

3.00000E+00 = 2 <sup>2</sup> x 7.50000E-01

3.14159E+00 = 2 <sup>2</sup> x 7.85398E-01



# Get HexaDecimal displayed Rust

```
fn extract_components(x: f32) -> (char, i32, u32) {
    let bits = x.to_bits();
    let sign = if (bits >> 31) & 1 == 0 { '+' } else { '-' };
    let mantissa = (bits & ((1 << 23) - 1)) * 2;
    let exponent = (((bits >> 23) & 0xFF) as i32) - 127;
    (sign, exponent, mantissa)
}

pub fn main() {
    let x: f32 = 1.0;
    let (sign, exponent, mantissa) = extract_components(x);
    println!("{}", = {}0x1.{:x}p{}", x, sign, mantissa, exponent);
    let x: f32 = 2.0;
    let (sign, exponent, mantissa) = extract_components(x);
    println!("{}", = {}0x1.{:x}p{}", x, sign, mantissa, exponent);
    let x: f32 = 3.0;
    let (sign, exponent, mantissa) = extract_components(x);
    println!("{}", = {}0x1.{:x}p{}", x, sign, mantissa, exponent);
    let x: f32 = 3.141592653589793;
    let (sign, exponent, mantissa) = extract_components(x);
    println!("{}", = {}0x1.{:x}p{}", x, sign, mantissa, exponent);
    let x: f32 = -0.3141592653589793;
    let (sign, exponent, mantissa) = extract_components(x);
    println!("{}", = {}0x1.{:x}p{}", x, sign, mantissa, exponent);
}
```

1 = +0x1.0p0

2 = +0x1.0p1

3 = +0x1.800000p1

3.1415927 = +0x1.921fb6p1

-3.1415927 = -0x1.921fb6p1



$$\pm d_0.d_1d_2 \cdots d_{p-1} \times \beta^e$$

<https://www.netlib.org/blas/machar.f>

[https://www.netlib.org/misc/machar\(C\)](https://www.netlib.org/misc/machar(C))

[https://people.sc.fsu.edu/~jburkardt/f\\_src/machar/machar.html](https://people.sc.fsu.edu/~jburkardt/f_src/machar/machar.html) (f90, C++, Py,...)

	Machar	Half fp16	Single fp32	Double fp64	Extended fp80	Quad fp128	Octuple fp256	Cephes† fp384
$\beta$	Ibeta	2	2	2	2	2	2	2
$p$	It	11	24	53	64	113	237	352
	Irnd	5	5	5	5	5	5	1
	Ngrd	0	0	0	0	0	0	0
$\log_2 \epsilon$	Machep	-10	-23	-52	-63	-112	-236	-352
	Negep	-11	-24	-53	-64	-113	-237	-353
	Iexp	5	8	11	15	15	19	16
$e_{min}$	Minexp	-14	-126	-1022	-16382	-16382	-262142	-32768
$e_{max}$	Maxexp	16	128	1024	16384	16384	262144	32767
$\epsilon$	Eps	9.766...e-4	1.192...e-7	2.220...e-16	1.084...e-19	1.926...e-34	9.056...e-72	1.090...e-106
	Epsneg	4.883...e-4	5.960...e-8	1.110...e-16	5.421...e-20	9.630...e-35	4.528...e-72	5.450...e-107
	Xmin	6.104...e-5	1.175...e-38	2.225...e-308	3.362...e-4932	3.362...e-4932	2.482...e-78913	7.065...e-9865
	Xmax	6.550...e+4	3.403...e+38	1.798...e+308	1.190...e+4932	1.190...e+4932	1.611...e+78913	7.077...e+9865

† <https://www.netlib.org/cephes/>



# Interface *leaky abstraction*

Machar	C / C++ (/ Python)	Fortran'90	ieee_arithmetic	Ada
Ibeta	numeric_limits::radix	radix (x)		F'Machine_Radix
It	numeric_limits::digits	digits (x)		F'Machine_Mantissa
		range (x)		F'Digits
Eps	numeric_limits::epsilon ()	epsilon (x)		F'Model_Epsilon
		precision (x)		
Minexp	numeric_limits::min_exponent	minexponent (x)		F'Machine_Emin
Maxexp	numeric_limits::max_exponent	maxexponent (x)		F'Machine_Emax
Xmin	numeric_limits::min_exponent	tiny (x)		F'Model_Small
Xmax	numeric_limits::max_exponent	huge (x)		F'Safe_Last
	copysign (d x, d y)	sign (x, y)	ieee_copy_sign (x, y)	F'Copy_Sign (value, sign)
	frexp (d x, i *exp)	exponent (x)	ieee_logb (x)	F'Exponent (x)
		fraction (x)		F'Fraction (x)
	ldexp (d x, i exp)			
	scalbn (d x, i exp)	set_exponent (x, i)	ieee_scalb (x, i)	F'Scaling (x, adjustment)
	nextafter(d x, d y)	nearest (x, s)	ieee_next_after (x, y)	F'Adjacent (x, towards)
		spacing (x)		
		rrspacing (x)		
	nearbyint (d x)			
	rint(d x)	nint (x)	ieee_rint (x)	F'Rounding (x)
	floor (d x)	floor (x)		F'Floor (x)
	ceil (d x)	ceiling (x)		F'Ceiling (x)
			ieee_rem (x, y)	F'Remainder (x, y)

Unfortunately the C/C++ API doesn't vectorise well.

You might need to extract exponent and mantissa in non standard way for performance



$$1 + 2 = 3, \quad 0.1 + 0.2 \neq 0.3?$$

$$\Sigma = a + b \stackrel{?}{=} c \quad \Delta = a + b - c$$

with

$$a = 0.1 \quad b = 0.2 \quad c = 0.3$$



$$1 + 2 = 3, \quad 0.1 + 0.2 \neq 0.3?$$

$$\Sigma = a + b \stackrel{?}{=} c \quad \Delta = a + b - c$$

with

$$a = 0.1 \quad b = 0.2 \quad c = 0.3$$

	$a$	$b$	$c$	$\Sigma$	$\Delta$
fp32	0.100000001	0.200000003	0.300000012	0.300000012	0
fp64	0.100000000000000001	0.200000000000000001	0.29999999999999999	0.300000000000000004	$5.551 \dots 10^{-17}$
fp80	0.1000000000000000000001	0.2000000000000000000003	0.3000000000000000000011	0.3000000000000000000011	0
fp16	0.099976	0.19995	0.30005	0.29980	$2.4414 \dots 10^{-4}$



$1 + 2 = 3, \quad 0.1 + 0.2 \neq 0.3?$

$$\Sigma = a + b \stackrel{?}{=} c \quad \Delta = a + b - c$$

with

$$a = 0.1 \quad b = 0.2 \quad c = 0.3$$

	$a$	$b$	$c$	$\Sigma$	$\Delta$
fp32	0.100000001	0.200000003	0.300000012	0.300000012	0
fp64	0.100000000000000001	0.200000000000000001	0.29999999999999999	0.30000000000000004	$5.551 \dots 10^{-17}$
fp80	0.1000000000000000000001	0.200000000000000000003	0.300000000000000000011	0.300000000000000000011	0
fp16	0.099976	0.19995	0.30005	0.29980	$2.4414 \dots 10^{-4}$

cf. <https://0.3000000000000000004.com/> in many languages



$$1 + 2 = 3, \quad 0.1 + 0.2 \neq 0.3?$$

$$\Sigma = a + b \stackrel{?}{=} c \quad \Delta = a + b - c$$

with

$$a = 0.1 \quad b = 0.2 \quad c = 0.3$$

	a	b	c	$\Sigma$	$\Delta$
fp32	0.100000001	0.200000003	0.300000012	0.300000012	0
fp64	0.100000000000000001	0.200000000000000001	0.29999999999999999	0.300000000000000004	$5.551 \dots 10^{-17}$
fp80	0.1000000000000000000001	0.200000000000000000003	0.300000000000000000011	0.3000000000000000000011	0
fp16	0.099976	0.19995	0.30005	0.29980	$2.4414 \dots 10^{-4}$

cf. <https://0.3000000000000000004.com/> in many languages

$\Rightarrow \mathbb{D} \not\subseteq \mathbb{B}$ : some decimal are not binary

$\Rightarrow$  binary conversion needs some rounding

$$\frac{1}{5} = 0.2_{10} = 0.00\overline{1100}_2 \dots \ominus 13421773 \times 2^{-26} = 0.2 + 2,98 \times 10^{-9}$$

*God created the integers, all else is the work of man.*

KRONECKER



# Decimal vs. binary ...and binary vs. floating

$$\mathbb{D} = \left\{ \frac{n}{10^p}, n \in \mathbb{Z}, p \in \mathbb{N} \right\} = \mathbb{Z}[1/10] \text{ (decimal)}$$

$$\mathbb{B} = \left\{ \frac{n}{2^p}, n \in \mathbb{Z}, p \in \mathbb{N} \right\} = \mathbb{Z}[1/2] \text{ (binary)}$$

$\mathbb{B} \subset \mathbb{D}$  but  $\mathbb{D} \not\subset \mathbb{B}$ :  $\frac{1}{5} \in \mathbb{D}$ ,  $\frac{1}{5} \notin \mathbb{B} \Rightarrow 0.1 + 0.2 \neq 0.3$  ( $\frac{1}{5} = 0.00\overline{1100}_2 \dots$ )  $\Rightarrow$  not good for financial computations...

● closure:

$$\forall (x, y) \in \mathbb{B}^2, \quad x + y \in \mathbb{B},$$

$$\forall (x, y) \in \mathbb{B}^2, \quad x \times y \in \mathbb{B}$$

● commutativity  $\forall (x, y) \in \mathbb{B}^2, \quad x + y = y + x,$

$$\forall (x, y) \in \mathbb{B}^2, \quad x \times y = y \times x$$

● associativity:

$$\forall (x, y, z) \in \mathbb{B}^3, \quad x + (y + z) = (x + y) + z,$$

$$\forall (x, y, z) \in \mathbb{B}^3, \quad x \times (y \times z) = (x \times y) \times z$$

● distributivity:

$$\forall (x, y, z) \in \mathbb{B}^3, \quad x \times (y + z) = x \times y + x \times z$$

● total order:

$$\forall (x, y, z) \in \mathbb{B}^3, \quad x \leq y \text{ and } y \leq z \Rightarrow x \leq z \quad (\text{transitivity});$$

$$\forall (x, y) \in \mathbb{B}^2, \quad x \leq y \text{ and } y \leq x \Rightarrow x = y \quad (\text{antisymmetry});$$

$$\forall x \in \mathbb{B}, \quad x \leq x \quad (\text{reflexivity});$$

$$\forall (x, y) \in \mathbb{B}^2, \quad x \leq y \text{ or } y \leq x \quad (\text{totality}).$$

● topology:

$\mathbb{B} \subset \mathbb{D} \subset \mathbb{Q}$  are dense in  $\mathbb{R} \Rightarrow$  arbitrarily close approximations to the real numbers



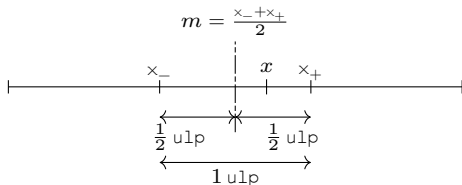
# Decimal vs. binary ...and binary vs. floating

- closure:  
 $\forall (x, y) \in \mathbb{F}^2, \quad x + y \notin \mathbb{F},$   
 $\forall (x, y) \in \mathbb{F}^2, \quad x \times y \notin \mathbb{F}$   
 $\Rightarrow$  rounding and extension  $\overline{\mathbb{F}} = \mathbb{F} \cup \{\pm \text{Inf}\} \cup \{\text{NaN}\} \cup \{0\_ \}$  overflow, underflow, inexact
- commutativity  $\forall (x, y) \in \mathbb{F}^2, \quad x + y = y + x,$   
 $\forall (x, y) \in \mathbb{F}^2, \quad x \times y = y \times x$
- associativity:  
 $\forall (x, y, z) \in \mathbb{F}^3, \quad x + (y + z) \neq (x + y) + z,$   
 $\forall (x, y, z) \in \mathbb{F}^3, \quad x \times (y \times z) \neq (x \times y) \times z$
- distributivity:  
 $\forall (x, y, z) \in \mathbb{F}^3, \quad x \times (y + z) \neq x \times y + x \times z$
- total order:  
 $\forall (x, y, z) \in \mathbb{F}^3, \quad x \leq y \wedge y \leq z \Rightarrow x \leq z \quad (\text{transitivity}) ;$   
 $\forall (x, y) \in \mathbb{F}^2, \quad x \leq y \wedge y \leq x \Rightarrow x = y \quad (\text{antisymmetry}) ;$   
 $\forall x \in \mathbb{F}, \quad x \leq x \quad (\text{reflexivity}) ;$   
 $\exists (x, y) \in \overline{\mathbb{F}}^2, \quad x \leq y \wedge y \leq x \quad (\text{NaN}).$
- topology:  
 $\mathbb{B} \subset \mathbb{D} \subset \mathbb{Q}$  are dense in  $\mathbb{R} \Rightarrow$  arbitrarily close approximations to the real numbers  
but  
 $\mathbb{F}$ : floating point numbers, finite parts of  $\mathbb{B}$  (or  $\mathbb{D}$ ) are dense nowhere



# Rounding

$\forall x \in \mathbb{R}, \exists (x_-, x_+) \in \mathbb{F}^2 \mid x_- \leq x \leq x_+$  (closest representable neighbours: *Unit in the Last Place*)



$\Rightarrow$  correct rounding requires at least 2 extra bits beyond target accuracy (*cf guard bit, round bit, "sticky" bit*)

or even more (*table maker's dilemma*)

correct rounding, faithful rounding, happy-go-lucky rounding

rounding is non-linear but completely deterministic!



## Mainstream definition

(in language constants in Ada, C, C++, Fortran, Python and Rust etc.,)

- **Definition:** The machine epsilon  $\epsilon_{\text{mach}}$  is the distance between 1.0 and the next representable floating-point number.

*a.k.a.* 1 ulp after 1

Equivalently: “the largest relative interval between two nearest numbers in finite-precision”

In other words,

$$\epsilon_{\text{mach}} = \min\{\delta > 0 : 1.0 + \delta \neq 1.0\}.$$

- **Typical values (IEEE 754):**

fp	$\epsilon$	approx.
fp16	$2^{-10}$	9.766...e-4
fp32	$2^{-23}$	1.192...e-7
fp64	$2^{-52}$	2.220...e-16
fp80	$2^{-63}$	1.084...e-19
fp128	$2^{-112}$	1.926...e-34

- **Significance:**  $\epsilon_{\text{mach}}$  gives us a rough measure of how close two floating-point numbers can be before they “merge” into one representable value.  $\Rightarrow$  Resolution



# Conversion

- $\mathbb{D} \not\subset \mathbb{B}$ : every decimal is not a binary

⇒ conversion to binary relies on rounding

$$\frac{1}{5} = 0.2_{10} = 0.001100_2 \dots \ominus 13421773 \times 2^{-26} = 0.2 + 2,98 \times 10^{-9}$$

---

4 byte	float	25.4E0 = 25.399999619...
8 byte	double	25.4D0 = 25.39999999999999858...
10 byte	long-double	25.4T0 = 25.3999999999999999653...
16 byte	quadruple	25.4Q0 = 25.399999999999999999999999999877...
2 byte	half	25.4_2 = 25.406...

---

- $\mathbb{B} \subset \mathbb{D}$ : every binary is a decimal

However, converting a binary, usually from a computation, usually for display or storage, is not toward the exactly corresponding decimal: it would require too many meaningless decimal digits.

$$\frac{1}{8} = 0.001_2 = 0.125_{10} \ominus 0.1_{10} \dots$$

⇒ conversion to decimal also relies on rounding



# Decimal conversion

- ⇒ use decimal floating points: `_Decimal32`, `_Decimal64`, `_Decimal128`  
(starting from C23)
- ⇒ program in SQL or COBOL...
- ⇒ change scale:  
count integer hundredth if you need 2 exact places
- ⇒ fixed point instead of floating point



# pi $\neq$ $\pi$ ?

	exact	fp32	fp64	fp16*
$\sin \pi$	0	-8.7422777e-8	1.2246467991473532e-16	9.6750e-4
$\cos \pi$	-1	-1.000000	-1.0000000000000000	-1.000
$\sin \frac{\pi}{6}$	$\frac{1}{2}$	0.5000000	0.4999999999999999	0.4998
$\cos \frac{\pi}{3}$	$\frac{1}{2}$	0.5000000	0.5000000000000001	0.5005
$\sin \frac{\pi}{3}$	$\frac{\sqrt{3}}{2}$	0.8660254	0.8660254037844386	0.8657
$\cos \frac{\pi}{6}$	$\frac{\sqrt{3}}{2}$	0.8660254	0.8660254037844387	0.8662

$$\begin{array}{lcl} \sin 0 = 0 & \Leftrightarrow & \sin(0.0) = 0 \\ \sin \pi = 0 & \text{but} & \sin(\pi) \neq 0 \quad \text{no finite representation...} \\ \sin \pi = 0 & = & \sin(\pi - \eta) = \sin \eta \sim \eta \end{array}$$

$$|\eta| < \pi \varepsilon / 2, \Rightarrow |\sin \pi| < \frac{\pi}{2} \varepsilon$$



$\text{pi} \neq \pi?$

If it is a problem

- ⇒ use half-turn trig functions: `sinpi`, `cospi`,... (starting from C23 / F23...)
- ⇒ use degrees trig functions: `sind`, `cosd`,... (all good Fortran compilers... + F23)



# Addition: absorption

$$H_N = \sum_{n=1}^N 1/n \sim \ln N + \gamma + \frac{1}{2N} - \frac{1}{12N^2} + \frac{1}{120N^4} - \frac{1}{252N^6} + \dots$$

Table – *Harmonic sum*

fp	N	up sum	down sum	theoretical sum
fp16	250	6.063	6.098	6.098
fp16	500	7.039	6.793	6.793
fp16	1000	7.086	7.477	7.484
fp16	2000	7.086	8.188	8.180
fp16	4000	7.086	8.789	8.875
fp16	8000	7.086	9.797	9.563
fp16	16000	7.086	9.797	10.26
fp16	32000	7.086	9.797	10.95
fp32	32000	10.95073	10.95072	10.95071
fp32	3200000	15.55911	15.55588	15.55588



# Addition: absorption

```
impure elemental subroutine harmonique (nbmax)
  integer (8), intent (in) :: nbmax
  integer (8) :: idx
  real (pr) :: somme_croissante = 0, somme_decroissante = 0, somme_theorique

  somme_croissante = 0
  somme_decroissante = 0
  somme_theorique = euler + log (real (nbmax, sp)) &
    + 1.0_pr / real (2 * nbmax, sp) &
    - 1.0_pr / (12 * real (nbmax, sp) **2) &
    + 1.0_pr / (120 * real (nbmax, sp) **4)

  do idx = 1, nbmax
    somme_croissante = somme_croissante + 1.0_pr / real (idx, pr)
  end do

  do idx = nbmax, 1, -1
    somme_decroissante = somme_decroissante + 1.0_pr / real (idx, pr)
  end do

  write (*, *) nbmax, somme_croissante, somme_decroissante, somme_theorique, &
    int ((somme_croissante/somme_theorique -1.0_pr)/epsilon (1.0_pr)), &
    int ((somme_decroissante/somme_theorique -1.0_pr)/epsilon (1.0_pr))
end subroutine harmonique
```



## Addition: KAHAN-BABUŠKA summation

```
!...  
accumulator = 0  
compensation = 0  
  
do jdx = 1, size (harmonic_array)  
  adjustedValue = harmonic_array (jdx) - compensation  
  tempSum      = accumulator + adjustedValue  
  compensation = (tempSum - accumulator) - adjustedValue  
  accumulator  = tempSum  
end do  
!...
```



# Addition: absorption

Table – *Harmonic sum*

fp	N	up sum	down sum	theoretical sum	Kahan up sum	Kahan down sum	fp64
fp16	62	4.707	4.711	4.715	4.711	4.711	4.712
fp16	125	5.398	5.414	5.410	5.410	5.410	5.410
fp16	250	6.062	6.098	6.102	6.102	6.102	6.101
fp16	500	7.039	6.793	6.793	6.793	6.793	6.793
fp16	1000	7.086	7.477	7.484	7.484	7.484	7.485
fp16	2000	7.086	8.188	8.180	8.180	8.180	8.178
fp16	4000	7.086	8.789	8.875	8.875	8.867	8.871
fp16	8000	7.086	9.797	9.562	9.562	9.562	9.564
fp16	16000	7.086	9.797	10.258	10.258	10.258	10.258
fp16	32000	7.086	9.797	10.953	10.953	10.953	10.951
fp16	64000	7.086	9.797	11.648	11.641	11.641	11.644
fp16	128000	7.086	9.797	12.336	11.664	11.664	12.337
fp16	65536000	7.086	9.797	18.578	11.664	11.664	18.575



# Hierarchy of operations

- **arithmetic:**  $+$ ,  $-$ ,  $\times$ ,  $/$ , integer powers
- **algebraic:**  $\sqrt{\quad}$ ,  $\sqrt[n]{\quad}$ , fractional powers and roots of polynomials
- **elementary (transcendental) functions:**  
exp, ln, sin, cos, irrational powers, all circular and hyperbolic trigonometry
- **higher transcendental functions *a.k.a.* special functions:**  
BESSEL, AIRY, Polylogarithm, elliptic integral, EULER  $\Gamma$  function, RIEMANN  $\zeta$  function,...

Correct rounding is guaranteed by the standard for:

- **arithmetic**
- **square root**



## transcendental functions

- ... costly
- ... **correct rounding** not guaranteed

Rounding is **non-linear**

- ⇒ mixing various scales
- ⇒ start runoff (butterfly effect)

to get correct rounding with  $n$  digits/bits... <https://members.loria.fr/PZimmermann/wc/decimal32.html>

$$\exp(0.5091077534282133) = \underbrace{1.663806007261509}_{16 \text{ digits}} \underbrace{5000000000000000}_{16 \text{ digits}} 49 \dots$$

$$\exp(0.7906867968553504) = \underbrace{2.204910231771509}_{16 \text{ digits}} \underbrace{4999999999999999}_{16 \text{ digits}} \dots$$

**Double rounding** (rounding from high precision to intermediate precision, then to low precision) can also give worse final rounding than expected.



glibc-2.41 ( $\Rightarrow$  libm-2.41) released on 2025-01-30!!!

<https://sourceware.org/glibc/wiki/Release/2.41>

$\Rightarrow$  correct rounding!

thanks Paul ZIMMERMANN



## Anyway, what physics needs double?

```
program stefan_constant ! ratio of total power radiated per unit surface area
                        ! to fourth power of absolute temperature of black body
use, intrinsic :: iso_fortran_env, only: pr => real32, real64 !, real128
real (pr), parameter :: pi = acos (-1.0_pr) !  $\pi = \text{Arccos}(-1)$ 
real (pr) ::
    lightspeed = 299792458.0_pr, & !  $c = 299\,792\,458 \text{ m}\cdot\text{s}^{-1}$ 
    planck      = 6.62607015e-34_pr, & !  $h = 6.626\,070\,15 \times 10^{-34} \text{ J}\cdot\text{Hz}^{-1}$ 
    boltzmann   = 1.380649e-23_pr, & !  $k = 1.380\,649 \times 10^{-23} \text{ J}\cdot\text{K}^{-1}$ 
    stefan_boltzmann = 5.670374419e-8_pr !  $\sigma = 5.670\,374\,419\dots \times 10^{-8} \text{ W}\cdot\text{m}^{-2}\cdot\text{K}^{-4}$ 
real (pr) :: sigma, sigma_bis !  $\sigma = (2\pi^5 / 15) k^4 / (c^2 h^3)$ 

sigma      = 2 * pi**5 / 15 * boltzmann**4 / (lightspeed**2 * planck**3)
sigma_bis  = 2 * pi**5 / 15 * boltzmann * (boltzmann / planck)**3 / &
    lightspeed**2

print *, sigma          ! NaN
print *, sigma_bis      ! 5.67037581e-8 = (1+2 $\epsilon$ ) $\sigma$ 
print *, stefan_boltzmann !  $\sigma = 5.67037439e-8$ 
                        ! 5.670374419...  $\times 10^{-8} \text{ W}\cdot\text{m}^{-2}\cdot\text{K}^{-4}$ 

print *, (sigma_bis / stefan_boltzmann - 1)
print *, (sigma_bis / stefan_boltzmann - 1) / epsilon (stefan_boltzmann)
end program stefan_constant
```



# Anyway, what physics needs double? ... fp16

```

program stefan_constant ! ratio of total power radiated per unit surface area
                        ! to fourth power of absolute temperature of black body
use, intrinsic :: iso_fortran_env, only: real32, real64 !, real128
integer, parameter :: real16 = 2, pr = real16
real (pr), parameter ::
  pi      = 3.141592653589793_pr      !  $\pi = \text{Arccos}(-1)$ 
real (pr) ::
  lightspeed = 299792458.0_pr,      & !  $c = 299\,792\,458 \text{ m}\cdot\text{s}^{-1}$ 
  planck     = 6.62607015e-34_pr,   & !  $h = 6.626\,070\,15 \times 10^{-34} \text{ J}\cdot\text{Hz}^{-1}$ 
  boltzmann  = 1.380649e-23_pr,     & !  $k = 1.380\,649 \times 10^{-23} \text{ J}\cdot\text{K}^{-1}$ 
  stefan_boltzmann = 5.670374419e-8_pr !  $\sigma = 5.670\,374\,419... \times 10^{-8} \text{ W}\cdot\text{m}^{-2}\cdot\text{K}^{-4}$ 
real (pr) :: sigma, sigma_bis      !  $\sigma = (2\pi^5 / 15) k^4 / (c^2 h^3)$ 

sigma      = 2 * pi**5 / 15 * boltzmann**4 / (lightspeed**2 * planck**3)
sigma_bis  = 2 * pi**5 / 15 * boltzmann * (boltzmann / planck)**3 / &
  lightspeed**2

print *, sigma      ! gfortran ! ifx ! nvfortran ! nv fp16
print *, sigma_bis ! NaN ! NaN ! NaN ! NaN
print *, stefan_boltzmann !  $\sigma = 5.67037581\text{e-}8$  !  $5.6703765\text{e-}8$  !  $5.6703765\text{e-}8$  !  $5.9605\text{e-}8 \Rightarrow$  smallest denorma

print *, 'c: ', lightspeed      ! c:      Inf
print *, 'h: ', planck          ! h:      0.000
print *, 'k: ', boltzmann       ! k:      0.000
print *, 'e: ', epsilon (stefan_boltzmann) ! e:      9.7656E-4
print *, '2: ', tiny (stefan_boltzmann)
print *, '3: ', tiny (stefan_boltzmann) * epsilon (stefan_boltzmann)
end program stefan_constant

```



By way of exception in base 10 (not in binary)! mantissa: 3 decimal digits

For  $a = 3.34$  and  $b = 3.33$

- $a \ominus b = 0.01 \Rightarrow$  **cancellation** (reducing relative accuracy)  
but a **benign** one (the floating point result is exact:  $a \ominus b = a - b$ )
- $$\begin{cases} a^2 - b^2 & = 0.0667 = 6.67 \times 10^{-2} \\ a \otimes a \ominus b \otimes b & = 0.1 = 1.00 \times 10^{-1} \end{cases}$$
 50% of relative error on the result, or  
333 ulp, no digit is even correct: **catastrophic cancellation**
- When does this occur?
- How many digits are lost?

Plus, there is an **overflow** risk

$\Rightarrow$  Let's **factorize** this!

$$(a \oplus b) \otimes (a \ominus b) = 6.67 \otimes 0.01 = 6.67 \times 10^{-2} \quad \text{exact}$$

$\Rightarrow$  The Right Way™



## Catastrophic Cancellation? Compensation Calamiteuse ?

⚠ Now in binary! With  $\text{fp}_{16}$ , mantissa: 10 binary digits

For  $a = 0 \times 1.6a8p + 0 = 1 + \frac{213}{512} = 1.4160\dots$  and  $b = 0 \times 1.6a4p + 0 = 1 + \frac{425}{1024} = 1.4150\dots$

- $a \ominus b = 0 \times 1p - 10 = 1 \epsilon = \frac{1}{1024} = 9.765\dots e - 4 \Rightarrow$  **cancellation** (reducing relative accuracy)

but a **benign** one (the floating point result is exact:  $a \ominus b = a - b$ )

- $\begin{cases} a^2 - b^2 & = \frac{2899}{2^{20}} = 2.7647\dots e - 3 \\ a \otimes a \ominus b \otimes b = 0 \times 1p - 8 & = \frac{1}{2^{56}} = 3.9063\dots e - 3 \end{cases}$

41 % of relative error, or 422 ulp, no digit is even correct: **catastrophic cancellation**

$\Rightarrow$  Let's **factorize** this!

$$(a \oplus b) \otimes (a \ominus b) = \frac{725}{256} \otimes \frac{1}{1024} = 0 \times 1.6a8p - 9 = \frac{725}{2^{18}} = 2.7657\dots e - 3 \quad \text{within } 0.35 \text{ ulp}$$

$\Rightarrow$  The Right Way™



# Cursed Cancellation

⇒ higher degree polynomials can degrade high resolutions (cf RUMP)

$$P = 333.75y^6 + x^2(11x^2y^2 - y^6 - 121y^4 - 2) + 5.5y^8 + \frac{x}{2y}$$

with  $x = 77617$ ,  $y = 33096$  (coprime integers)

[S.M. RUMP, 1988, "Algorithms for Verified Inclusions..." <https://www.tuhh.de/ti3/paper/rump/Ru88a.pdf>]



# Cursed Cancellation

⇒ higher degree polynomials can degrade high resolutions (*cf* RUMP)

$$P = 333.75y^6 + x^2(11x^2y^2 - y^6 - 121y^4 - 2) + 5.5y^8 + \frac{x}{2y}$$

with  $x = 77617$ ,  $y = 33096$  (coprime integers)

[S.M. RUMP, 1988, "Algorithms for Verified Inclusions..." <https://www.tuhh.de/ti3/paper/rump/Ru88a.pdf>]

float:  $P = -6.33825300e + 29$



# Cursed Cancellation

⇒ higher degree polynomials can degrade high resolutions (*cf* RUMP)

$$P = 333.75y^6 + x^2(11x^2y^2 - y^6 - 121y^4 - 2) + 5.5y^8 + \frac{x}{2y}$$

with  $x = 77617$ ,  $y = 33096$  (coprime integers)

[S.M. RUMP, 1988, "Algorithms for Verified Inclusions..." <https://www.tuhh.de/ti3/paper/rump/Ru88a.pdf>]

float:  $P = -6.33825300e + 29$   
double:  $P = -1.1805916207174113e + 021$



# Cursed Cancellation

⇒ higher degree polynomials can degrade high resolutions (*cf* RUMP)

$$P = 333.75y^6 + x^2(11x^2y^2 - y^6 - 121y^4 - 2) + 5.5y^8 + \frac{x}{2y}$$

with  $x = 77617$ ,  $y = 33096$  (coprime integers)

[S.M. RUMP, 1988, "Algorithms for Verified Inclusions..." <https://www.tuhh.de/ti3/paper/rump/Ru88a.pdf>]

float:  $P = -6.33825300e + 29$   
double:  $P = -1.1805916207174113e + 021$   
long double:  $P = +5.76460752303423489188e + 17$



# Cursed Cancellation

⇒ higher degree polynomials can degrade high resolutions (*cf* RUMP)

$$P = 333.75y^6 + x^2(11x^2y^2 - y^6 - 121y^4 - 2) + 5.5y^8 + \frac{x}{2y}$$

with  $x = 77617$ ,  $y = 33096$  (coprime integers)

[S.M. RUMP, 1988, "Algorithms for Verified Inclusions:..." <https://www.tuhh.de/ti3/paper/rump/Ru88a.pdf>]

float:	$P = -6.33825300e + 29$
double:	$P = -1.1805916207174113e + 021$
long double:	$P = +5.76460752303423489188e + 17$
quad:	$P = +1.17260394005317863185883490452018380$



# Cursed Cancellation

⇒ higher degree polynomials can degrade high resolutions (*cf* RUMP)

$$P = 333.75y^6 + x^2(11x^2y^2 - y^6 - 121y^4 - 2) + 5.5y^8 + \frac{x}{2y}$$

with  $x = 77617$ ,  $y = 33096$  (coprime integers)

[S.M. RUMP, 1988, "Algorithms for Verified Inclusions:..." <https://www.tuhh.de/ti3/paper/rump/Ru88a.pdf>]

float:	$P = -6.33825300e + 29$
double:	$P = -1.1805916207174113e + 021$
long double:	$P = +5.76460752303423489188e + 17$
quad:	$P = +1.17260394005317863185883490452018380$
fp16:	$P = \text{NaN}$



# Cursed Cancellation

⇒ higher degree polynomials can degrade high resolutions (cf RUMP)

$$P = 333.75y^6 + x^2(11x^2y^2 - y^6 - 121y^4 - 2) + 5.5y^8 + \frac{x}{2y}$$

with  $x = 77617$ ,  $y = 33096$  (coprime integers)

[S.M. RUMP, 1988, "Algorithms for Verified Inclusions..." <https://www.tuhh.de/ti3/paper/rump/Ru88a.pdf>]

float:	$P = -6.33825300e + 29$
double:	$P = -1.18059162071741113e + 021$
long double:	$P = +5.76460752303423489188e + 17$
quad:	$P = +1.17260394005317863185883490452018380$
fp16:	$P = \text{NaN}$
exact:	$P \approx -0.827396059946821368141165095479816292$
	$P = -\frac{54767}{66192}$

- the result can be represented accurately in any precision

even in fp16:  $-0 \times 1.a7cp-1 = -\frac{1695}{2048} = -0.8276 \dots$

- the result can't be computed in any precision

## How to control rounding errors?



# Cursed Cancellation

⇒ higher degree polynomials can degrade high resolutions (*cf* RUMP)

$$P = 333.75y^6 + x^2(11x^2y^2 - y^6 - 121y^4 - 2) + 5.5y^8 + \frac{x}{2y}$$

with  $x = 77617$ ,  $y = 33096$  (coprime integers)

[S.M. RUMP, 1988, "Algorithms for Verified Inclusions..." <https://www.tuhh.de/ti3/paper/rump/Ru88a.pdf>]

float:	$P = -6.33825300e + 29$
double:	$P = -1.1805916207174113e + 021$
long double:	$P = +5.76460752303423489188e + 17$
quad:	$P = +1.17260394005317863185883490452018380$
fp16:	$P = \text{NaN}$
exact:	$P \approx -0.827396059946821368141165095479816292$
	$P = -\frac{54767}{66192}$

- the result can be represented accurately in any precision

even in fp16:  $-0 \times 1.a7cp-1 = -\frac{1695}{2048} = -0.8276 \dots$

- the result can't be computed in any precision

## How to control rounding errors?

dec32:  $-4.000000e + 30$

dec64:  $-3.0000000000000000e + 21$

dec128: 1.172603940053178631858834904520184



# Cursed Cancellation

⇒ higher degree polynomials can degrade high resolutions (*cf* RUMP)

$$P = 333.75y^6 + x^2(11x^2y^2 - y^6 - 121y^4 - 2) + 5.5y^8 + \frac{x}{2y}$$

with  $x = 77617$ ,  $y = 33096$  (coprime integers)

[S.M. RUMP, 1988, "Algorithms for Verified Inclusions..." <https://www.tuhh.de/ti3/paper/rump/Ru88a.pdf>]

float:	$P = -6.33825300e + 29$	$-0 \times 1p + 99$
double:	$P = -1.1805916207174113e + 021$	$-0 \times 1p + 70$
long double:	$P = +5.76460752303423489188e + 17$	$+0 \times 8.00000000000000013p + 56$
quad:	$P = +1.17260394005317863185883490452018380$	$+0 \times 9.617e2cad835f5bap - 3$
fp16:	$P = \text{NaN}$	
exact:	$P \approx -0.827396059946821368141165095479816292$	
	$P = -\frac{54767}{66192}$	

- the result can be represented accurately in any precision

even in fp16:  $-0 \times 1.a7cp-1 = -\frac{1695}{2048} = -0.8276 \dots$

- the result can't be computed in any precision

## How to control rounding errors?

dec32:  $-4.000000e + 30$

dec64:  $-3.0000000000000000e + 21$

dec128: 1.172603940053178631858834904520184



# Cursed Cancellation

⇒ higher degree polynomials can degrade high resolutions (*cf* RUMP)

$$P = 333.75y^6 + x^2(11x^2y^2 - y^6 - 121y^4 - 2) + 5.5y^8 + \frac{x}{2y}$$

with  $x = 77617$ ,  $y = 33096$  (coprime integers)

[S.M. RUMP, 1988, "Algorithms for Verified Inclusions..." <https://www.tuhh.de/ti3/paper/rump/Ru88a.pdf>]

float:	$P = -6.33825300e + 29$	$-0 \times 1p + 99$
double:	$P = -1.18059162071741113e + 021$	$-0 \times 1p + 70$
long double:	$P = +5.76460752303423489188e + 17$	$+0 \times 8.0000000000000013p + 56$
quad:	$P = +1.17260394005317863185883490452018380$	$+0 \times 9.617e2cad835f5bap - 3 = \frac{x}{2y}$
fp16:	$P = \text{NaN}$	
exact:	$P \approx -0.827396059946821368141165095479816292$	
	$P = -\frac{54767}{66192}$	

- the result can be represented accurately in any precision

even in *fp16*:  $-0 \times 1.a7cp-1 = -\frac{1695}{2048} = -0.8276 \dots$

- the result can't be computed in any precision

## How to control rounding errors?

dec32:  $-4.000000e + 30$

dec64:  $-3.0000000000000000e + 21$

dec128:  $1.172603940053178631858834904520184 = \frac{x}{2y}$



# Cursed Cancellation

⇒ higher degree polynomials can degrade high resolutions (cf RUMP)

$$\begin{aligned} P &= 333.75y^6 + x^2(11x^2y^2 - y^6 - 121y^4 - 2) + 5.5y^8 + \frac{x}{2y} \\ &= \frac{1}{4}(1335y^6 + 4x^2(11x^2y^2 - y^6 - 121y^4 - 2) + 22y^8) + \frac{x}{2y} \\ &= \frac{1}{4}(1335Y^3 + 4X(11XY - Y^3 - 121Y^2 - 2) + 22Y^4) + \frac{x}{2y} \quad X = x^2, Y = y^2 \end{aligned}$$

with  $x = 77617$ ,  $y = 33096$  (coprime integers)

Factor the polynomial part?

⇒ No way! Irreducible polynomial in  $\mathbb{R}$ ,  $\mathbb{C}$ ,  $\mathbb{P}\mathbb{C}$  (checked with `macaulay2` M2)



## HORNER-RUFFINI

- computational cost of all the exponentiation,
- accuracy loss it represents.

$$\begin{aligned} p(x) &= a_0 + a_1x + \dots + a_{n-1}x^{n-1} + a_nx^n \\ &= a_n(x - x_1) \times \dots \times (x - x_n) \quad (\mathbb{C}) \\ &= a_n(x - x_1) \times \dots \times (x^2 + b_mx + c_m) \quad (\mathbb{R}) \\ &= a_0 + x(a_1 + x(\dots + x(a_{n-1} + xa_n)\dots)) \end{aligned} \tag{1}$$

- gain in speed, (saving of operations)
  - also in accuracy, partly for the same reason,
  - guarantee of stability of the result and safety against intermediate overshoots
- ⇒ “multiply-accumulate” machine instructions (fma).
- + compensation summation techniques, such as the summation algorithm of W. KAHAN



- A difference...
  - ...of squares...
  - ...of sums...
  - ...and sums...
  - ...of squares...
- 
- two passes approach
  - arbitrary data shift towards some expected average value
  - 1-pass online Welford's algorithm (one more division per iteration)



# Typical computation

- dot product
- convolution product (“backwards” dot product)
- Fourier transform
- matrix product is a matrix of dot products

turns out to be a sum of simple products (quadratic in essence).

⇒ we expect to encounter problems similar to difference of squares and variance computation.  
But here we can't use the factorisation trick...

- mixed precision
- `fma` (fused multiply accumulate)
- `fma` used to extract exact product
- combined with Kahan or other compensated sums



# Comparing floats

Comparing starts by subtracting

$$a = b \Leftrightarrow a - b = 0 \stackrel{?}{\Leftrightarrow} a - b == 0$$

so in order to compare, we look for the ultimate catastrophic cancelation  
But the numbers we compare are fraught with error.

$$a - b = a \ominus b + \delta_{\text{rounding}} + \delta a - \delta b$$

For instance:

- $a == 0.2 \neq \frac{1}{5}$ : representation error  $\Rightarrow 5 * a == 1$ ?
- $a == \text{pi} \neq \pi$ : representation error
- $a == 2$ : no representation errors, but there are rounding errors on  $a$

... and in the case of floating-point numbers, it's always a relative error.  
and so comparison with an absolute threshold is not the most relevant.

$$a = b \Leftrightarrow |a - b| < \eta$$

$\Rightarrow \eta$  can't be arbitrarily small  
and so we look for relative comparisons

$$a = b \Leftrightarrow |a - b| < \epsilon \cdot \max(\text{abs}(a), \text{abs}(b))$$

where  $\epsilon$  will be a multiple of the  $\epsilon$  machine.



# Comparing floats to zero

$$a - 0 = a + \delta a$$

when we look for relative comparisons to zero

$$a == 0 \Leftrightarrow \text{abs}(a) < \epsilon \cdot \text{abs}(a)!!!$$

... we get an impossible constraint that only an exact zero can match.

⇒ we are back to absolute uncertainties

$$a == 0 \Leftrightarrow |a| < \eta$$

⇒  $\eta$  can't be arbitrarily small:

each floating point type has a tiniest positive non-null normal value

... as well as a tiniest positive non-null denormal

... which is not epsilon machine

fp	$\epsilon$	tiny	1/Huge	$\epsilon \times \text{tiny}$
fp16	9.766...e-4	6.104...e-5	1.526...e-5	5.961...e-8
fp32	1.192...e-7	1.175...e-38	2.939...e-39	1.401...e-45
fp64	2.220...e-16	2.225...e-308	5.563...e-309	4.941...e-324
fp80	1.084...e-19	3.362...e-4932	8.405...e-4933	3.645...e-4951
fp128	1.926...e-34	3.362...e-4932	8.405...e-4933	6.475...e-4966

Table – *Lilliputian floats*



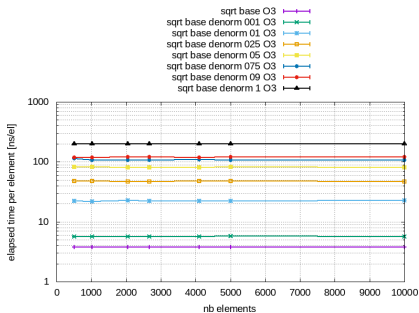
What is the meaning of this comparison?

- test whether the argument is in the domain of the function
  - ▶  $a \neq 0$  to avoid computing  $1/a$
  - ▶  $a \geq 0$  to avoid computing  $\text{sqrt}(a)$
  - ▶  $a > 0$  to avoid computing  $\log(a)$
  - ▶  $\text{abs}(a) \leq 1$  to avoid computing  $\text{asin}(a)$
- (unit) test expected values of a function at a given point is the value representable? is the argument representable?
  - ▶  $\sin \pi = 0$ ?



# Denormals / Subnormals

- below  $1.17 \times 10^{-38}$  for fp32
- below  $2.22 \times 10^{-308}$  for fp64
- below  $6.09 \times 10^{-5}$  for fp16  
(up to  $5.96 \times 10^{-8}$ )
- Why?  
⇒ allow for “gradual underflow”
- Why not?  
⇒ 100× slower  
(thanks Pierre AUBERT)
- How?
  - ▶ float difference around the minimum normal threshold
  - ▶ decreasing geometric progression  
(thanks Hadrien GRASLAND!)



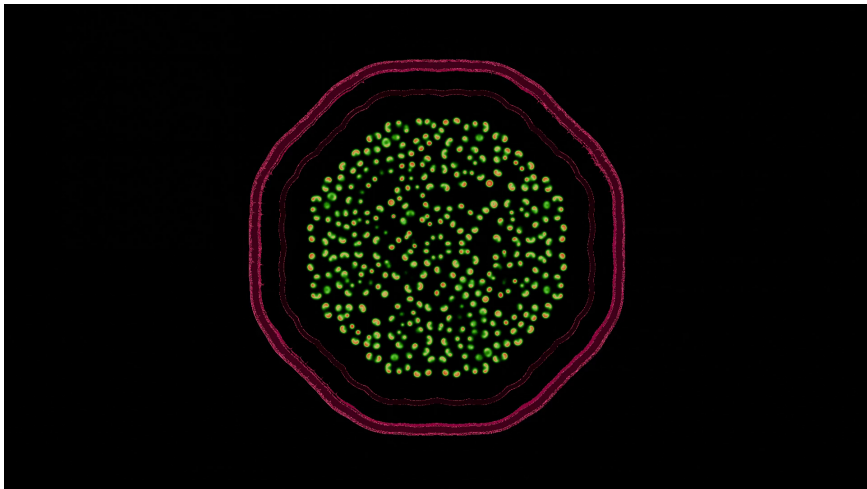


## Movie time

(thanks Pierre AUBERT & Sébastien VALAT)



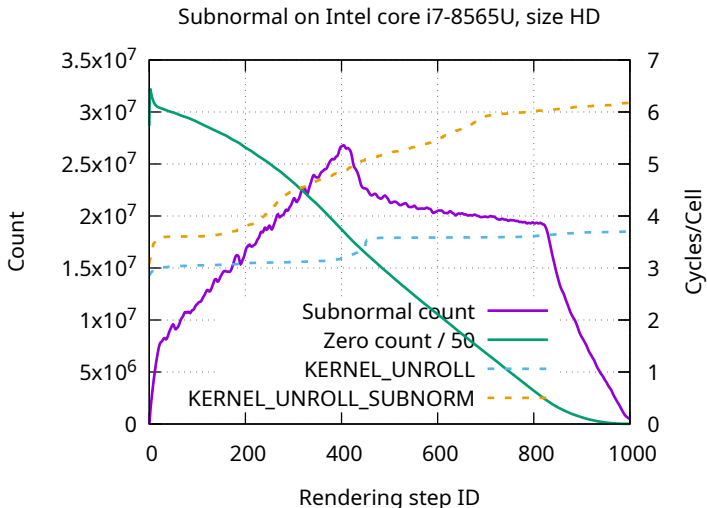
# Denormals / Subnormals





# Denormals / Subnormals

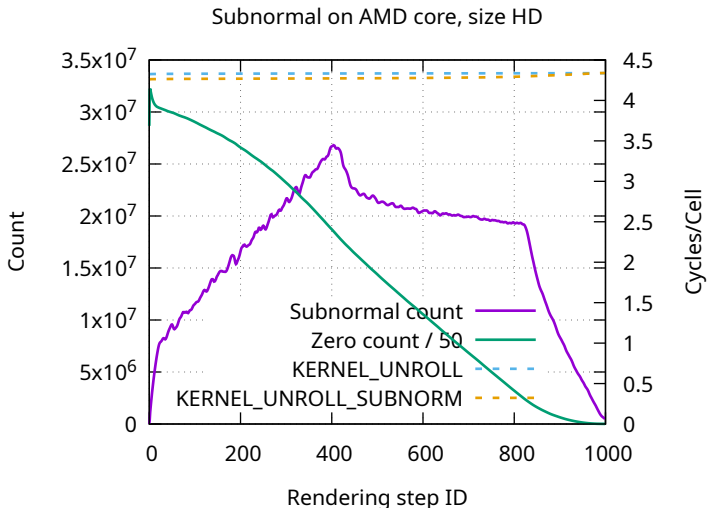
(thanks Pierre AUBERT & Sébastien VALAT)





# Denormals / Subnormals

(thanks Pierre AUBERT & Sébastien VALAT)





# Denormals / Subnormals

Some math:  $\Delta u = \vec{\nabla}^2 u$

Wave equation (D'ALEMBERT)

$$\frac{\partial^2}{c^2 \partial t^2} u = \Delta u$$

Heat equation (FOURIER)

$$\frac{\partial}{\partial t} u = \alpha \Delta u$$

Thermal diffusivity:

$$\alpha = \frac{\lambda}{\rho c_p} \quad (\text{m}^2 \text{s}^{-1})$$



# Denormals / Subnormals

Some math:  $u = e^{i\omega t} e^{i\mathbf{k}\cdot\mathbf{x}}$

Wave equation (D'ALEMBERT)

$$\frac{\partial^2}{c^2 \partial t^2} u = \Delta u$$

$$\frac{\omega^2}{c^2} = \mathbf{k} \cdot \mathbf{k} = k^2$$

$$k = \pm \frac{\omega}{c}$$

Heat equation (FOURIER)

$$\frac{\partial}{\partial t} u = \alpha \Delta u$$

$$i\omega = -\alpha k^2$$

$$k = \pm \frac{1-i}{2} \sqrt{\frac{2\omega}{\alpha}}$$

⇒ decaying sine wave !



Gray Scott reaction:

$$\begin{aligned}\frac{\partial u}{\partial t} &= r_u \nabla^2 u - uv^2 + f_r(1 - u) \\ \frac{\partial v}{\partial t} &= r_v \nabla^2 v + uv^2 - (f_r + k_r)v\end{aligned}$$

- Starting from the drop of  $v$  in the middle of  $u$ ...
- ... the linear part of the equation will send away decaying waves of  $v$
- ... (considering all the frequencies in the sudden impulse of  $v$ )
- ... and the front wave of  $v$  is exponentially small
- ... reaching subnormal values



# Quadratic

$$\begin{aligned}ax^2 + bx + c &= 0 \quad (a \neq 0) \\ \Delta &= b^2 - 4ac \\ x_{\pm} &= \frac{-b \pm \sqrt{\Delta}}{2a}\end{aligned}$$

2 possible *catastrophic cancelation* (« *compensations calamiteuses* »)

- $-b$  &  $\sqrt{\Delta}$

$$\Rightarrow q = -b - \text{sgn}(b)\sqrt{\Delta} = -\text{sgn}(b) (|b| + \sqrt{\Delta})$$

$$\begin{cases} x_1 = \frac{q}{2a} \\ x_2 = \frac{2c}{q} = \frac{c}{ax_1} \end{cases}$$

- discriminant  $\Delta = b^2 - 4ac \Rightarrow \text{fma}$

4 possible *overflow*:

- $b^2$  : *spurious overflow* (if  $|b| > 10^{19}$ ,  $\Delta = \text{Inf}$ ,  $|q| = \text{Inf}$  while  $|q| \sim 2 \times 10^{19}$ )
- $ac$ ,  $b/a$ ,  $c/b$

BAKER, Henry G., "You Could Learn a Lot from a Quadratic: Overloading Considered Harmful", doi:10.1145/609742.609746



Giulio FAGNANO (1682–1766)

$$ax^2 + bx + c = 0 \quad (a \neq 0)$$

$$Y = 1/x \quad \rightarrow \quad cY^2 + bY + a = 0$$

$$\Delta = b^2 - 4ac \quad \text{same invariant for both equations}$$

$$x_{\pm} = \frac{-b \pm \sqrt{\Delta}}{2a}$$

$$Y_{\pm} = \frac{-b \mp \sqrt{\Delta}}{2c}$$

$$x_{\pm} = \frac{2c}{-b \mp \sqrt{\Delta}}$$

Mentzel, Numerical Recipes 86



$$Q = \frac{\sqrt{|ac|}}{b} \quad \text{is dimensionless}$$

$$F = \frac{1}{2} (1 + \sqrt{1 - 4\sigma Q^2})$$

$$x_1 = -\frac{b}{a} F \quad x_2 = -\frac{c}{b} \frac{1}{F}$$

$$\text{for } \sigma = +1 \quad F_+ = \frac{1}{2} (1 + \sqrt{(1 - 2Q)(1 + 2Q)})$$

$$\text{for } \sigma = -1 \quad F_- = \frac{1}{2} (1 + \sqrt{\text{fma}(2Q, 2Q, 1)}) \quad Q \rightarrow +\infty \quad Q$$

$$F_- = \frac{Q}{2} \left( \frac{1}{Q} + \sqrt{\text{fma}\left(\frac{1}{Q}, \frac{1}{Q}, 4\right)} \right) \quad \forall Q > \frac{1}{2\sqrt{\epsilon}}$$

- Low entropy formula
- Importance of dimensional analysis (dimensionless numbers implementation)

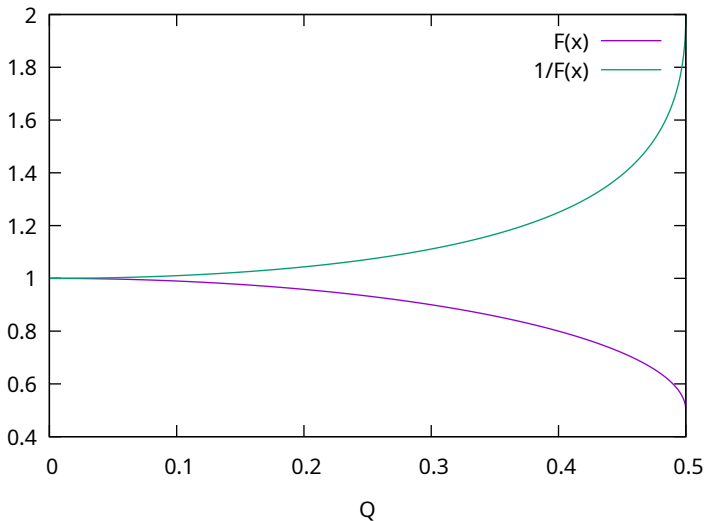


Figure –  $F = \frac{1}{2} (1 + \sqrt{1 - 4Q^2}) \forall Q \in [0; 1/2[$ .



$$F = \frac{1}{2} (1 + \sqrt{1 - 4Q^2}) \quad \forall Q \in [0; 1/2[$$

$$\kappa = -\frac{(1 - \sqrt{1 - 4Q^2})}{\sqrt{1 - 4Q^2}} \quad \forall Q \in [0; 1/2[$$

Table – Quadratic roots

A	-B/2	C	true $\Delta$	true roots	computed $\Delta$	computed roots
10.27	29.61	85.37	0.0022	2.88772... 2.87859...	0.1000	2.914 2.852 dec4
10.28	29.62	85.34	0.0492	2.90290... 2.86075...	0	2.881 2.881 dec4
10.27	29.61	85.37	0.0022	2.88772... 2.87859...	0.1000	2.883 fp16
10.28	29.62	85.34	0.0492	2.90290... 2.86075...	0	2.881 fp16
94906265.625	94906267.000	94906268.375	1.89...	1.000000028975958... 1.0	0.0	1.000000014487979 1.000000014487979
94906266.375	94906267.375	94906268.375	1.0	1.000000021073424... 1.0	2.0	1.000000025437873 0.999999995635551

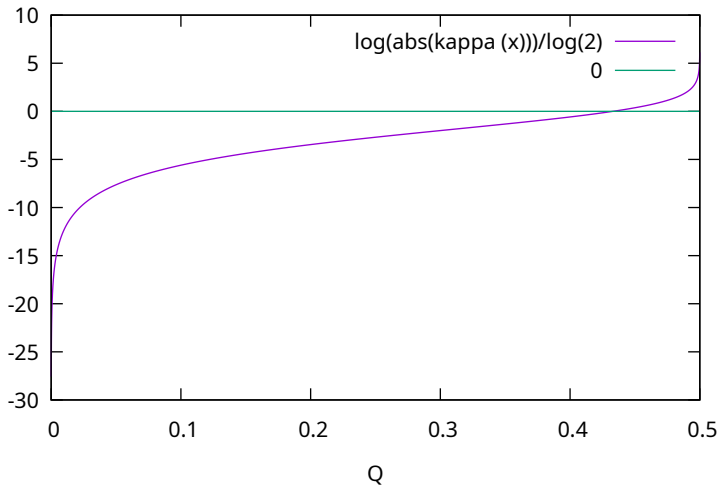


Figure –  $\kappa = \frac{(1-\sqrt{1-4Q^2})}{\sqrt{1-4Q^2}} \quad \forall Q \in [0; 1/2[.$



## Dimensional analysis, split

- scale parameters, or problem's **characteristic scales**
- ... **dimensionless** shape parameters (pure numbers)
- lower formulas entropy
- often many ways to do it
  - ▶ problem's symmetries,
  - ▶ limit computation complexity,
  - ▶ limit computation exceptions.
- if the math solution has no float representation, we should allow intermediate results not to be representable as well
- bring values close to unity  
**where the floating point density is highest!**



# COMPTON Scattering

$$\theta = \arccos \left[ 1 + m_e c^2 \left( \frac{1}{E_1 + E_2} - \frac{1}{E_2} \right) \right]$$

- 2 (same sign) substractions

$$\theta = \arccos \left[ 1 - m_e c^2 \left( \frac{1}{E_2} - \frac{1}{E_1 + E_2} \right) \right]$$

- basic algebra:

$$\theta = \arccos \left[ 1 - \frac{m_e c^2 E_1}{E_2 (E_1 + E_2)} \right] \quad \frac{m_e c^2 E_1}{E_2 (E_1 + E_2)} \text{ is dimensionless}$$

- ...one remaining (same sign) subtraction
- basic trigonometry:  $\cos 2\alpha = 1 - 2 \sin^2 \alpha \Leftrightarrow \sin^2 \frac{\theta}{2} = \frac{1 - \cos \theta}{2} = \frac{\text{versin } \theta}{2} = \text{haversin } \theta$

$$\theta = 2 \arcsin \sqrt{\frac{m_e c^2 E_1}{2 E_2 (E_1 + E_2)}}$$

- ...no remaining (same sign) subtraction



$$\theta = \arccos \left[ 1 + m_e c^2 \left( \frac{1}{E_1 + E_2} - \frac{1}{E_2} \right) \right]$$

- 2 (same sign) substractions

$$\theta = \arccos \left[ 1 - m_e c^2 \left( \frac{1}{E_2} - \frac{1}{E_1 + E_2} \right) \right]$$

- basic algebra:

$$\theta = \arccos \left[ 1 - \frac{m_e c^2 E_1}{E_2 (E_1 + E_2)} \right] \quad \frac{m_e c^2 E_1}{E_2 (E_1 + E_2)} \text{ is dimensionless}$$

- ...one remaining (same sign) subtraction

- basic trigonometry:  $\cos 2\alpha = 1 - 2 \sin^2 \alpha \Leftrightarrow \sin^2 \frac{\theta}{2} = \frac{1 - \cos \theta}{2} = \frac{\text{versin } \theta}{2} = \text{haversin } \theta$

$$\theta = 2 \arcsin \sqrt{\frac{m_e c^2 E_1}{2 E_2 (E_1 + E_2)}}$$

- ...no remaining (same sign) subtraction



$$\theta = \arccos \left[ 1 + m_e c^2 \left( \frac{1}{E_1 + E_2} - \frac{1}{E_2} \right) \right]$$

- 2 (same sign) subtractions

$$\theta = \arccos \left[ 1 - m_e c^2 \left( \frac{1}{E_2} - \frac{1}{E_1 + E_2} \right) \right]$$

- basic algebra:

$$\theta = \arccos \left[ 1 - \frac{m_e c^2 E_1}{E_2 (E_1 + E_2)} \right] \quad \frac{m_e c^2 E_1}{E_2 (E_1 + E_2)} \text{ is dimensionless}$$

- ...one remaining (same sign) subtraction
- basic trigonometry:  $\cos 2\alpha = 1 - 2 \sin^2 \alpha \Leftrightarrow \sin^2 \frac{\theta}{2} = \frac{1 - \cos \theta}{2} = \frac{\text{versin } \theta}{2} = \text{haversin } \theta$

$$\theta = 2 \arcsin \sqrt{\frac{m_e c^2 E_1}{2 E_2 (E_1 + E_2)}}$$

- ...no remaining (same sign) subtraction



$$\theta = \arccos \left[ 1 + m_e c^2 \left( \frac{1}{E_1 + E_2} - \frac{1}{E_2} \right) \right]$$

- 2 (same sign) subtractions

$$\theta = \arccos \left[ 1 - m_e c^2 \left( \frac{1}{E_2} - \frac{1}{E_1 + E_2} \right) \right]$$

- basic algebra:

$$\theta = \arccos \left[ 1 - \frac{m_e c^2 E_1}{E_2 (E_1 + E_2)} \right] \quad \frac{m_e c^2 E_1}{E_2 (E_1 + E_2)} \text{ is dimensionless}$$

- ...one remaining (same sign) subtraction

- basic trigonometry:  $\cos 2\alpha = 1 - 2 \sin^2 \alpha \Leftrightarrow \sin^2 \frac{\theta}{2} = \frac{1 - \cos \theta}{2} = \frac{\text{versin } \theta}{2} = \text{haversin } \theta$

$$\theta = 2 \arcsin \sqrt{\frac{m_e c^2 E_1}{2 E_2 (E_1 + E_2)}}$$

- ...no remaining (same sign) subtraction



$$\theta = \arccos \left[ 1 + m_e c^2 \left( \frac{1}{E_1 + E_2} - \frac{1}{E_2} \right) \right]$$

- 2 (same sign) subtractions

$$\theta = \arccos \left[ 1 - m_e c^2 \left( \frac{1}{E_2} - \frac{1}{E_1 + E_2} \right) \right]$$

- basic algebra:

$$\theta = \arccos \left[ 1 - \frac{m_e c^2 E_1}{E_2 (E_1 + E_2)} \right] \quad \frac{m_e c^2 E_1}{E_2 (E_1 + E_2)} \text{ is dimensionless}$$

- ...one remaining (same sign) subtraction

- basic trigonometry:  $\cos 2\alpha = 1 - 2 \sin^2 \alpha \Leftrightarrow \sin^2 \frac{\theta}{2} = \frac{1 - \cos \theta}{2} = \frac{\text{versin } \theta}{2} = \text{haversin } \theta$

$$\theta = 2 \arcsin \sqrt{\frac{m_e c^2 E_1}{2 E_2 (E_1 + E_2)}}$$

- ...no remaining (same sign) subtraction



# Compound interest

$$A = P \left(1 + \frac{r}{n}\right)^{nt}$$

$$I = P \left(1 + \frac{r}{n}\right)^{nt} - P$$

$$I = P \left[ \left(1 + \frac{r}{n}\right)^{nt} - 1 \right]$$

$$I = P \left[ \text{pow} \left( \left(1 + \frac{r}{n}\right), nt \right) - 1 \right]$$

$$I = P \left[ \exp \left( nt \ln \left(1 + \frac{r}{n}\right) \right) - 1 \right]$$

$$I = P \left[ \exp \left( nt \log_{1p} \left( \frac{r}{n} \right) \right) - 1 \right]$$

$$I = P \left[ \text{expm1} \left( nt \log_{1p} \left( \frac{r}{n} \right) \right) \right]$$



# Compound interest

$$A = P \left(1 + \frac{r}{n}\right)^{nt}$$

$$I = P \left(1 + \frac{r}{n}\right)^{nt} - P$$

$$I = P \left[ \left(1 + \frac{r}{n}\right)^{nt} - 1 \right]$$

$$I = P \left[ \text{pow} \left( \left(1 + \frac{r}{n}\right), nt \right) - 1 \right]$$

$$I = P \left[ \exp \left( nt \ln \left(1 + \frac{r}{n}\right) \right) - 1 \right]$$

$$I = P \left[ \exp \left( nt \log_{1p} \left( \frac{r}{n} \right) \right) - 1 \right]$$

$$I = P \left[ \text{expm1} \left( nt \log_{1p} \left( \frac{r}{n} \right) \right) \right]$$



# Compound interest

$$A = P \left(1 + \frac{r}{n}\right)^{nt}$$

$$I = P \left(1 + \frac{r}{n}\right)^{nt} - P$$

$$I = P \left[ \left(1 + \frac{r}{n}\right)^{nt} - 1 \right]$$

$$I = P \left[ \text{pow} \left( \left(1 + \frac{r}{n}\right), nt \right) - 1 \right]$$

$$I = P \left[ \exp \left( nt \ln \left(1 + \frac{r}{n}\right) \right) - 1 \right]$$

$$I = P \left[ \exp \left( nt \log_{1p} \left( \frac{r}{n} \right) \right) - 1 \right]$$

$$I = P \left[ \text{expm1} \left( nt \log_{1p} \left( \frac{r}{n} \right) \right) \right]$$



# Compound interest

$$A = P \left(1 + \frac{r}{n}\right)^{nt}$$

$$I = P \left(1 + \frac{r}{n}\right)^{nt} - P$$

$$I = P \left[ \left(1 + \frac{r}{n}\right)^{nt} - 1 \right]$$

$$I = P \left[ \text{pow} \left( \left(1 + \frac{r}{n}\right), nt \right) - 1 \right]$$

$$I = P \left[ \exp \left( nt \ln \left(1 + \frac{r}{n}\right) \right) - 1 \right]$$

$$I = P \left[ \exp \left( nt \log_{10} \left( \frac{r}{n} \right) \right) - 1 \right]$$

$$I = P \left[ \text{expm1} \left( nt \log_{10} \left( \frac{r}{n} \right) \right) \right]$$



# Compound interest

$$A = P \left(1 + \frac{r}{n}\right)^{nt}$$

$$I = P \left(1 + \frac{r}{n}\right)^{nt} - P$$

$$I = P \left[ \left(1 + \frac{r}{n}\right)^{nt} - 1 \right]$$

$$I = P \left[ \text{pow} \left( \left(1 + \frac{r}{n}\right), nt \right) - 1 \right]$$

$$I = P \left[ \text{exp} \left( nt \ln \left(1 + \frac{r}{n}\right) \right) - 1 \right]$$

$$I = P \left[ \text{exp} \left( nt \log_{1p} \left( \frac{r}{n} \right) \right) - 1 \right]$$

$$I = P \left[ \text{expm1} \left( nt \log_{1p} \left( \frac{r}{n} \right) \right) \right]$$



# Compound interest

If  $\log_1 p$  is not available (*cf.* GOLDBERG)

$$\ln(1+x) = \begin{cases} x & \text{if } 1 \oplus x = 1 \\ \frac{x \ln(1+x)}{(1+x)-1} & \text{else.} \end{cases}$$

$$\exp(x) - 1 = \begin{cases} x & \text{if } 1 \oplus x = 1 \\ \frac{x(\exp(x)-1)}{(1+x)-1} & \text{else.} \end{cases}$$



# Area of triangle

area  $S$  as a function of lengths  $a$ ,  $b$  and  $c$  of edges

$$S = \sqrt{p(p-a)(p-b)(p-c)} \quad (\text{HERON of ALEXANDRIA, } \textit{Stereometrica})$$

$$p = \frac{a+b+c}{2} \quad \text{half-perimeter}$$

Symmetric, but numerically unstable, for needle-like triangles (when large and small values meet in the same formula)

KAHAN Re-labelling:  $a > b > c$

$$\frac{1}{4} \sqrt{[a + (b + c)] [c - (a - b)] [c + (a - b)] [a + (b - c)]}$$

Apparent Symmetry is lost, but the formula is way more robust  
Originating from a determinantal expression

$$S = \frac{1}{4} \sqrt{\begin{vmatrix} 0 & a^2 & b^2 & 1 \\ a^2 & 0 & c^2 & 1 \\ b^2 & c^2 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{vmatrix}}$$

⇒ exercise: code and test data from <https://people.eecs.berkeley.edu/~wkahan/Triangle.pdf>



# Volume of the tetrahedron

$$V = \sqrt{\frac{1}{288} \begin{vmatrix} 0 & a^2 & b^2 & c^2 & 1 \\ a^2 & 0 & C^2 & B^2 & 1 \\ b^2 & C^2 & 0 & A^2 & 1 \\ c^2 & B^2 & A^2 & 0 & 1 \\ 1 & 1 & 1 & 1 & 0 \end{vmatrix}}$$

$$X = (c - A + b)(A + b + c)$$

$$x = (A - b + c)(b - c + A)$$

$$Y = (a - B + c)(B + c + a)$$

$$y = (B - c + a)(c - a + B)$$

$$Z = (b - C + a)(C + a + b)$$

$$z = (C - a + b)(a - b + C)$$

$$\xi = \sqrt{xYZ} \quad \eta = \sqrt{yZX} \quad \zeta = \sqrt{zXY} \quad \lambda = \sqrt{xyz}$$

$$V = \frac{1}{192abc} \sqrt{(\xi + \eta + \zeta - \lambda)(\lambda + \xi + \eta - \zeta)(\eta + \zeta + \lambda - \xi)(\zeta + \lambda + \xi - \eta)}$$



$$P(x) = ax^2 + bx + c \quad (a \neq 0)$$

$$\Delta = b^2 - 4ac \quad \text{invariant... that looks much like a determinant}$$

$$\Delta = \text{Disc}(P) \quad \text{discriminant of polynomial } P$$

$$= \frac{(-1)^{n(n-1)/2}}{a_n} \text{Res}(P, P') \quad \text{resultant of polynomial } P \text{ and } P'$$

$$= \frac{(-1)^{n(n-1)/2}}{a_n} |M(P, P')| \quad \text{determinant of SYLVESTER matrix}$$

$$= a_n^{2n-2} \prod_{i < j} (r_i - r_j)^2 = (-1)^{n(n-1)/2} a_n^{2n-2} \prod_{i \neq j} (r_i - r_j)$$

$$\text{Res}(P, Q) = \begin{array}{c} \underbrace{\hspace{10em}}_{P(x) \times q \text{ copies}} \hspace{2em} \underbrace{\hspace{10em}}_{Q(x) \times p \text{ copies}} \\ \left| \begin{array}{cccccccc} a_n & 0 & \cdots & 0 & b_m & 0 & \cdots & 0 \\ a_{n-1} & a_n & \ddots & \vdots & \vdots & b_m & \ddots & \vdots \\ \vdots & a_{n-1} & \ddots & 0 & \vdots & \vdots & \ddots & 0 \\ \vdots & \vdots & \ddots & a_n & b_1 & \vdots & \vdots & b_m \\ a_0 & \vdots & \ddots & a_{n-1} & b_0 & \ddots & \vdots & \vdots \\ 0 & \ddots & \vdots & \vdots & 0 & \ddots & b_1 & \vdots \\ \vdots & \ddots & a_0 & \vdots & \vdots & \ddots & b_0 & b_1 \\ 0 & \cdots & 0 & a_0 & 0 & \cdots & 0 & b_0 \end{array} \right| \end{array} \quad \begin{array}{l} \text{deg } P(x) = p, \\ \text{deg } Q(x) = q \end{array}$$



$$\begin{aligned}P(x) &= ax^2 + bx + c \quad (a \neq 0) \\ \Delta &= b^2 - 4ac \quad \text{invariant... indeed a determinant} \\ \Delta &= \text{Disc}(P) \quad \text{discriminant of polynomial } P \\ &= a^2(r_1 - r_2)^2 \\ &= \frac{-1}{a} \begin{vmatrix} a & 2a & 0 \\ b & b & 2a \\ c & 0 & b \end{vmatrix}\end{aligned}$$

Invariant...under what?

$$\begin{aligned}Q(x) &= P(x + \delta) \\ &= a(x + \delta)^2 + b(x + \delta) + c \\ &= ax^2 + 2a\delta x + a\delta^2 + bx + b\delta + c \\ &= ax^2 + (b + 2a\delta)x + (c + b\delta + a\delta^2) \\ \Delta_Q &= b'^2 - 4a'c' \\ &= (b + 2a\delta)^2 - 4a(c + b\delta + a\delta^2) \\ &= (b^2 + 4ab\delta + 4a^2\delta^2) - 4ac - 4ab\delta - 4a^2\delta^2 \\ &= b^2 + 4ab\delta - 4ab\delta + 4a^2\delta^2 - 4a^2\delta^2 - 4ac \\ &= b^2 - 4ac \\ &= \Delta_P\end{aligned}$$



$$P(x) = ax^3 + bx^2 + cx + d = 0 \quad (a \neq 0)$$

$$\Delta = \text{Disc}(P) = \frac{1}{a} \text{Res}(P, P')$$

$$= \frac{1}{a} \begin{vmatrix} a & 0 & 3a & 0 & 0 \\ b & a & 2b & 3a & 0 \\ c & b & c & 2b & 3a \\ d & c & 0 & c & 2b \\ 0 & d & 0 & 0 & c \end{vmatrix}$$

$$= 18abcd - 4b^3d + b^2c^2 - 4ac^3 - 27a^2d^2$$

$$= \frac{4(b^2 - 3ac)^3 - (2b^3 - 9abc + 27a^2d)^2}{27a^2}$$

$$= \frac{4\Delta_0^3 - \Delta_1^2}{27a^2} = -(4p^3 + 27q^2)$$

$$\Delta_1 = \frac{-1}{8a} \text{Res}(P, P'')$$

$$= \frac{-1}{8a} \begin{vmatrix} a & 6a & 0 & 0 \\ b & 2b & 6a & 0 \\ c & 0 & 2b & 6a \\ d & 0 & 0 & 2b \end{vmatrix} = 2b^3 - 9abc + 27a^2d = 27q$$

$$\Delta_0 = \frac{-1}{12a} \text{Res}(P', P'')$$

$$= \frac{-1}{12a} \begin{vmatrix} 3a & 6a & 0 \\ 2b & 2b & 6a \\ c & 0 & 2b \end{vmatrix} = b^2 - 3ac = -3p$$



For the depressed cubic  $P(x) = x^3 + px + q = 0$   
when  $\Delta < 0$  (i.e.  $4p^3 + 27q^2 > 0$ ), one gets a single real root:

$$r = \sqrt[3]{\underbrace{-\frac{q}{2} + \sqrt{\frac{q^2}{4} + \frac{p^3}{27}}}_{u_1}} + \sqrt[3]{\underbrace{-\frac{q}{2} - \sqrt{\frac{q^2}{4} + \frac{p^3}{27}}}_{u_2}}$$



one of the two radicands can suffer catastrophic cancellation!



For the depressed cubic  $P(x) = x^3 + px + q = 0$   
 when  $\Delta < 0$  (i.e.  $4p^3 + 27q^2 > 0$ ), one gets a single real root:

$$r = \sqrt[3]{\underbrace{-\frac{q}{2} + \sqrt{\frac{q^2}{4} + \frac{p^3}{27}}}_{u_1}} + \sqrt[3]{\underbrace{-\frac{q}{2} - \sqrt{\frac{q^2}{4} + \frac{p^3}{27}}}_{u_2}}$$



one of the two radicands can suffer catastrophic cancellation!

$$\sqrt[3]{u_1 \cdot u_2} = \sqrt[3]{\left(\frac{q}{2}\right)^2 - \left(\frac{q^2}{4} + \frac{p^3}{27}\right)} = -\sqrt[3]{\frac{p^3}{27}} = -\frac{p}{3}$$

$$u_+ = -\left(\frac{q}{2} + \operatorname{sgn} q \sqrt{\frac{q^2}{4} + \frac{p^3}{27}}\right)$$

$$r = \sqrt[3]{u_+} - \frac{p/3}{\sqrt[3]{u_+}}$$



For the depressed cubic  $P(x) = x^3 + px + q = 0$   
 when  $\Delta < 0$  (i.e.  $4p^3 + 27q^2 > 0$ ), one gets a single real root:

$$r = \sqrt[3]{\underbrace{-\frac{q}{2} + \sqrt{\frac{q^2}{4} + \frac{p^3}{27}}}_{u_1}} + \sqrt[3]{\underbrace{-\frac{q}{2} - \sqrt{\frac{q^2}{4} + \frac{p^3}{27}}}_{u_2}}$$



one of the two radicands can suffer catastrophic cancellation!

$$\sqrt[3]{u_1 \cdot u_2} = \sqrt[3]{\left(\frac{q}{2}\right)^2 - \left(\frac{q^2}{4} + \frac{p^3}{27}\right)} = -\sqrt[3]{\frac{p^3}{27}} = -\frac{p}{3}$$

$$u_+ = -\left(\frac{q}{2} + \operatorname{sgn} q \sqrt{\frac{q^2}{4} + \frac{p^3}{27}}\right)$$

$$r = \sqrt[3]{u_+} - \frac{p/3}{\sqrt[3]{u_+}}$$



for  $p > 0$ , the two radicals can cancel each other! catastrophically?



For the depressed cubic  $P(x) = x^3 + px + q = 0$

with  $x = |q|^{1/3}Y$ ,  $P(x) = |q|Y^3 + p|q|^{1/3}Y + \text{sgn } q|q| = |q|Q(Y)$

$$Q(Y) = Y^3 + \frac{p}{|q|^{2/3}}Y + \text{sgn } q$$

$$= Y^3 + sY + \text{sgn } q \quad \text{where} \quad s = \frac{p}{|q|^{2/3}} = \frac{p}{\sqrt[3]{|q|^2}} = \frac{-3\Delta_0}{\sqrt[3]{|\Delta_1|^2}} \quad \text{is dimensionless}$$

$$u_+ = -\text{sgn } q \left( \frac{1}{2} + \sqrt{\frac{1}{4} + \left(\frac{s}{3}\right)^3} \right)$$

$$r = \sqrt[3]{u_+} - \frac{s/3}{\sqrt[3]{u_+}} = -\text{sgn } q \left( \sqrt[3]{\frac{1}{2} + \sqrt{\frac{1}{4} + \left(\frac{s}{3}\right)^3}} - \frac{s/3}{\sqrt[3]{\frac{1}{2} + \sqrt{\frac{1}{4} + \left(\frac{s}{3}\right)^3}}} \right)$$

$$= -\text{sgn } q \frac{\sqrt[3]{4}}{2} \left( \sqrt[3]{1 + \sqrt{1 + \zeta^3}} - \frac{\zeta}{\sqrt[3]{1 + \sqrt{1 + \zeta^3}}} \right) \quad \text{with } \zeta = \sqrt[3]{4} \frac{s}{3} : \text{ no cancelation for } \zeta < 0$$

$$= -\text{sgn } q \frac{\sqrt[3]{4}}{2} \frac{\left(\sqrt[3]{1 + \sqrt{1 + \zeta^3}} - \sqrt{\zeta}\right) \left(\sqrt[3]{1 + \sqrt{1 + \zeta^3}} + \sqrt{\zeta}\right)}{\sqrt[3]{1 + \sqrt{1 + \zeta^3}}} \underset{\zeta \rightarrow +\infty}{\sim} -\text{sgn } q \frac{\sqrt[3]{4}}{3\zeta} = \frac{-\text{sgn } q}{s}$$

$$\sqrt{1 + \zeta^3} = \sqrt{(1 + \zeta)(1 - \zeta + \zeta^2)}$$

when  $\Delta < 0$  (i.e.  $4s^3 + 27 > 0$ ), i.e.  $\left(\frac{s}{3}\right)^3 > -\frac{1}{4}$ , one gets a single real root:



$$\kappa = \frac{\zeta^4}{2\sqrt{\zeta^3+1}(\sqrt{\zeta^3+1}+1)^{\frac{5}{3}} - 2\zeta\sqrt{\zeta^3+1} - 2\zeta^4 - 2\zeta} - \frac{\zeta}{(\sqrt{\zeta^3+1}+1)^{\frac{2}{3}} - \zeta}$$

$$+ \frac{\zeta^3}{-2\zeta\sqrt{\zeta^3+1}(\sqrt{\zeta^3+1}+1)^{\frac{1}{3}} + 2\sqrt{\zeta^3+1} + 2\zeta^3 + 2}$$

$\zeta \xrightarrow{\sim} +\infty$   $-1$

$$\kappa = \frac{\zeta^3}{2\sqrt{\zeta}\sqrt{\zeta^3+1}(\sqrt{\zeta^3+1}+1)^{\frac{2}{3}} + 2\sqrt{\zeta^3+1} + 2\zeta^3 + 2}$$

$$+ \frac{\zeta^3}{-2\sqrt{\zeta}\sqrt{\zeta^3+1}(\sqrt{\zeta^3+1}+1)^{\frac{2}{3}} + 2\sqrt{\zeta^3+1} + 2\zeta^3 + 2}$$

$$+ \frac{\zeta}{2\sqrt{\zeta}(\sqrt{\zeta^3+1}+1)^{\frac{1}{3}} + 2\zeta} - \frac{\zeta}{2\sqrt{\zeta}(\sqrt{\zeta^3+1}+1)^{\frac{1}{3}} - 2\zeta} - \frac{\zeta^3}{2\sqrt{\zeta^3+1} + 2\zeta^3 + 2}$$



# Testing precision with BASIC'85

```
for (unsigned nbTot = NBITERMIN; nbTot < NBITERMAX; nbTot++) {  
    float x = X0;  
    for (unsigned nbIter = 0; nbIter < nbTot; nbIter++) x = sqrt (x);  
    float bottomRadix = x;  
    for (unsigned nbIter = 0; nbIter < nbTot; nbIter++) x = x * x;  
    printf ("%d %f %f (%+e) %f (%+e)\n", nbTot, X0, x, x-X0, bottomRadix, bottomRadix-1.0);  
}
```

iter	X0	x	x - X0	btmRdx	btmRdx - 1
10	2.000000	1.999958	(-4.184246e-05)	1.000677	(+6.771088e-04)
11	2.000000	2.000196	(+1.962185e-04)	1.000339	(+3.385544e-04)
12	2.000000	2.000196	(+1.962185e-04)	1.000169	(+1.692772e-04)
13	2.000000	2.000196	(+1.962185e-04)	1.000085	(+8.463860e-05)
14	2.000000	2.000196	(+1.962185e-04)	1.000042	(+4.231930e-05)
15	2.000000	1.996286	(-3.713965e-03)	1.000021	(+2.110004e-05)
16	2.000000	1.988545	(-1.145530e-02)	1.000010	(+1.049042e-05)
17	2.000000	1.988545	(-1.145530e-02)	1.000005	(+5.245209e-06)
18	2.000000	1.988545	(-1.145530e-02)	1.000003	(+2.622604e-06)
19	2.000000	1.988545	(-1.145530e-02)	1.000001	(+1.311302e-06)
20	2.000000	1.868132	(-1.318680e-01)	1.000001	(+5.960464e-07)
21	2.000000	1.648514	(-3.514862e-01)	1.000000	(+2.384186e-07)
22	2.000000	1.648514	(-3.514862e-01)	1.000000	(+1.192093e-07)
23	2.000000	1.000000	(-1.000000e+00)	1.000000	(+0.000000e+00)



What is the relative sensitivity of a function with respect to input argument fluctuation?  
⇒ *condition number* or absolute value of *elasticity*

$$\kappa(x) = \frac{\left| \frac{f(x_a) - f(x)}{f(x)} \right|}{\left| \frac{x_a - x}{x} \right|} = \frac{\left| \frac{f(x_a) - f(x)}{(x_a - x)} \right|}{\left| \frac{f(x)}{x} \right|} \sim \left| \frac{x f'(x)}{f(x)} \right| = \left| \frac{d(\ln |f(x)|)}{d \ln |x|} \right| \quad (2)$$

$\kappa$  is dimensionless, a pure number (*doubly logarithmic derivative*)



- Power law  $x \rightarrow C \times x^n$  (with  $C$  and  $n$  real constants) are the functions with uniform condition number:  $\forall x, \kappa(x) = n$ .
- $\log_2 \kappa$ : number of accuracy bits lost *in the best case, with correct rounding*
- $f : x \rightarrow x^2 \Rightarrow \kappa = \frac{2x \cdot x}{x^2} = 2$ : no singularity, relative error doubles on each iteration  
 $f : x \rightarrow \sqrt{x} \Rightarrow \kappa = \frac{1}{2}$ : no singularity, relative error is halved on each iteration (but can't really get below  $\frac{1}{2}$  ulp)

Very few uncertainty caused by iterations of  $\sqrt{\quad}$ , still the last half ulp is responsible for losing 100% of accuracy

then iterations of  $x \rightarrow x^2$  amplify this generally negligible error to a macroscopic one.



# Elasticity and condition number: more examples

$$\kappa_{f \circ g} = \kappa_f \times \kappa_g$$

$$\kappa_{f \times g} = \kappa_f + \kappa_g$$

$$\kappa_{f^n} = n\kappa_f$$

- $f : x \rightarrow x - c \Rightarrow \kappa = \frac{x}{x-c}$ : singularity  $x = c$  (*catastrophic cancellation*)
- $f : x \rightarrow \ln x \Rightarrow \kappa(x) = \frac{1}{\ln x}$ : singularity  $x = 1$ ,  $f(x = 1 + h) = \ln(1 + h)$

$$\kappa(h) = \frac{h}{(1+h)\ln(1+h)} \underset{h \rightarrow 0}{\sim} \frac{1}{(1+h)} \quad \text{hence the importance of } \log_2 p$$

- $f : x \rightarrow \exp x - 1 \Rightarrow \kappa(x) = \frac{x \exp x}{\exp x - 1}$ : indeterminate form  $x = 0$ ,  $\kappa(h) \underset{h \rightarrow 0}{\sim} 1$   
hence the importance of  $\exp m 1$

- $f : x \rightarrow \cos x - 1 \Rightarrow \kappa(x) = \frac{-x \sin x}{\cos x - 1}$ : indeterminate form  $x = 0$ ,  $\kappa(h) = \frac{h \cos \frac{h}{2}}{\sin \frac{h}{2}} \underset{h \rightarrow 0}{\sim} 2$   
hence the importance of trigonometry

To bypass cleanly this «tower of roots» problem (even in single precision), one needs to change the naive approach and use  $\log_2 p$  and  $\exp m 1 \Rightarrow$  *exercise: do it!*



$$P : x \mapsto a_n x^n + \dots + a_1 x + a_0$$

$$x \mapsto a_n \prod_{i=1}^n (x - x_i)$$

$$xP' : x \mapsto na_n x^n + \dots + a_1 x$$

$$\mapsto a_n x \sum_{j=1}^n \prod_{i \neq j} (x - x_i)$$

$$\kappa = xP'/P : x \mapsto a_n \sum_{j=1}^n \frac{x}{(x - x_j)}$$

Condition number has a pole around each root



Matrix condition number

$$\kappa(A) = \|A^{-1}\| \|A\| \geq \|A^{-1}A\| = 1$$

maximum ratio in relative error on  $x$  to relative error on  $b$  in  $Ax = b$

condition number for multi-variables functions

$$\kappa = \frac{\|J(x)\|}{\|f(x)\|/\|x\|} = \frac{\|x\| \|J(x)\|}{\|f(x)\|} = \frac{\|x_k\| \left\| \frac{\partial f_i}{\partial x_j} \right\|}{\|f(x)\|}$$



## Elasticity and condition number: dot product

$$f_y : x \mapsto y^T \cdot x$$

$$J : x \mapsto y^T$$

$$\kappa : x \mapsto \frac{\|J(x)\|}{|f(x)|/\|x\|} = \frac{\|x\| \|J(x)\|}{|f(x)|} = \frac{\|x\| \|y^T\|}{|y^T \cdot x|} = \frac{1}{\cos \angle(x, y)} > 1$$

Condition number has a pole around orthogonal vectors

<https://math.stackexchange.com/questions/3147927/condition-number-of-dot-product-of-vectors>

$$f : x \mapsto \sum_i^N x_i = y^T \cdot x \quad \text{with } y = (1 \dots 1)$$

$$\kappa : x \mapsto \frac{\|x\| \|y^T\|}{|y^T \cdot x|} = \frac{\|x\| \sqrt{N}}{|\sum_i x_i|} = \frac{1}{\cos \angle(x, y)} > 1$$

Condition number has a pole around collapsing sum

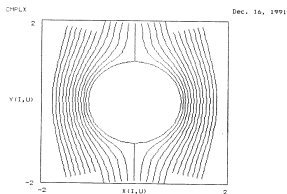


Figure 1 : Eluding Flow Past the Unit Disk

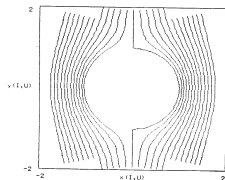


Figure 2 : Eluding Flow Past the Unit Disk, Almost

7

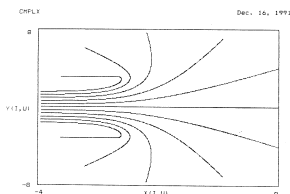


Figure 3 : Borda's Mouthpiece

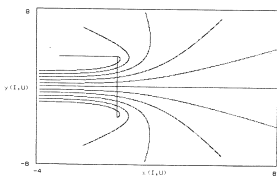


Figure 4 : Borda's Mouthpiece, Almost

8

**Eluding Flow past a Disk:**  $f: Z \mapsto (Z - 1/Z)/2$  and  $g: W \mapsto W - i\sqrt{iW - 1}\sqrt{iW + 1}$

Do not "simplify"  $g(W)$  to  $W - i\sqrt{-W^2 - 1}$  nor to  $W - \sqrt{W^2 + 1}$  since they behave differently. Though  $\forall W, f(g(W)) = W, \forall |Z| > 1, g(f(Z)) = Z$  only, and some  $|Z| = 1$ ; otherwise  $g(f(Z)) = -1/Z$ . Deducing where these identities hold is tricky.

**Borda's Mouthpiece:**  $W \mapsto 1 + W^2 + W\sqrt{W^2 + 1} + \ln(W^2 + W\sqrt{W^2 + 1})$

as  $W$  runs on radial straight lines through 0 in the right half-plane, including the imaginary axis.



<https://dlmf.nist.gov/>    <https://dlmf.nist.gov/4>    <https://sourceforge.net/projects/cmplx>

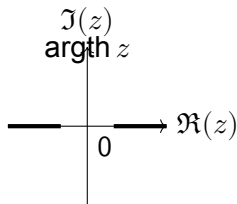
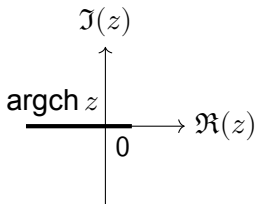
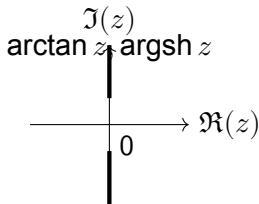
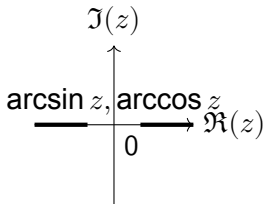
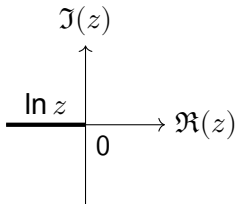
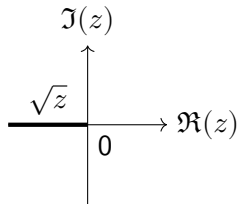
$$\arcsin z = -i \ln \left( iz + \sqrt{1 - z^2} \right)$$

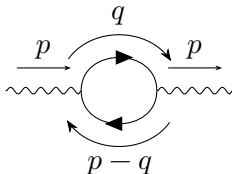
$$\arccos z = -i \ln \left( z + i\sqrt{1 - z^2} \right)$$

$$\operatorname{argsh} z = \ln \left( z + \sqrt{1 + z^2} \right)$$

$$\operatorname{argch} z = 2 \ln \left( \sqrt{(z+1)/2} + \sqrt{(z-1)/2} \right)$$

$$\operatorname{argth} z = \ln \left( (1+z)\sqrt{1/(1-z^2)} \right)$$





$$\begin{aligned}
 B_0(p, m_1, m_2) &= 16\pi^2 Q^{4-n} \int \frac{d^n q}{i(2\pi)^n} \frac{1}{\left[ q^2 - m_1^2 + i\varepsilon \right] \left[ (q-p)^2 - m_2^2 + i\varepsilon \right]} \\
 &= \frac{1}{\varepsilon} - \int_0^1 dx \ln \frac{(1-x)m_1^2 + xm_2^2 - x(1-x)p^2 - i\varepsilon}{Q^2} \\
 &= \frac{1}{\varepsilon} - \ln \left( \frac{p^2}{Q^2} \right) - f_B(x_+) - f_B(x_-)
 \end{aligned}$$

$$s = p^2 - m_2^2 + m_1^2, x_{\pm} = \frac{s \pm \sqrt{s^2 - 4p^2(m_1^2 - i\varepsilon)}}{2p^2}, f_B(x) = \ln(1-x) - x \ln(1-x^{-1}) - 1$$

- ⇒ the (microscopic) difference of  $\varepsilon$  induces a (macroscopic) difference of  $2\pi$  on the imaginary part
- ⇒ the analytic functions<sup>1</sup> of complex analysis are sharply discontinuous at the crossing of their *branch cut*

<sup>1</sup> a *minima* indefinitely continuously differentiable

# Discrete Stochastic Arithmetic (DSA) [Vignes'04]

Classic arithmetic

$$A \oplus B \rightarrow R$$

$R = 3.14237654356891$

DSA

Random  
rounding

$$A_1 \oplus B_1 \rightarrow R_1$$

$$A_2 \oplus B_2 \rightarrow R_2$$

$$A_3 \oplus B_3 \rightarrow R_3$$

$R_1 = \mathbf{3.141354786390989}$

$R_2 = \mathbf{3.143689456834534}$

$R_3 = \mathbf{3.142579087356598}$

- each operation executed 3 times with a random rounding mode
- number of correct digits in the results estimated using Student's test with the confidence level 95%
- operations executed synchronously
  - ⇒ detection of numerical instabilities  
Ex: `if (A>B)` with A-B numerical noise
  - ⇒ optimization of stopping criteria



- implements stochastic arithmetic for C/C++ or Fortran codes
- few code rewriting
- all operators and mathematical functions overloaded
- support for MPI, OpenMP, GPU, vectorised codes
- supports emulated ou native half precision
- in one CADNA execution: accuracy of any result, complete list of numerical instabilities

## CADNA cost

- memory: 4
- run time  $\approx 10$



# Executing CADNA

Before modifying the precisions used, we want to explore the current accuracy.



# Executing CADNA

Before modifying the precisions used, we want to explore the current accuracy.

To execute CADNA, we essentially change the types.



Before modifying the precisions used, we want to explore the current accuracy.

To execute CADNA, we essentially change the types.

This execution exposed multiple numerical instabilities that hide potential massive loss of accuracy.

```
-----  
CADNA_C 3.1.11 software
```

```
CRITICAL WARNING: the self-validation detects major problem(s).  
The results are NOT guaranteed.
```

```
There are 538393974 numerical instabilities  
10409 UNSTABLE DIVISION(S)  
40122229 UNSTABLE MULTIPLICATION(S)  
267297 UNSTABLE BRANCHING(S)  
448561143 UNSTABLE INTRINSIC FUNCTION(S)  
266 UNSTABLE MATHEMATICAL FUNCTION(S)  
49432630 LOSS(ES) OF ACCURACY DUE TO CANCELLATION(S)  
-----
```

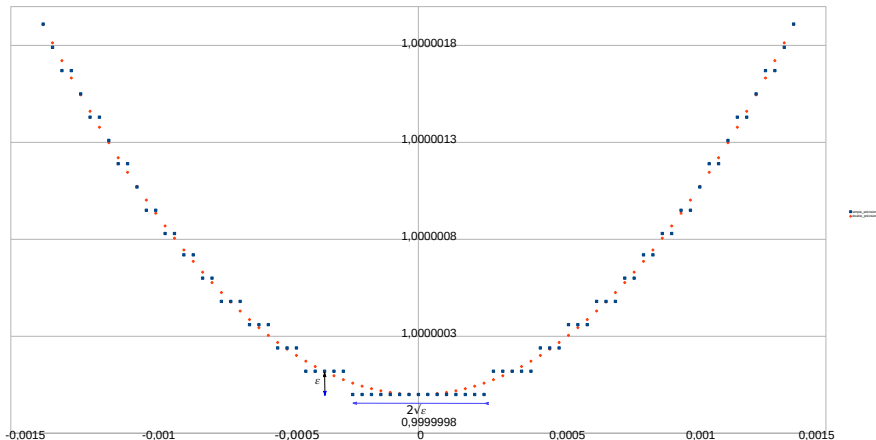


# Minimisation

Numerical evaluation of derivatives / gradients / Jacobian / Hessian

$$x \mapsto 1 + (x - 1)^2 \Rightarrow f(x = x_0 + h) = f(x_0) + \underbrace{h \cdot \frac{\partial f}{\partial x}}_{=0 \text{ at extremum}} + {}^t h \cdot \frac{\partial^2 f}{\partial x^2} \cdot h + o(h^2) \dots \text{TAYLOR}$$

Forme Quadratique





# Neural Network

## Exploration of Machine learning for Polynomial Root Finding

Vitaliy Gyrya, Mikhail Shashkov, Alexei Skurikhin  
(T-5) Applied Mathematics & Plasma Physics, (XCP-4) Methods & Algorithms, (ISR-3) Space Data Science & Systems



Machine Learning for Computational Fluid and Solid Dynamics  
February 19-21, 2019



### Motivation

We are interested in application of Machine Learning (ML) for improving numerical methods for solving partial differential equations (PDEs). One example of such an improvement is the optimization of the parameters of artificial viscosity for Lagrangian and arbitrary-Lagrangian-Eulerian methods. Another example is solving the Riemann problem, which is at the core of many numerical methods for computational gas and solid dynamics. To build confidence in ML methods and understand their strengths and weaknesses we decided to start by applying ML to solve simple quadratic equations of one variable.

### Problem

Consider a quadratic equation,  $ax^2 + bx + c = 0$ , whose roots are  $r_L$  and  $r_R$ . We would like to learn the function

$$(a, b, c) \rightarrow (r_L, r_R)$$

without relying on our knowledge of the underlying processes. Instead we will consider a number of observations observations (training set)

$$(\hat{a}^i, \hat{b}^i, \hat{c}^i) \rightarrow (r_L^i, r_R^i), \quad i = 1, \dots, N.$$

From which we will try to predict

$$(\hat{a}^j, \hat{b}^j, \hat{c}^j) \rightarrow (\tilde{r}_L^j, \tilde{r}_R^j) \approx (r_L^j, r_R^j), \quad j = N + 1, \dots, N + K.$$

The goal is to minimize

$$\text{COST} = \sum_i (r_L^i - \tilde{r}_L^i)^2 + \sum_j (r_R^j - \tilde{r}_R^j)^2.$$

### Challenges

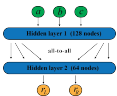
The quadratic equation was selected as a proxy for the following reasons that are relevant to many complex practical problems:

- There are several branches in the solution: if  $a = 0$ , the quadratic equation becomes a linear equation, with one root – this is a qualitative change from one regime to a different one; depending on the discriminant the number of roots as well as the nature of the roots changes (real vs. complex).
- Finding solution involves different arithmetic operations some of which can be difficult to model by machine learning techniques. For example, division and square root are a challenge for neural networks to represent as activation functions.
- Probably, the most significant challenge is that for a small range of input parameters for which output values are increasingly large.

### Feed-forward Neural Network

#### NN Architecture:

Input Layer: 3 nodes  
 Hidden Layer 1: 128 ReLU  
 Hidden Layer 2: 64 ReLU  
 Output Layer: 2 Linear  
 Connectivity: full.



#### NN Training:

Batch size: 200  
 Training epochs: under 500  
 Optimizer: Adam (<https://arxiv.org/abs/1412.6980v8>)

### Gauss Process Regression (GPR)

Probabilistic Bayesian generalization of linear regression approach.

- Built in model of uncertainty estimator.
- Need to specify a covariance kernel.

Our choice of kernel:  
 ConstantKernel() +  
 Matern(length\_scale = 2, nu = 3/2) +  
 WhiteKernel(noise\_level = 1)



### Test & training sets

We considered a number of distributions for the coefficients  $(a, b, c)$ . In all these cases we assumed that

$$a \in [r, 1], \quad b \in [-1, -1], \quad c \in [-1, -1], \quad \epsilon = 1/20$$

and the roots  $(r_L, r_R)$  are real, i.e.  $D = b^2 - 4ac \geq 0$ .

We considered the following distributions for  $(a, b, c)$

- Uniform random distribution.
  - Regular distribution for  $(a, b, c)$ , i.e. distribution on a grid.
  - Regular distribution for  $(1/a, b, c)$ , i.e. distribution on a grid.
- The sizes of the training and test sets were approximately equal and were on the order of 40K to 50K data points.

### GPR for large datasets

- GPR performance degrades quickly (scaling  $\sim N^3$ ).
- Depending on the machine the threshold of tractable training sets was between 5K and 50K sample points.
- More advanced techniques are needed for larger data sets.
- Ensembles of smaller GPR could be used.

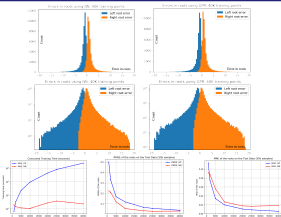
### Adaptive sampling with GPR

#### Adaptation procedure:

- Consider the pool of uniformly distributed parameters  $(\hat{a}^i, \hat{b}^i, \hat{c}^i)$ .
- Select an initial training set of points (50) at random. Generated GPR based on these points.
- For the given GPR consider the "uncertainty"  $\sigma$  at all of the sample points. Find the triples  $(\hat{a}^i, \hat{b}^i, \hat{c}^i)$  with the largest uncertainty and add them to the training set.
- Generate a new GPR for the updated training set.
- Repeat steps 3-4 until stopping criteria is satisfied, e.g. training set reached predefined size.



### Results



### Conclusions

- For small data sets (2K points) GPR is more accurate
- GPR can utilize adaptive sampling
- GPR does not scale well to larger data sets (~2K points).
- NN scales well for large data sets and has better accuracy over GPR (more than 5K points).





# Floating Point Types

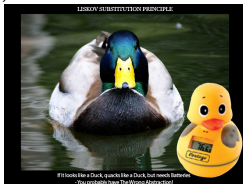
`float` and `double` are identified as simple or even primitive types, but they are much richer than it seems.

**Object** point of view: do these types fit into a hierarchy of classes?

⇒ Violation of the LISKOV's substitution principle (LSP)

if S subtypes T, what holds for T-objects holds for S-objects.

If S is a subtype of T, objects of type T in a program can be replaced by objects of type S without changing any of the desirable properties of that program (e.g. correct results)



A poorly encapsulated abstraction (*leaky*): we can measure the smallest positive non-zero float, the largest one, the machine epsilon, the base: we can access the implementation details



Not metrology: we do not seek “precision for precision’s sake”

The **functional** paradigm invites us to write computer function approaching mathematical functions, and we tend to focus on the aspect of **purity**.

But a mathematical function also seeks **totality** (being defined on the largest domain of definition):

*the function should be calculable for any argument for which it is defined.*

- removing non-jump and non-essential discontinuity:  $\Rightarrow \frac{\sin x}{x} \Big|_{x=0} = 1$  (naively  $\sin(0.0) / 0.0 = \text{NaN}$ )
- analytic continuation: factorial  $\Rightarrow \Gamma$ , or RIEMANN  $\zeta$  function
  - $\Rightarrow$  maximal extension of function domain
  - $\Rightarrow$  piecewise function definition, *casuistry*

Using IEEE-754 exceptional values, we can reach a “weak totality”:

- $\log(0.0) = -\text{Inf}$  (mathematically correct)
- $\log(-1.0) = \text{NaN}$  (mathematically correct? more precisely  $\text{NaN}$ )

Precision limitations lead to a gray zone in this kind of totality:

- $\text{expf}(88.72284) = +\text{Inf}$  (but mathematically it's  $2^{128} \Rightarrow \text{domainException}$ )
- $\text{expf}(-103.972084) = 0.0\text{f}$  (but mathematically it's just below  $2^{-150} \Rightarrow \text{domainException}$ )
- $\text{gammaf}(35.0401001) = +\text{Inf}$  (but mathematically it's  $2^{128}$ )

OK with `double`, but not with `float`.

Not all `Inf` have the same meaning, not all `NaN` have the same meaning, *cf* `null` in SQL



## « Why aiming for precision? »

⇒ Implicit **contract**: the fonction will

- 1 (if the argument is inside the mathematical domain of the mathematical function)
- 2 (if the type representation of the argument is inside the domain of the function that has a representable image in the return type)
- 3 return a result
- 4 this result is relevant(?)
- 5 (ideally the returned value is the representation of the image of the mathematical function applied on the represented argument)



# « Why aiming for precision? »

totality (mathematical) vs. representable totality

A representable solution resulting from representable arguments CAN go through a non-representable intermediate calculation. IEEE-754 exceptional values are not the value of the function, relative error of 100%, as in catastrophic cancelation.

## least surprise principle

- we agree to compute erroneous results, because we know that we cannot compute exact results: exact results are rarely (= almost never) representable:  $\pi$ ,  $e$ ,  $\sqrt{2}$ ,  $1/3$ ,  $1/5$  in base 2...
- On the other hand, we don't want things to be very wrong: mathematical result 2 but the function returns NaN

If the calculation is badly carried out, we can end up with

- infinite roots, where they exist and can be represented
- to an absence of roots, where they exist and are representable
- to a presence of roots, where they do not exist

**a difference of degree generates a difference of nature** (catastrophe theory, bifurcation, chaos)

The relative size of the danger zone in the parameter space will be much larger in low precision.

Annex for a less costly nondimensionalization:

« *You Could Learn a Lot from a Quadratic* » doi:10.1145/609742.609746, shows how to nondimensionalize with binary, much less costly in time and accuracy than divisions (and roots) in physicist nondimensionalization. Easy when knowing IEEE-754

API.



# « precision? » a take-away

$$\text{PRECISE NUMBER} + \text{PRECISE NUMBER} = \text{SLIGHTLY LESS PRECISE NUMBER}$$

$$\text{PRECISE NUMBER} \times \text{PRECISE NUMBER} = \text{SLIGHTLY LESS PRECISE NUMBER}$$

$$\text{PRECISE NUMBER} + \text{GARBAGE} = \text{GARBAGE}$$

$$\text{PRECISE NUMBER} \times \text{GARBAGE} = \text{GARBAGE}$$

$$\sqrt{\text{GARBAGE}} = \text{LESS BAD GARBAGE}$$

$$(\text{GARBAGE})^2 = \text{WORSE GARBAGE}$$

$$\frac{1}{N} \sum (\text{N PIECES OF STATISTICALLY INDEPENDENT GARBAGE}) = \text{BETTER GARBAGE}$$

$$\left( \frac{\text{PRECISE NUMBER}}{\text{GARBAGE}} \right) = \text{MUCH WORSE GARBAGE}$$

$$\text{GARBAGE} - \text{GARBAGE} = \text{MUCH WORSE GARBAGE}$$

$$\frac{\text{PRECISE NUMBER}}{\text{GARBAGE} - \text{GARBAGE}} = \text{MUCH WORSE GARBAGE, POSSIBLE DIVISION BY ZERO}$$

$$\text{GARBAGE} \times 0 = \text{PRECISE NUMBER}$$

<https://xkcd.com/2295/>