

UNITED STATES

SAN JOSE, CA, USA FEBRUARY 24-27, 2025

Beyond Integers and Floating Point Designing and Verifying with Alternate Number Representations

Russell Klein

Siemens EDA





Agenda

- Numbers and their representation
 - Integer, floats, fixed points, and more
- Why change representations?
 - Power, Performance, Area
- Impacts on algorithms
 - Edge conditions rounding, overflow, and saturation
- Impacts on verification



A Bit of Number Theory

- Numbers are infinite
 - Infinitely large
 - For both the positive and negative values
 - Infinitely small
- Modern (and even some ancient) math relies on this
 - Algebra, Calculus, Trigonometry are all deeply rooted in this concept



Representations

- Any representation we create will be practically finite
 - Only so much room on the paper
- Arabic was the most extensible
 - But still limited for very large and very small numbers

Arabic	Egyptian	Babylonian	Greek	Roman	Chinese	Aztec	Cretan	Mayan
1		γ	α	Ι	-	•	'	•
2		ŶΫ	β	II	1	••	"	••
3		$\gamma\gamma\gamma$	γ	III	111	•••	***	•••
4		$\gamma\gamma\gamma\gamma$	δ	Ш	四			••••
5		77777	3	V	Æ			—
6		$\overline{A}\overline{A}\overline{A}\overline{A}$	د	VI	六	•••••		-
7		$\begin{array}{c} \overline{\gamma}\overline{\gamma}\overline{\gamma}\overline{\gamma}\overline{\gamma}\\ \overline{\gamma}\overline{\gamma}\end{array}$	ζ	VII	七	••••		
8	1111111	<u> </u>	η	VIII	八	••••		
9		$\begin{array}{c} \overline{\mathbf{A}} \overline{\mathbf{A}} \overline{\mathbf{A}} \overline{\mathbf{A}} \overline{\mathbf{A}} \\ \overline{\mathbf{A}} \overline{\mathbf{A}} \overline{\mathbf{A}} \overline{\mathbf{A}} \end{array}$	θ	VIIII	九	••••		<u></u>
10	\cap	A	ι	X	+		۲	
20	$\cap \cap$	AA	к	XX	<u>_</u> +	Ρ		•
30	$\cap\cap\cap$	444	λ	XXX	三十	:::: P		$\cdot =$
40	$\cap\cap\cap\cap$	AAA A	μ	XXXX	四十	PΡ		••
50	$\cap\cap\cap\cap\cap$		ν	L	五十	:::: PP		=
60	$\begin{array}{c} \cap \cap \cap \cap \cap \\ \cap \end{array}$	Ŷ	ŝ	LX	六十	ppp		•••
70		۶۹	0	LXX	七十			=
80		₽₽₽	π	LXXX	八十	PPPP		••••
90		₽ <mark>4</mark> 4	Ŷ	LXXXX	九十			=
100	୨	AA AAA	ρ	С	一百	PPPPP	/	
200	୭୭	ALL ALL	σ	CC	二百	ppppp ppppp	//	=

J. Zhang, D. Norman, "Cognition," Vol 57, Issue 3, December 1995, Pages 271-295





Representations

- Sciences and engineering require representations for very big and very small numbers
- Scientific notations
 - Exponentiation
- Examples

Earth's mass	5,972,400,000,000,000,000,000,000,000 grams	5.9724 x 10 ²⁷
Electron's mass	0.000000000000000000000000000000000000	9.1094 x 19 ⁻²⁸



SI Units

Prefix		Desimal	Adoption	
Symbol	10	Decima	[nb 1]	
Q	10 ³⁰	1 000 000 000 000 000 000 000 000 000 0	2022[3]	
R	10 ²⁷	1 000 000 000 000 000 000 000 000 000	2022	
Y	Y 10 ²⁴ 1 000 000 000 000 000 000 000 000		1001	
Z 10 ²¹		1 000 000 000 000 000 000 000	1991	
Е	10 ¹⁸	1 000 000 000 000 000 000	1075[4]	
Р	10 ¹⁵	1 000 000 000 000 000	1975	
Т	10 ¹²	1 000 000 000 000	1060	
G	10 ⁹	1 000 000 000	1900	
М	10 ⁶	1 000 000	1873	
k	10 ³	1 000		
h	10 ²	100	1795	
da 10 ¹		10		
	efix Symbol Q R Y Z Z E P T G M K M k h da	Base Symbol 10 Q 10 ³⁰ Q 10 ²⁷ R 10 ²⁴ Y 10 ²⁴ Z 10 ²¹ E 10 ¹⁵ P 10 ¹² Q 10 ⁹ M 10 ⁶ K 10 ³ 10 10	efix Base Decimal Symbol 10 Decimal Q 10 ³⁰ 1 000 000 000 000 000 000 000 000 000 R 10 ²⁷ 1 000 000 000 000 000 000 000 000 000 Y 10 ²⁴ 1 000 000 000 000 000 000 000 000 000 Z 10 ²¹ 1 000 000 000 000 000 000 000 000 E 10 ¹⁸ 1 000 000 000 000 000 000 000 P 10 ¹⁵ 1 000 000 000 000 000 000 G 10 ¹² 1 000 000 000 000 000 000 M 10 ¹² 1 000 000 000 000 000 M 10 ⁶ 1 000 000 000 000 k 10 ³ 1 000 000 000 h 10 ² 1 000 000 000 h 10 ² 1 000 000 h 10 ² 1 000 h 10 ² 1 000	

deci	d	10 ⁻¹	0.1		
centi	С	10 ⁻²	0.01	1795	
milli	m	10 ⁻³	0.001		
micro	μ	10 ⁻⁶	0.000 001	1873	
nano	n	10 ⁻⁹	0.000 000 001	1060	
pico	р	10 ⁻¹²	0.000 000 000 001	1900	
femto	f	10 ⁻¹⁵	0.000 000 000 000 001	1064	
atto	а	10 ⁻¹⁸	0.000 000 000 000 000 001	1904	
zepto	z	10 ⁻²¹	0.000 000 000 000 000 000 001	1001	
yocto	y 10 ⁻²⁴ 0.000 000 000 000 000 000 000 001		1991		
ronto	r	10 ⁻²⁷	0.000 000 000 000 000 000 000 000 001	2022[3]	
quecto	q	10 ⁻³⁰	0.000 000 000 000 000 000 000 000 000 0	2022	



https://en.wikipedia.org/wiki/Metric_prefix



Computing Machinery Representations



• Integer- 32 bits and 64 bits



Integers

- Consistent precision across the entire range
- Discrete steps between represented values



Computing Machinery Representations



• Floating point – IEEE standard 32 bits and 64 bits



Floating point

- Not consistent precision, and also discrete steps between values
 - High precision near 0
 - Lower precision for bigger numbers
 - For IEEE float 32, numbers above 2²⁴ (16,777,216) are separated by 2.0





Problems

• Fun with Excel (running on a 64-bit CPU, in 2025)

	1.23F-100		1.23E-100		3.45E+100		1.23E-05
	1 23E-100		1.23E-100		7.89E+120		1.23E-05
sum.	2.46E-100		3.45E+100		-3.45E+100		3.45E+05
sum.	2.402-100		7.89E+120		-7.89E+120		7.89E+09
			-3 45E+100		1.23E-100		-3.45E+05
			-7 89E+120		1.23E-100		-7.89E+09
		sum:	0	sum:	2.46E-100	sum:	2.48E-05
		Sum.	0				



Problems

• Fun with Python (really numpy)

```
>>>
[>>>
                                                                           [>>>
                                    [>>> import numpy
[>>> import numpy
                                                                           >> import numpy
                                    [>>> a = numpy.float32(12345)
[>>> a = numpy.float32(123)
                                                                           >>> b = numpy.float32(123*10**9)
                                    [>>> b = numpy.float32(123*10**9)
[>>> b = numpy.float32(123*10**9)
                                                                           >>> b - numpy.float32(4096.0)
                                    >>> a
>>> a
                                                                           123000000000.0
123.0
                                    12345.0
                                                                           >>> b - numpy.float32(4097.0)
|>>> a = a + b
                                    |>> a = a + b
                                                                           122999990000.0
|>>> a = a - b
                                    |>> a = a - b
                                                                           >>>
                                                                           1.
>>> a
                                    >>> a
0.0
                                    16384.0
                                    >>>
```



Algorithms Implemented in Hardware

- Fast Fourier transformations
- Audio filters
- Gaussian blurring
- Convolutions
- Inferencing
- Trigonometry
- And many more

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

$$\cos \alpha + \cos \beta = 2\cos \frac{1}{2}(\alpha + \beta)\cos \frac{1}{2}(\alpha - \beta)$$

$$(x+a)^n = \sum_{k=0}^n \binom{n}{k} x^k a^{n-k}$$

$$f(x) = a_0 + \sum_{n=1}^{\infty} \left(a_n \cos \frac{n\pi x}{L} + b_n \sin \frac{n\pi x}{L} \right)$$



Algorithms Implemented in Hardware

- Based on math
- Generally, there will be some software reference implementation
 - Presumably, well verified for some set of inputs
- Usually, computationally complex
 - Simple, lightweight computations can be left in software



The "Range" Problem

• An average function

```
int average(int data[], int count)
{
    int sum = 0;
    int i;
    if (count == 0) return sum;
    for (i=0; i<count; i++) {
        printf("data[%d] = %d \n", i, data[i]);
        sum += data[i];
    }
    return sum/count;
}</pre>
```

```
data[0] = 10
data[1] = 20
data[2] = 30
data[3] = 40
data[4] = 50
data[5] = 60
data[5] = 60
data[6] = 70
data[7] = 80
data[8] = 90
data[9] = 100
average is = 55
```



The "Range" Problem

• An average function

```
int average(int data[], int count)
{
    int sum = 0;
    int i;
    if (count == 0) return sum;
    for (i=0; i<count; i++) {
        printf("data[%d] = %d \n", i, data[i]);
        sum += data[i];
    }
    return sum/count;
}</pre>
```

```
data[0] = 217897144
data[1] = 417876878
data[2] = 422077472
data[3] = 418769530
data[4] = 439176704
data[5] = 505877146
data[6] = 596749270
data[7] = 406971490
data[8] = 439084336
data[9] = 430487396
average is = 7
```



Math in Hardware

- In a design, the representation can be anything we want
 - Obviously, integer registers are sized to be just big enough
- Why deviate from the reference implementation?
 - Power
 - Performance
 - Area





Source: Nvidia DAC2017





Source: Nvidia DAC2017





Source: Nvidia DAC2017





SAN JOSE, CA, USA FEBRUARY 24-27, 2025





Source: Nvidia DAC2017



Source: Nvidia DAC2017

Source: Nvidia DAC2017

Source: Nvidia DAC2017

Smaller is Better

A one-bit integer multiplier is an "and" gate

Smaller is better

Bit Width vs Relative Delay

- Propagation delay is shorter when there are fewer gates
- 32-bit adder has a carry chain through 32 operators

Source: Nvidia DAC2017

Energy Cost of Data Movement

Source: Nvidia DAC2017

Energy Cost of Data Movement

- Reducing the size of the data means:
 - Less data to store
 - Less data to move
- Algorithms that need to be in hardware often have large data sets
 - This is most pronounced in AI algorithms
- Reduce memory footprint by 2-3X
- Reduce movement costs by 2-3X

Alternate Representations

- Smaller integers
- Fixed point
- Varying sized floating point
 - bfloat_16
- Posits
- Indexes

Varying Sized Integers

- Obvious (included for completeness)
- You don't use a 32-bit integer when you have data values from 1 to 10
- Synthesis tools sometimes trim integers for you
 - If it can work out that bits are not needed

Range: -256 to 255 Precision: 1

Fixed Point

• Integer with an implied "binary point"

Fixed Point

- Uses almost the same operators as integers
 - Multiplication and division operations need to mind the binary point

0.1 x 0.1 = 0.01 1 x 1 = 1 10 x 10 = 100

• Consistent precision across the entire range of value represented

Varying Sized Floating-Point

• Arbitrary sized exponent and mantissa

Range: -64K to 64K Precision: Varies, best 0.0078

Varying Sized Floating-Point

- Better precision near 0 and larger range than equivalent integer or fixed point
- For many computations standard sized floats are overkill
 - Excess bits in data storage, data movement, and operators

Bfloat-16 (or Google's Brain Float)

- Inferencing algorithms usually have numbers near 0
- As numbers get big precision is less important
 - Just knowing that the quantity is "big" is often sufficient
- Al researchers at Google came up with 16-bit floating point optimized for Al training and inferencing
 - Smaller data, less data movement, and smaller silicon have a big impact on training costs
- Google's TPUs support this format
- RISC-V extensions can support this in vector unit

Bfloat-16

• Brain float

Range: -10³⁸ to 10³⁸ Precision: Varies

Bfloat-16 to IEEE Float-32

• Key value is fast conversion to and from native float-32

Other Notable Floating-Point Formats

Posits

- Gustafson and Yonemoto developed these in 2017 for AI
- Adds "regime" bits
 - Essentially, an exponent for the exponent
- Gives better precision near zero
 - And larger range than equivalently sized floating point
- Claim significantly better accuracy for training and inferencing with smaller representation

Posits

Source: https://spectrum.ieee.org/floating-point-numbers-posits-processor

Indexes (or Look-up Table)

- For "lumpy" non-contiguous irregular data
- Can potentially save space (and movement) by loading values into a memory and storing/moving an index

Alternate Representations

- Smaller integers
- Fixed point
- Varying sized floating point
 - bfloat_16
- Posits
- Indexes

- This was sample of a few ways to represent numbers for computation
- Limited only by designer creativity

Computing with Alternate Formats

- Any change to the representation will change the math
 - And can affect the results of computations
 - There will be less precision and a smaller range
- Key issues
 - Overflows
 - Rounding

Overflow

- Overflow happens when the result of the computation is larger than can be represented by the register/variable
- For integers, high order bits are dropped
 - Result is incorrect, all subsequent calculations are corrupted
- For floats, value becomes INF (infinity)
 - Subsequent operations result in either INF or NaN

Overflow

- With smaller representations, overflow is a bigger concern
- With 32-bit and 64-bit CPUs overflow is uncommon
 - Older 8/16 bit CPUs have overflow flags and "add with carry" instructions
- Overflow handling at the HW level is rarely propagated to SW
- Fun Fact
 - RISC-V architecture has no ADD with CARRY instruction
 - And no carry/overflow/borrow bit in ISA

Saturation

- Saturating math sets the output value to the maximum representable value
 - Does not just drop bits
- In some cases, this produces acceptable results

Saturating Math

- Saturating math assigns largest possible value instead of dropping bits
- With 10-bit signed fixed-point number with 7 integer bits:

	Overflow:	Saturation:				
62.5 + 2.0	$\begin{array}{c} 0 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 0 \ 0 \\ 1 \ 0 \ 0 \ 0 \ 0 \end{array}$	62.5 0 1 1 1 1 1 1 . 1 0 0 + 2.0 1 0 . 0 0 0				
- 1.5	1 1 1 1 1 0 1 . 1 0 0	63.875 0111111.111				
REALLY WI	RONG!	Close to correct				

Rounding

- Integer adders truncate results:
 - In C code: printf(" %d \n", 1000 * (99/100));
 - Results in "0"
- Assigning integers from any other type truncates low order bits
 - Consistent in most programming languages and with bit-vector types in Verilog
- As representation's range shrinks, truncation error becomes more significant

Rounding Models

- Many different rounding models
- Rounding will affect results
- All verification models must be bit-level consistent

Image source: https://en.wikipedia.org/wiki/Rounding

Examples

- Inferencing accelerator
 - MNIST handwritten character recognition
 - Conversion from float32 (Python) to 10 bit fixed point (Verilog)
- Audio processing
 - MFCC calculation (MEL Frequency Coefficient Calculation)
 - Fast Fourier transformation of PCM audio to
 - Conversion from double (C++) to bespoke 8-bit float (Verilog)

MNIST Handwritten Character Recognition

- "Hello, World!" example for machine learning
 - Small, but runs millions of multiply accumulate operations
 - Performs many 2d convolutions and vector/matrix multiplications
 - Good platform for exploring implementation and optimization methods

Fixed Point Size Analysis

- High-water mark of data and intermediate values showed range of values was -37 to 56
 - Float32 (+/-10³⁸ is excessive)
- Sensitivity analysis performed across varying fixed-point representations

Fixed Point Conversion Results

SY

					Integer E	Bits				
		8	7	6	5	4	3	2	1	0
	8	98.05	98.05	98.05	97.55	76.75	28.70	18.00	16.80	14.90
(0 –	7	97.85	97.85	97.85	97.25	75.39	27.90	17.50	16.60	15.40
Sits	6	97.13	97.95	97.91	97.45	75.15	28.30	17.30	15.90	13.90
	5	97.21	98.08	98.10	97.40	72.57	24.50	16.90	15.20	14.90
0 U	4	96.94	97.79	97.76	95.71	59.90	21.40	16.20	13.10	15.10
tic	3	95.56	96.37	96.35	90.08	38.83	16.70	14.00	11.50	12.70
rao	2	82.31	83.13	83.13	64.73	22.70	14.90	12.30	10.50	8.50
ш	1	30.15	30.97	30.92	33.72	32.07	24.60	34.90	12.30	8.50
	0	9.53	9.33	9.50	9.37	9.37	8.50	8.50	8.50	10.00
	Accuracy									
Computational resources: $\overleftrightarrow{1}$ Area = -78% $\overleftrightarrow{2}$ Speed = 2.4X faster $\vcenter{2}$ Energy = -92%										
SYSTEMS INITIATIVE	•									DESIGN AND DVC CONFERENCE A UNITED BAN JOB FEBRUARY 2

, USA 7, 2025

MFCC – MEL Frequency Cepstrum Coefficients

- Coverts audio waveform into an audio 2D spectrogram
 - "MEL" frequencies correspond to human ear sensitivities
 - Computes audio energy in a range of frequencies and their changes over time
 - An FFT is computed for each frequency and timeslice
- Spectrograms are used in voice to text, wakewords, and voiceprinting

Float32 to 8-bit Float

Verification of Modified Algorithm

- Need to fully verify operators, computational operations
- Also need to verify that algorithm produces correct results
 - with modified representation
- Algorithms are likely too large to analytically size representation
 - MNIST had > 4 million MAC operations
 - MFCC performed 2,000 FFT calculations on input waveform
- Need to run large data sets
 - Both real and synthetic data
 - Verification will only be as good as the stimuli

Key Verification Points

- High and low watermarks for inputs, output, and intermediate data
 - Need to understand largest and smallest values represented at each point in the algorithm
- When overflows occur
 - Need to know anytime operation results in overflow of the representation
- These observations will change with each alteration to the representation
- Will require multiple runs at different representations
- Must have bit level accuracy matching ultimate RTL implementation

Run in Tandem with Original Algorithm

- Run modified algorithm in parallel with the original algorithm
- Perform on-the-fly comparisons
 - Ensure needed accuracy is maintained
 - Monitor both outputs and intermediate values

Performance Needed

- Algorithms are likely large and complex
 - Else would be left in software
- Logic simulation is probably too slow
 - Yolo-2 small object recognition algorithm at RTL takes ~5 hours 1 inference
 - About 6 CPU years for logic simulation of 10,000 inferences
- Raise the abstraction level
 - Can run 1,000X faster
 - But need to ensure bit level activity of math is identical to intended implementation

AC Data Types

- Bit accurate C++ library of numeric representations and operators
 - Templated typedefs for various representations
 - Extensions for bit level manipulation, for hardware operations
- Supports bit level accurate representation of various data formats
 - Overloaded operators for all types perform bit level accurate operations
 - Easy migration from C++ or other high-level languages
- Supports multiple rounding models and saturating math

http://hlslibs.org

AC Data Types

- Variable sized int
- Fixed point numbers
- Variable sized float
 - Bfloat-16
- IEEE float 32 and 64
- complex numbers

- ac_int<size, signedness>;
- ac_fixed<size, integer_bits, signedness>
- ac_float<size, exponent_size>;
- ac_bfloat<>;
- ac_ieee_float32<>, ac_ieee_float64<>;
- ac_complex<type>;

Value Range Analysis

- Records high and low watermarks for variables throughout algorithm
- Identifies overflow and underflow points
- Essentially, assertions and monitors in the ac_type classes
- This information enables zeroing in on the correct representation quickly

Synthesis from C++ for AC types

Catapult Hight-Level Synthesis

Architectural Exploration

User controls implementation with constraints and pragmas

Verification of RTL against C++

- Compare C++ against RTL with
 - Formal equivalency (SLEC)
 - Coverage

Technology Trends

- CPU single thread performance has stalled
- Performance demands continue to rise
- More algorithms moving to hardware

Conclusion

- Using alternative numeric representations when implementing designs can greatly improve PPA
- Need to carefully migrate from original algorithm to new representations, then to RTL
 - Continuous verification against the prior implementation is essential
 - Compare not just results but intermediate values as well
- Must ensure new representation maintains fidelity of the algorithm

Thank you

Questions or comments

?? || //

<u>Russell.Klein@siemens.com</u> <u>https://eda.sw.siemens.com/en-US/ic/catapult-high-level-synthesis/</u> http://hlslibs.org

