



MAX PLANCK INSTITUTE
FOR DYNAMICS OF COMPLEX
TECHNICAL SYSTEMS
MAGDEBURG



COMPUTATIONAL METHODS IN
SYSTEMS AND CONTROL THEORY

[20 YEARS]
[1998-2018]

An Introduction to Posit Arithmetic

On “Beating Floating Point at its Own Game:
Posit Arithmetics” by J.L. Gustafson and I. Yonemoto

C. Himpe

2018-12-05

CSC Reading Group

- Exascale Challenges: Memory & Power.
- Let's look at the fundamental layer: Numbers.
- Floats are completely bonkers!
- This is **not** a computer science, ...
- ... but a **scientific computing** (\rightarrow math) issue!



Coarse Dynamic range of 10^{-324} to 10^{308} (DP).

- Number of neurons in the brain: $\sim 10^{14}$
- Age of the universe [ns]: $\sim 10^{25}$
- Planck Scale: $\sim 10^{-35}$
- Atoms in the universe: $\sim 10^{82}$

Wasteful Many different bit patterns for NaN (Not-a-Number).

- Half Precision (16-bit): $\sim 10^3$
- Single Precision (32-bit): $\sim 10^7$
- Double Precision (64-bit): $\sim 10^{15}$

Ambivalent Why is there a -0 ? Should $\sqrt{-0} = -0$?

Weird $\frac{1}{0} = \infty$ and $\frac{1}{-0} = -\infty$, but $0 = -0$, hence $\infty = -\infty$?

Inconsistent $(NaN == NaN) = 0$, but $(0 == -0) = 1$? $bool(NaN) = 1$!

Arbitrary IEEE754 is more a guideline than a standard.

LOL "Extended Precision" (80-bit)

Unum¹ History:

2015 Unum I (Universal NUMber): Superset of floats

2016 Unum II: Float incompatible interval arithmetic, variable size

2017 Unum III (aka **Posit**): Float related and fixed size

Why Posits:

- Unum I & Unum II have varying (bit) length.
- Fix the worst flaws of floats and improve accuracy.
- A drop-in replacement for floats.
- A “modern” number format.
- No marketing or engineering hacks.

¹J.L. Gustafson. **The End of Error: Unum Computing**. Chapman and Hall/CRC, 2015

posit (verb) *Something that is posited.* (Wiktionary)

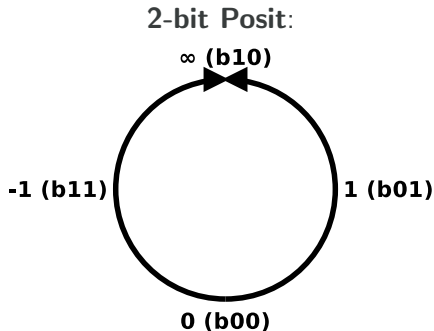
posit (noun) *To suggest something as a basic fact or principle from which a further idea is formed or developed.* (Cambridge Dictionary)

About Non-interval number format with encoding similar to floats.

Note POsit OPerations (POPs) instead of FLOPs.

Features

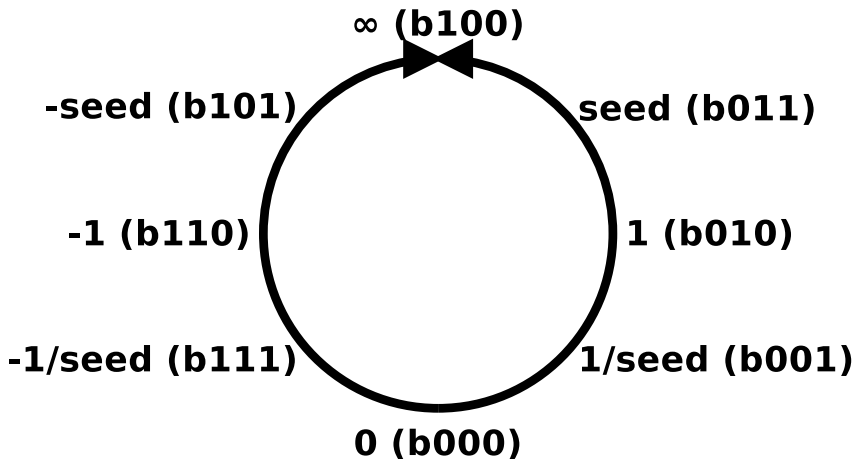
- More accurate.
- Less wasteful.
- More operations per Watt (\rightarrow Power Aware COmputing).
- Simple hardware implementation.
- Every bit pattern has a unique meaning.
- Only two special “numbers”, no NaNs.



Remember Complex Analysis:

- Stereographic projection of the complex plane.
- Projective completion: Point at infinity².
- Posits are a discrete representation of the real projective line.

²This also "fixes" division by zero (<https://www.1dividedby0.com>).



Float:

000000000000000000000000

Formula:

$$x = (-1)^{\text{sign}} \cdot 2^{\text{exponent}} \cdot (1 + \text{fraction})$$

Sign positive if zero, negative if one (one bit)

Exponent integer expressing power of two (fixed length)

Mantissa fractional remainder (fixed length)

(e_s Exponent Size) Posit:

00...000...

Formula:

$$x = (-1)^{\text{sign}} \cdot \left(2^{2^{e_s}}\right)^{\text{regime}} \cdot 2^{\text{exponent}} \cdot (1 + \text{fraction})$$

e_s maximum number of exponent bits,

Seed $2^{2^{e_s}}$

Sign positive if zero, negative if one (one bit),

Regime integer expressing power of power of two (varying length),

Exponent integer expressing power of two (varying length),

Mantissa (linear) fractional remainder (varying length),

Size sizeof(Regime) + sizeof(Exponent) + size(Fraction) = constant

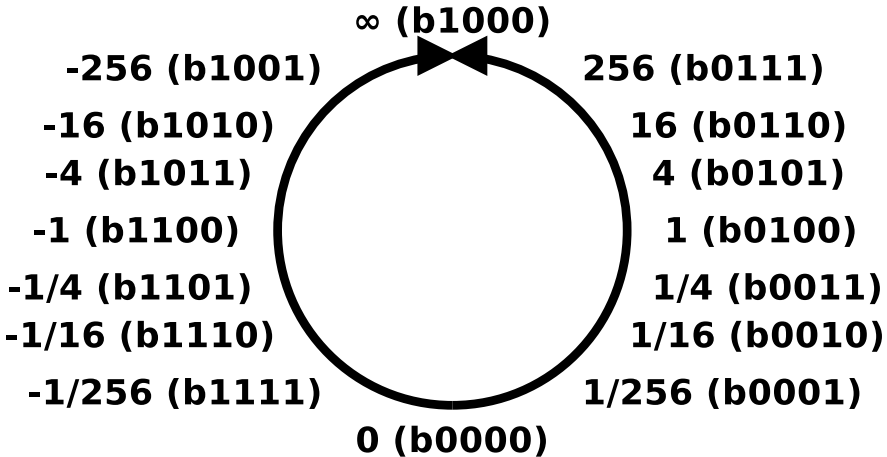
Rationale:

- Determines “larger order of magnitude”.
- Golomb-Rice type encoding.
- Algorithm: Count leading zeros. Available in:
 - Assembly (x86-32 `bsf`, x86-64 `lzcnt`, arm `clz`)
 - GNU (`__builtin_clz`)
 - Javascript (`Math.clz`)

Example (Four Regime Bits³):

0000 : -4 → regime = $(2^{2^{e_s}})^{-4}$	10ZZ : 0 → regime = $(2^{2^{e_s}})^0$
0001 : -3 → regime = $(2^{2^{e_s}})^{-3}$	110Z : 1 → regime = $(2^{2^{e_s}})^1$
001Z : -2 → regime = $(2^{2^{e_s}})^{-2}$	1110 : 2 → regime = $(2^{2^{e_s}})^2$
01ZZ : -1 → regime = $(2^{2^{e_s}})^{-1}$	1111 : 3 → regime = $(2^{2^{e_s}})^3$

³Remember: This is just an example, the regime field has variable length!





Posit (8-bit, 2-bit Exponent)

- 11000000 : $-1^1 \cdot \left(2^{2^2}\right)^0 \cdot 2^0 \cdot (1 + 0) = -1$
- 01001101 : $-1^0 \cdot \left(2^{2^2}\right)^0 \cdot 2^1 \cdot \left(1 + \frac{5}{8}\right) = 3.25 \approx \pi$
- 01001011 : $-1^0 \cdot \left(2^{2^2}\right)^0 \cdot 2^1 \cdot \left(1 + \frac{3}{8}\right) = 2.75 \approx e^1$

- Posits can be adapted in range.
- The (max) exponent size determines range (seed).
- Can be selected to be similar to floats.

Size (Bits)	Range (Float)	Range (Posit)	e_s (Posit)
16	$6 \cdot 10^{-8} \dots 7 \cdot 10^4$	$4 \cdot 10^{-9} \dots 3 \cdot 10^8$	1
32	$1 \cdot 10^{-45} \dots 3 \cdot 10^{38}$	$6 \cdot 10^{-73} \dots 2 \cdot 10^{72}$	3
64	$5 \cdot 10^{-324} \dots 2 \cdot 10^{308}$	$2 \cdot 10^{-299} \dots 4 \cdot 10^{298}$	4



Quantitative Comparison

See Figures 8, 9, 10, 11, 12, 14, 16 in:

J.L. Gustafson, I. Yonemoto. **Beating Floating Point at its Own Game: Posit Arithmetic**. Supercomputing Frontiers and Innovations: An International Journal Archive, 4(2): 71–86, 2017.



See Figure 13 in:

J.L. Gustafson, I. Yonemoto. **Beating Floating Point at its Own Game: Posit Arithmetic**. Supercomputing Frontiers and Innovations: An International Journal Archive, 4(2): 71–86, 2017.



Multiplication Closure

See Figure 15 in:

J.L. Gustafson, I. Yonemoto. **Beating Floating Point at its Own Game: Posit Arithmetic**. Supercomputing Frontiers and Innovations: An International Journal Archive, 4(2): 71–86, 2017.

Linear System Solver Package:

- Performance of dense linear system solve.
- Measures Flop (Pop) throughput.
- Determines the “Top 500” of super computers.

$$A \in U_{[0,1]}^{N \times N}$$

$$A \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = b$$

$$\rightarrow A x = b$$

32-bit Posits (3 regime bits) give a vector of exact ones!⁴

⁴(64-bit Floats don't)

- No more over- and underflow.
- Higher accuracy.
- Unique bit patterns.
- Sane dynamic range.
- Negation symmetry.
- Sorted like integers.
- Possibly halves memory demand.
- Posit design suggests fused operations for chained multiplication and additions⁵ in larger scratch registers, the so-called *quire*.
- Related formats are in development (Elias γ , δ , ω)

⁵See also FMA and MAD.

- No reciprocal symmetry⁶ (Compared to Unum II).
- No associativity.
- No hardware implementations (yet).
- Only signalling NaNs.
- No distinction between $-\infty$ and ∞ .

⁶Except for powers of two, or different fraction maps.

Soft Posits:

- **Octave:** <https://github.com/diegofgcoelho/positsoctave>
- **Matlab:** <https://bitbucket.org/kvasnica/munum>
- **Python:** <https://github.com/xman/numpy-posit>
- **Julia:** <https://github.com/Etaphase/FastSigmoids.jl>
- **C++:** <https://github.com/stillwater-sc/universal>
- **C:** <https://github.com/libcgb/bfp>
- and more⁷ at: https://posithub.org/khub_community

⁷FYI: I am dabbling in soft posits, too.



Summary

- less memory (use 32-bit posit instead of 64-bit float)
- more accurate (can be shown quantitatively and qualitatively)
- very elegant (projective mapping of reals)

More info: <https://posithub.org>

- D. Goldberg. **What Every Computer Scientist Should Know About Floating-Point Arithmetic**. Computing Surveys, 32(1): 6–48, 1991.
- J.L. Gustafson, I. Yonemoto. **Beating Floating Point at its Own Game: Posit Arithmetic**. Supercomputing Frontiers and Innovations: An International Journal Archive, 4(2): 71–86, 2017.
- J.L. Gustafson. **Posit Arithmetic**. Posithub, Knowledge Hub, 2017.
- J.L. Gustafson. **The End of Error: Unum Computing**. Chapman and Hall/CRC, 2015.
- Laslo Hunhold. **The Unum Number Format: Mathematical Foundations, Implementation and Comparison to IEEE 754 Floating-Point Numbers**. Bachelor Thesis, University of Cologne, 2016.
- P. Lindstrom, S. Lloyd, J. Hittinger. **Universal Coding of the Reals - Alternatives to IEEE Floating Point**. Proceedings of the Conference for Next Generation Arithmetic: 5, 2018.
- Wikipedia contributors. **Unum (number format)**. Wikipedia, The Free Encyclopedia, 2018.

Float and Posit Fraction Map: $\phi(f) = 1 + f$ (linear)

Reciprocal Symmetry Fraction Map⁸: $\phi(f) = \begin{cases} 1 + f & e > 0 \\ \frac{2}{2+f} & e < 0 \end{cases}$ (nonlinear)

Note:

- This means easy, but not necessarily faster, divisions.
- First bit (sign) is additive indicator (“left/right”)
- Second bit (exponent sign) is multiplicative indicator (“upper/lower”)
- One can construct other fraction maps for specific effects.

⁸P. Lindstrom, S. Lloyd, J. Hittinger. **Universal Coding of the Reals - Alternatives to IEEE Floating Point**. Proceedings of the Conference for Next Generation Arithmetic: 5, 2018.