Posit Arithmetic Lecture 1

Prof. John L. Gustafson

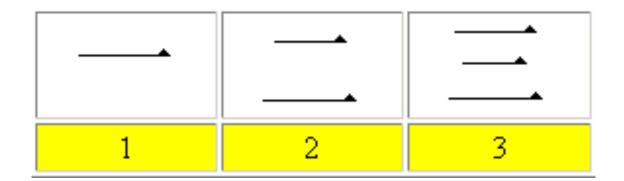
A*STAR and National University of Singapore





The Most Ancient of All Formats: *Unary*

Repeated symbols for counting numbers 1, 2, 3,...





1		6	11111
2	II	7	J##11
3	Ш	8	J## III
4	Ш	9	J##1111
5	##	10	ШШ

Regime Notation: Signed Unary, by Using Bits

```
000001
00001
0001
001
01
10
110
1110
11110
```

A run of *runlength* identical *r* bits is terminated by the opposite bit.

Represents an integer *k* which we will need later, so please remember this part!

```
k = -runlength if r is 0,

k = runlength - 1 if r is 1.
```

Numbers near 0 require very few bits.

Circuits exist for decoding regime bits

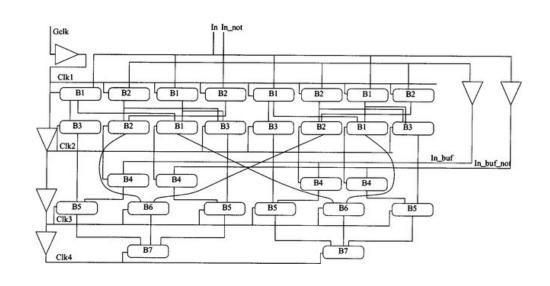
1.0005

-1.0003

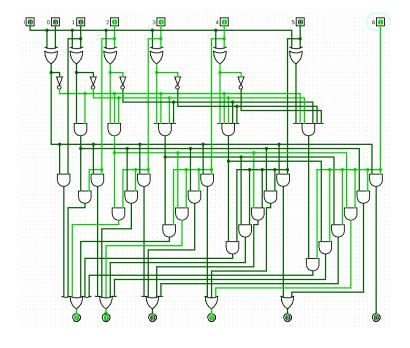
0.0002

⇒2×10⁻⁴

The CLZ (count leading zeros) instruction is already part of standard floating-point circuits.



CLZ logic (Wikipedia)



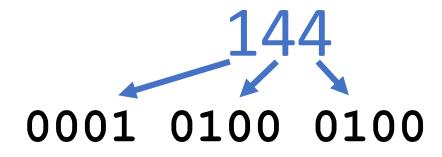
Regime shifter (I. Yonemoto)

Early computers used decimal internally!



ENIAC, 1946

This is an example of the mistake of imposing human tastes on hardware design.



Versus native binary, 8 bits:

$$\begin{array}{c}
 10010000 \\
 \hline
 2^7 + 2^4 = 144
 \end{array}$$

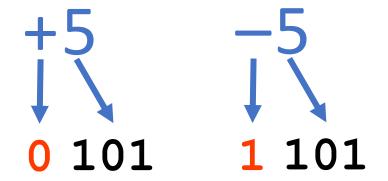
History's Next Mistake: A Separate "Sign Bit"



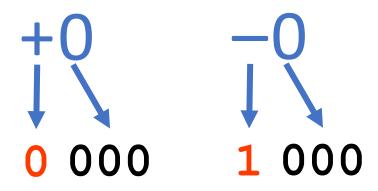


IBM 701, 1953

Negative integers were originally stored in "sign-magnitude" form, imitating the way humans write + and – before digit strings.



BAD idea. Why? Well, here's one reason:



Welcome to the joys of "negative zero."

Remember adding signed numbers in school?

To add nonzero signed integers *m* and *n*:

Are they the same sign or different sign?

If they are the same sign, add their magnitudes.

Apply that sign to the resulting sum, DONE.

Else if they have different signs,

Find out which magnitude is bigger.

If *m* has bigger magnitude,

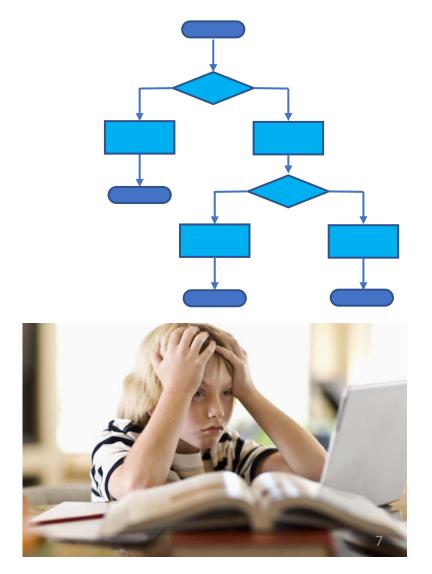
Subtract *n*'s magnitude from *m*'s magnitude.

Apply m's sign to the result. **DONE**.

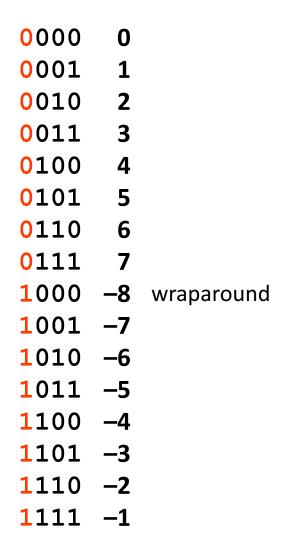
Else

Subtract m's magnitude from n's magnitude.

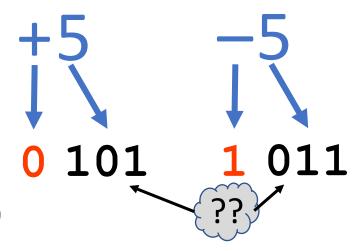
Apply *n*'s sign to the result. **DONE**.



Modern signed integers are 2's complement



Not as human friendly...



(flip the bits and add 1)

...but mathematically and computationally *far* better. Here's the addition algorithm for 2's complement:

To add signed integers *m* and *n*:

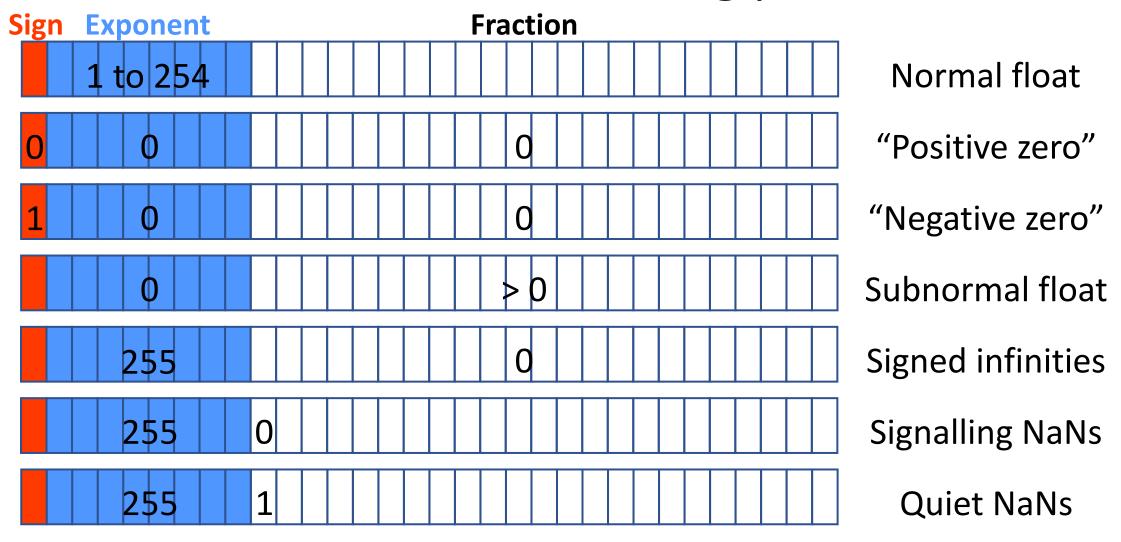
Add like unsigned integers. **DONE**

The fundamental idea of floating point format

Express values as $m \times 2^{j}$, where m and j are integers. Use a fixed set of bits for m and a fixed set for j.

- Notice we need both negative and positive m and j.
- Before the IEEE 754 Standard (1985), there were **many** different schemes for where the bits go in a word, and how to interpret the bits as signed integers.
- Larger dynamic range (more bits for *j*) means less accuracy (fewer bits for *m*). And vice versa.
- Note: in practice, the most common numbers have *j* near zero.

The IEEE 754 Standard greatly complicated the fundamental idea of a floating point format



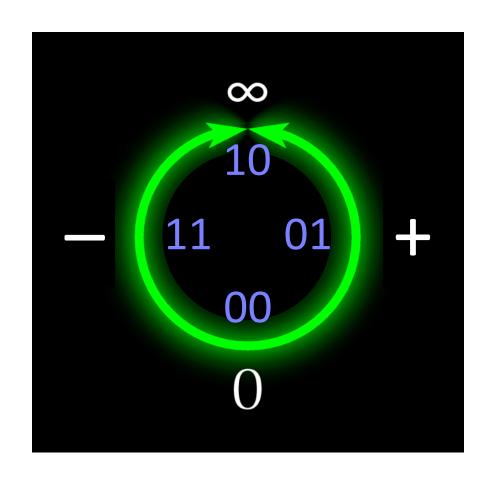
The fundamental idea of *posit* format

Express values as $m \times 2^j$, where m and j are integers. Use regime format for j, so a small j takes fewer bits.

- Creates 1-to-1 map of reals to integers, monotone
- The *j* is represented as the sum of two integers:
 - power-of-2 exponent e
 (unsigned integer ranging from 0 to 2^{es}-1)
 - a "regime exponent" k, the power of $useed = 2^{2^{es}}$.
 - Put more simply: the 2^{j} scaling is $2^{k2^{es}+e}$.

Some Prefer the Geometric Explanation

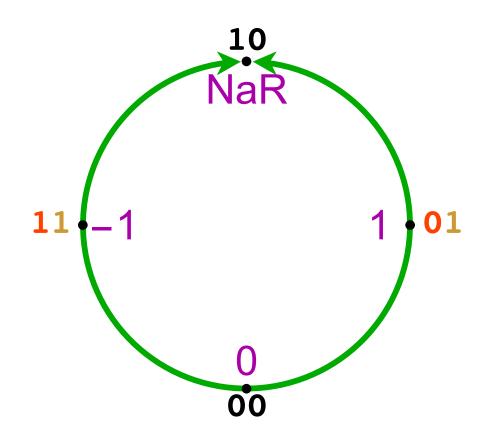
(This is how posits were invented)



Instead of the real number line, the *projective reals* put "the point at infinity" at the top of the circle... the first step to making an infinite line map to a finite-state computer format.

Big positives wrap to big negatives, just like (2's complement) signed integers.

Subtle change: Treat ∞ the way floats treat "NaN"

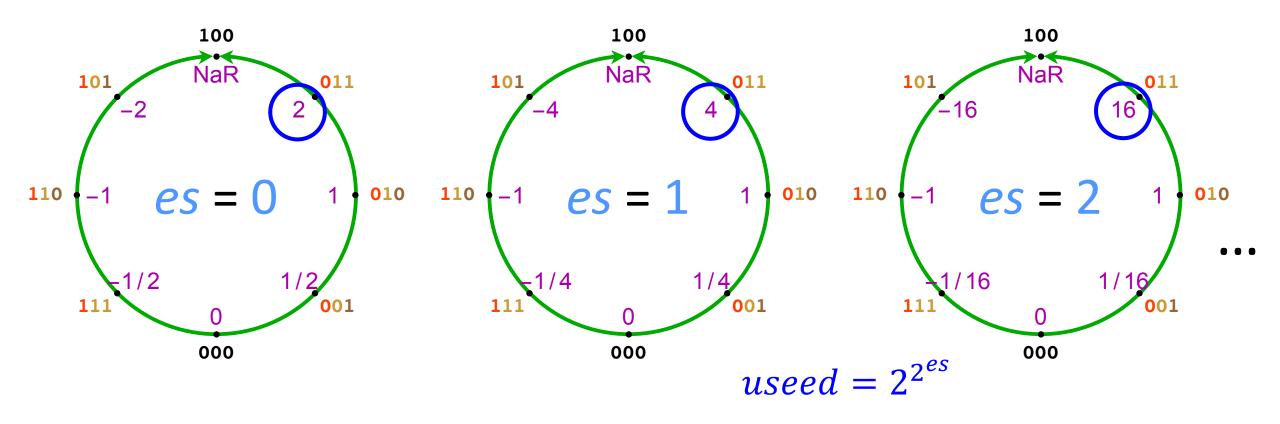


Not a Real (NaR) is the catch-all for $\sqrt{-1}$, 0/0, arcsin(3), etc.

Unlike $\pm \infty$, NaR always propagates, so 1/NaR = NaR, not 0.

Why "NaR" and not "NaN" (Not a Number"? Because imaginary and complex numbers are **numbers**! They're just not *real* numbers.

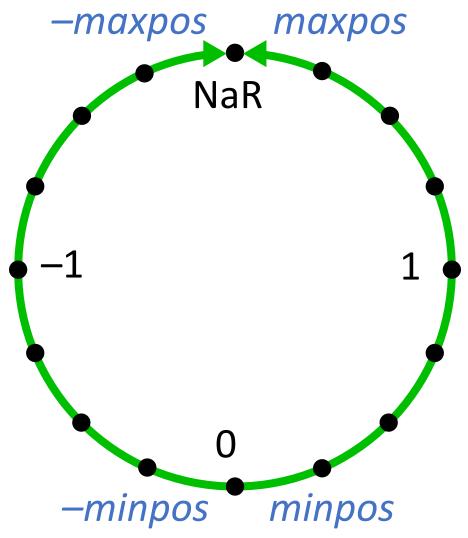
The es value controls the dynamic range.



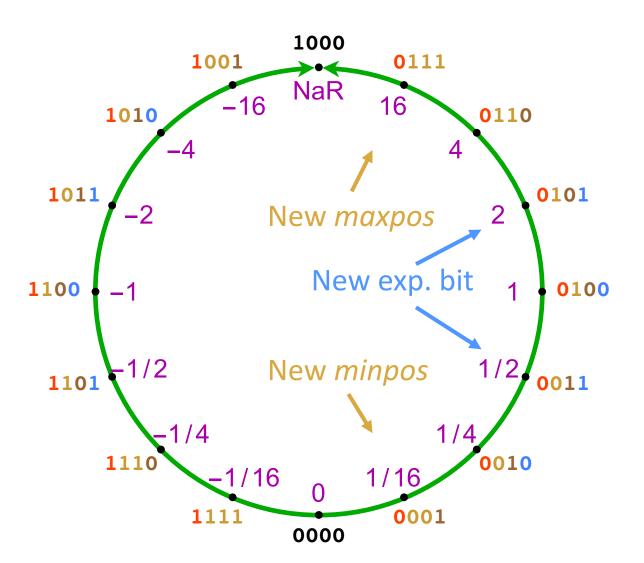
The northeast value is called the *useed*. Start with useed = 2, then do $useed \leftarrow useed^2$, es times.

Instead of overflow and underflow

- maxpos is the largest magnitude positive posit value.
- minpos is the smallest magnitude positive posit value.
- -maxpos is the largest magnitude negative posit value.
- -minpos is the smallest magnitude negative posit value.



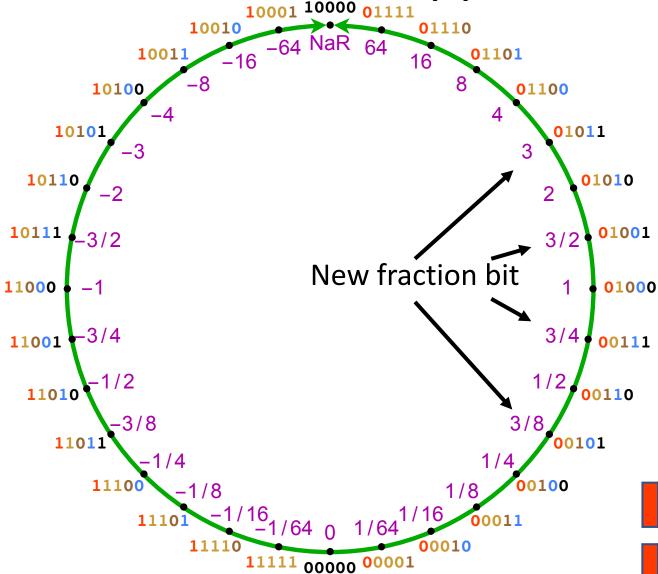
Pick es = 1 and append another bit.



If you can squeeze another integer power of 2 between two adjacent points of the previous, do so. The appended bit is an exponent bit.

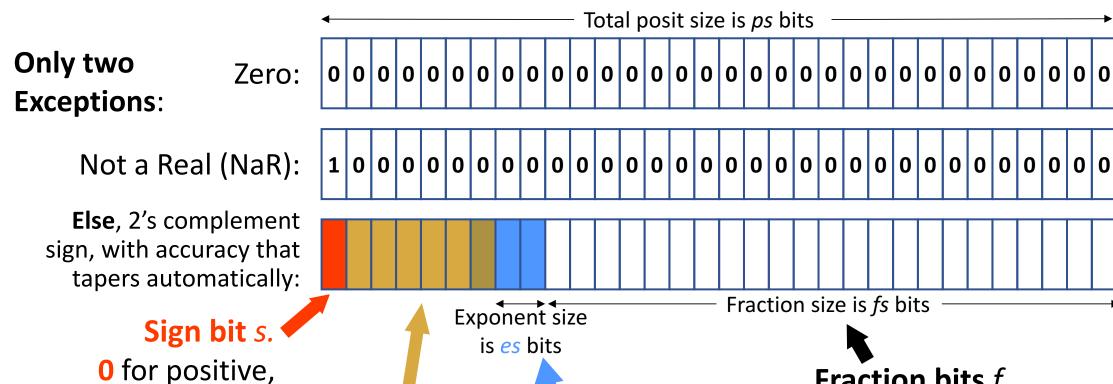
The new *maxpos* is the old *maxpos* times *useed*. The new *minpos* is the old *minpos* divided by *useed*. The appended bit is a regime bit.

Fraction bits appear



- If a new intermediate value cannot be an integer power of two, use the average. The appended bit is a fraction bit.
- Unusual property: Posits increase both dynamic range and decimals of accuracy by adding bits on the right.
- Notice that half of all posits use only two regime bits!

Making posits look more like IEEE floats:



"Regime" bits

Signed unary integer representing k, the power of $2^{2^{es}}$.

1 for negative.

Exponent bits *e*

unsigned integer. Size is es = 0, 1, 2,... for ps = 8, 16, 32,...

Fraction bits *f*

Fraction in positional notation. "Hidden bit" is **always** 1.

Human decoding or hardware decoding?



Easier to understand:

Check for 0 and NaR case first.

Else record s. If s = 1, find magnitude by 2's complement.

Find *k*, *e*, and *f* of the magnitude.

Posit represents the number $(-1)^s \cdot useed^k \cdot 2^e \cdot (1 + f/2^{fs})$.

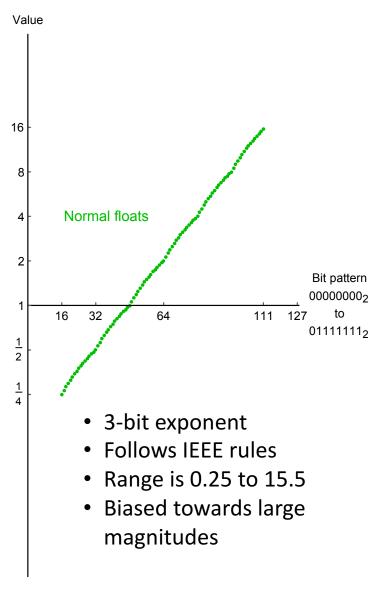
Better for circuit design, but cryptic:

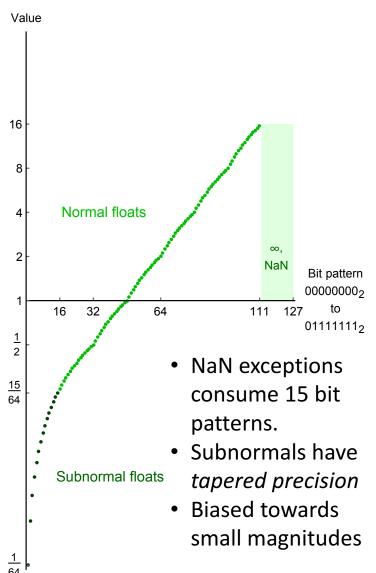
Check for 0 and NaR cases while finding s, k, e, f concurrently.

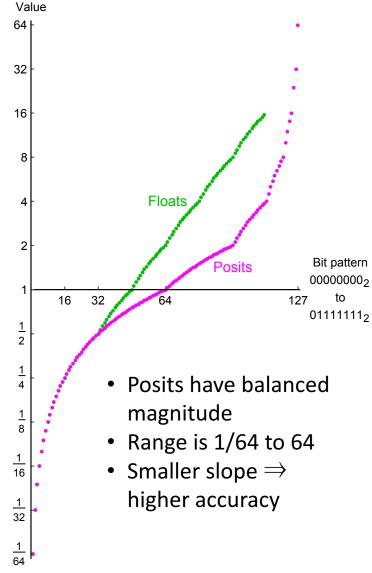
Posit represents
$$(1 - 3s + f) \cdot 2^{(-1)^s(k \cdot 2^{es} + e + s)}$$

From 1 to 2 if s = 0From -2 to -1 if s = 1

Visualization of 8-bit floats vs 8-bit posits







End of Posit Arithmetic Introduction