

NUMBER REPRESENTATIONS

Lecture-12

Binary numbers

Binary


0 and 1

Integers

Place values

$$\underline{0101}_b = 0 \times \underline{2^3} + 1 \times \underline{2^2} + 0 \times \underline{2^1} + 1 \times \underline{2^0} = 5_d$$

Negative Numbers

- No  sign: compensate using convention
 - Signed magnitude: 1st bit indicates sign
 - Two's complement

Negative Numbers

- No $-$ sign: compensate using convention
 - Signed magnitude: 1st bit indicates sign
 - Two's complement

2's complement

$$\underline{\underline{0101}}_b = 0 \times (-2^3) + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 = \underline{\underline{5}}_d$$

$$1101_b = 1 \times (-2^3) + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 = \underline{\underline{-3}}_d$$

■ Subtraction is now almost identical to addition

Dynamic range

- Ratio of largest magnitude to smallest (non-zero) magnitude

Dynamic range

■ Ratio of largest magnitude to smallest (non-zero) magnitude

eg Decimal 4 digits (+ve):

- Largest: 9999
- Smallest: 1
- Dynamic range = 9999
 - Usually expressed in dB: $20 \log D \approx 80dB$

↓
 \log_{10}

Dynamic range

Ratio of largest magnitude to smallest (non-zero) magnitude

eg Decimal 4 digits (+ve):

- Largest: 9999
- Smallest: 1
- Dynamic range = 9999
 - Usually expressed in dB: $20 \log D \approx 80dB$
- Binary:
 - Each bit doubles range: $20 \log 2 \approx 6dB$
 - 16 bits $\approx 96dB$

(no 2's c)
eg. 5-bit +ve : 11111 $\rightarrow 31$
00001 $\rightarrow 1$
 $DR = \frac{31}{1} = 31$
 $20 \log_{10} 31$

Real-valued numbers

- Convention: Choose location of *decimal* or *binary* point

Scaling

Choice of scaling factor determines **resolution**

Real-valued numbers

Convention: Choose location of *decimal* or *binary* point

Scaling

Choice of scaling factor determines **resolution**

- Decimal:
 - xxxxx => smallest increment = 1
 - xxx.xx => smallest increment = 0.01
 - x.xxxx => smallest increment = 0.0001

Real-valued numbers

Convention: Choose location of *decimal* or *binary* point

Scaling

Choice of scaling factor determines **resolution**

- Decimal:
 - $xxxxx \Rightarrow$ smallest increment = 1
 - $xxx.xx \Rightarrow$ smallest increment = 0.01
 - $x.xxxx \Rightarrow$ smallest increment = 0.0001

Dynamic range is the same for all:
 $99999/1 = 999.99/0.01 = 9.9999/0.00001 \approx 100dB$

Arithmetic: Addition

Handwritten diagram illustrating the scaling factor for binary addition:

- $0101.01 = 5.25$ (Scaling factor: $25/4$)
- $+ 01.0001 = 1.0625$ (Scaling factor: $17/2^4 = 17/16$)
- Result: $0110.0101 = 6.3125$

The diagram shows the scaling factors $25/4$ and $17/2^4 = 17/16$ being used to align the binary points for addition.

■ Align the binary points

Arithmetic: Addition

$$\begin{array}{r} \text{0101.01} \quad = 5.25 \\ + \text{01.0001} \quad = 1.0625 \\ \hline \text{0110.0101} \quad = 6.3125 \end{array}$$

Handwritten red annotations: A bracket above the first number is labeled '6'. A bracket below the second number is labeled '6'. A bracket below the result is labeled '8 bits'.

■ Align the binary points

- Truncate output if needed?
- Overflow possible, same as with integers

Arithmetic: Multiplication

| | |
|----------------|-----------------|
| <u>0101.01</u> | -> Multiplicand |
| 01.0001 | -> Multiplier |
| ----- | |
| .010101 | $\times 2^{-4}$ |
| -0.00000 | $\times 2^{-3}$ |
| 00.0000 | $\times 2^{-2}$ |
| 000.000 | $\times 2^{-1}$ |
| 0101.01 | $\times 2^0$ |
| 00000.0 | $\times 2^1$ |
| ----- | |
| 00101.100101 | |

Partial product.

Product.

Scaling factors:

Significand / Mantissa.

$$5.25_d = 0101.01_b = (\underline{010101})_b \times 2^{-2}$$

6-bit *mantissa* = 21 and scale factor = 2^{-2}

Scaling factors:

$$5.25_d = 0101.01_b = (010101)_b \times 2^{-2}$$

6-bit *mantissa* = 21 and scale factor = 2^{-2}

$$1.0625_d = 01.0001_b = (010001)_b \times 2^{-4}$$

6-bit *mantissa* = 17 and scale factor = 2^{-4}

VIVADO HLS

ap_fixed<>

ap_fixed

- Arbitrary precision fixed point
 - `ap_fixed<W,I,Q,0,N>`
-

| Identifier | Description | | |
|------------|---|----------------|--|
| W | Word length in bits | | |
| I | The number of bits used to represent the integer value (the number of bits above the binary point) | | |
| Q | Quantization mode | | |
| | This dictates the behavior when greater precision is generated than can be defined by smallest fractional bit in the variable used to store the result. | | |
| | SystemC Types | ap_fixed Types | Description |
| | SC_RND | AP_RND | Round to plus infinity |
| | SC_RND_ZERO | AP_RND_ZERO | Round to zero |
| | SC_RND_MIN_INF | AP_RND_MIN_INF | Round to minus infinity |
| | SC_RND_INF | AP_RND_INF | Round to infinity |
| | SC_RND_CONV | AP_RND_CONV | Convergent rounding |
| | SC_TRN | AP_TRN | Truncation to minus infinity (default) |
| | SC_TRN_ZERO | AP_TRN_ZERO | Truncation to zero |

| Identifier | Description | | |
|------------|---|----------------|----------------------------|
| O | Overflow mode. This dictates the behavior when the result of an operation exceeds the maximum (or minimum in the case of negative numbers) value which can be stored in the result variable. | | |
| | SystemC Types | ap_fixed Types | Description |
| | SC_SAT | AP_SAT | Saturation |
| | SC_SAT_ZERO | AP_SAT_ZERO | Saturation to zero |
| | SC_SAT_SYM | AP_SAT_SYM | Symmetrical saturation |
| | SC_WRAP | AP_WRAP | Wrap around (default) |
| | SC_WRAP_SM | AP_WRAP_SM | Sign magnitude wrap around |
| N | This defines the number of saturation bits in the overflow wrap modes. | | |

Scaling factors:

$$5.25_d = 0101.01_b = (010101)_b \times 2^{-2}$$

ap-fixed $\langle 6, 4 \rangle$

6-bit mantissa = 21 and scale factor = 2^{-2}

$$1.0625_d = 01.0001_b = (010001)_b \times 2^{-4}$$

ap-fixed $\langle 6, 2 \rangle$

6-bit mantissa = 17 and scale factor = 2^{-4}

$$5.578125_d = 00101.100101_b = (00101100101)_b \times 2^{-6}$$

ap-fixed $\langle 12, 6 \rangle$

Product: 12-bit mantissa = $21 \times 17 = 357$ and scale factor = $2^{-2} \times 2^{-4}$

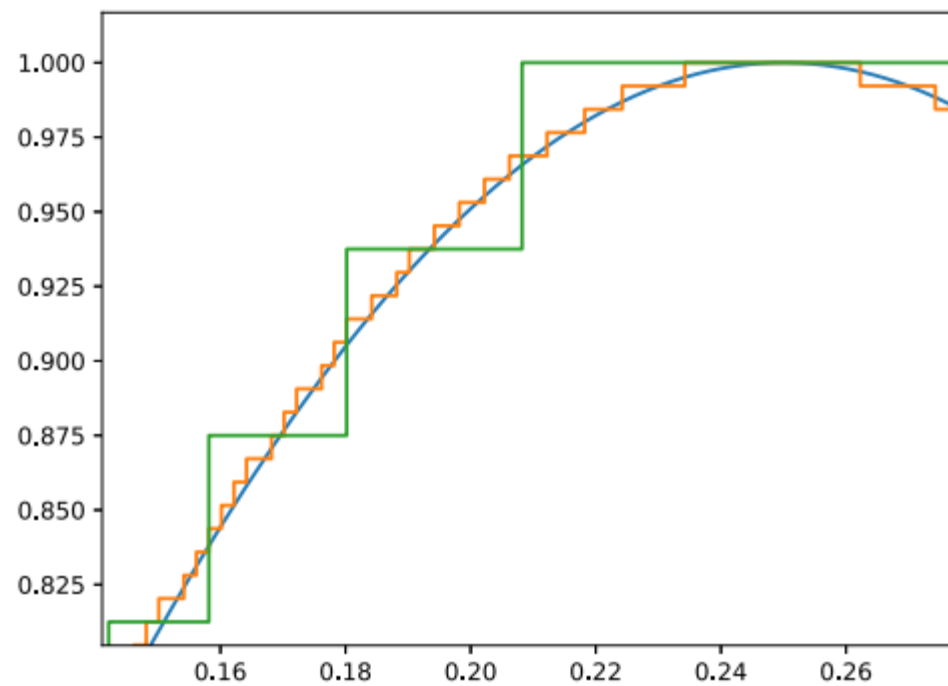
QUANTIZATION AND ERROR

Lecture 12

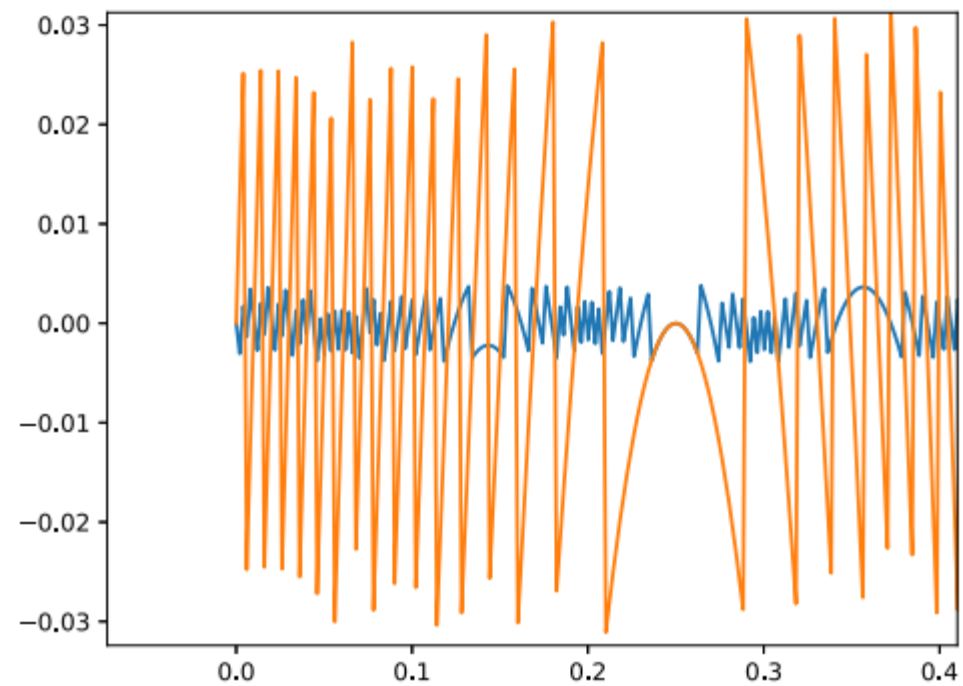
Quantization

- Magnitude
 - Largest value to be represented: integer
- Resolution
 - Smallest *delta* to be resolved: fraction

Quantization effects



Quantization Error



Error

- Assume uniformly distributed:
 - Model as **additive noise**
- * Signal to Noise Ratio (SNR)
- Roughly 6dB per bit of additional precision

FLOATING POINT AND OTHER NUMBER REPRESENTATIONS

Lecture 12

The scaling factor problem

- Example: SPICE circuit simulation

The scaling factor problem

Example: SPICE circuit simulation

- Typical ranges of values (for low-power electronic circuits)
 - Voltages: up to about 5.0V max, resolution in mV (for references etc.)
 - Currents: nA to mA (leakage currents to bias currents)
 - Impedance: $k\Omega$ to $M\Omega$

$\sim DR : 1000s.$

$\sim 10^6$

$\rightarrow 10^3$

The scaling factor problem

Example: SPICE circuit simulation

- Typical ranges of values (for low-power electronic circuits)
 - Voltages: up to about 5.0V max, resolution in mV (for references etc.)
 - Currents: nA to mA (leakage currents to bias currents) 10^{-9}
 - Impedance: $k\Omega$ to $M\Omega$ - 10^6
- Dynamic range
 - with individual units (μA), $k\Omega$ etc.) $\approx \underline{10^6}$ - $\underline{20}$ bits
 - SI units: nA to $M\Omega$ - $\underline{10^{15}}$ - $\underline{45}$ bits

Efficient use of bits

Efficient use of bits

- 0.000123
- 12.3
- 123000000

Efficient use of bits

- 0.000123
- 12.3
- 123000000

Significant

Scaling factors.

- 1.23×10^{-4}
- 1.23×10^1
- 1.23×10^8

Efficient use of bits

- 0.000123
- 12.3
- 123000000

- 1.23×10^{-4}
- 1.23×10^1
- 1.23×10^8

Significant bits / digits

Exponent

Efficient use of bits

- 0.000123
- 12.3
- 123000000

Significant bits / digits

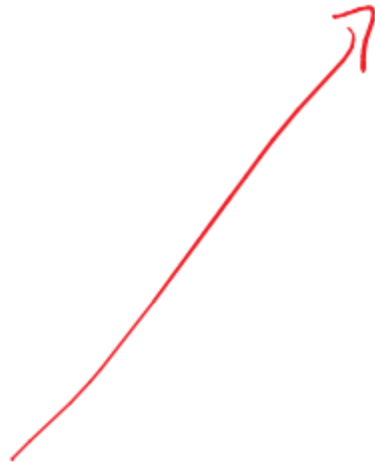
Exponent

Options

- Each block has own exponent
 - Each number has own exponent
-

- digit before point is not 0.
- 1.23×10^{-4}
 - 1.23×10^1
 - 1.23×10^8

Normalized
scientific
notation



Self-contained representation

- **Mantissa**: Significant bits - numerical value (1.23)
- **Exponent**: Scaling factor ← *Fixed pt : implicit.*
- **Sign**: Use a separate bit to represent this

Self-contained representation

- **Mantissa**: Significant bits - numerical value (1.23)
- **Exponent**: Scaling factor
- **Sign**: Use a separate bit to represent this

Example: IEEE 754 SP *Single Precision*

| | | |
|---|----------|----------|
| S | Exponent | Mantissa |
|---|----------|----------|

32b

| | | | | | |
|-----|-------------------------|-------------|--------------------------|----------|-------------------|
| 0/1 | e = 1 to 254 offset 127 | 1.mmm...mmm | Max: $2.0 \cdot 2^{-23}$ | Min: 1.0 | MinRes: 2^{-23} |
|-----|-------------------------|-------------|--------------------------|----------|-------------------|

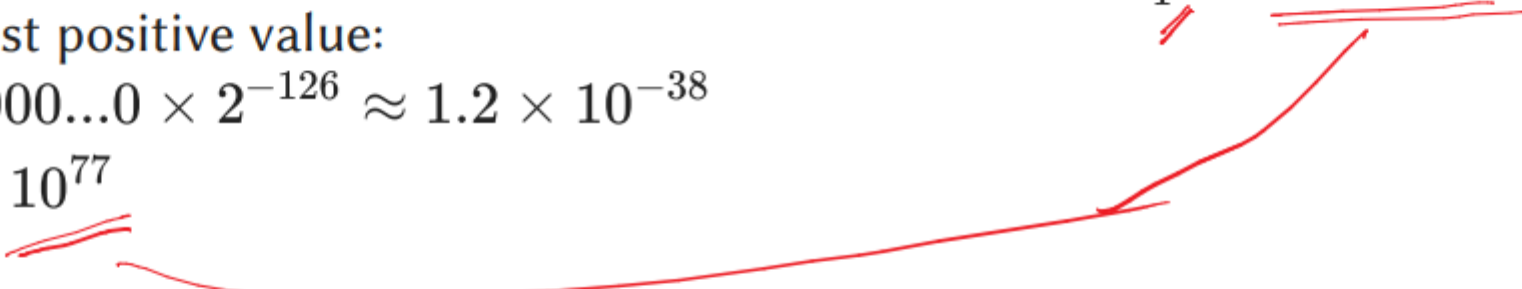
1
e ← *8*
Implicit ← *1* 0101110...²³ . . . 1

→ *1.010110...* *1 × 2^e*
24 bits.

Dynamic range

- Largest positive value:
 - $\underline{1.111\dots1} \times \underline{2^{+127}} \approx \underline{3.4 \times 10^{38}}$
- Smallest positive value:
 - $1.000\dots0 \times \underline{2^{-126}} \approx \underline{1.2 \times 10^{-38}}$
- DR: $\approx \underline{10^{77}}$

Dynamic range

- Largest positive value:
 - $1.111...1 \times 2^{+127} \approx 3.4 \times 10^{38}$
 - Smallest positive value:
 - $1.000...0 \times 2^{-126} \approx 1.2 \times 10^{-38}$
 - DR: $\approx 10^{77}$
 - Contrast 32-bit integer:
 - DR: $\frac{2^{31}}{1} \approx 2.2 \times 10^9$
- 

Dynamic range


- Largest positive value:
 - $1.111...1 \times 2^{+127} \approx 3.4 \times 10^{38}$
- Smallest positive value:
 - $1.000...0 \times 2^{-126} \approx 1.2 \times 10^{-38}$
- DR: $\approx 10^{77}$

- Contrast 32-bit integer:
 - DR: $\frac{2^{31}}{1} \approx 2.2 \times 10^9$

Problem?

- Representing 9 digit integers?
 - 123456789 \rightarrow 123456792

Special cases

- exponent = 255
 - mantissa = all zeros: $\pm\infty$  *Infinity.*
 - mantissa non-zero: NaN (not-a-number: eg 0/0, $\infty - \infty$ etc.)

Special cases

- exponent = 255
 - mantissa = all zeros: $\pm\infty$
 - mantissa non-zero: NaN (not-a-number: *eg* $0/0$, $\infty - \infty$ etc.

Addition

$$1.23 \times 10^{10}$$

$$3.45 \times 10^4$$

$$? \times 10^?$$

Addition

$$\begin{array}{r} 1.23 \times 10^{10} \\ 3.45 \times 10^4 \\ \hline ? \times 10^? \end{array}$$

- Convert to same exponent
 - Compare exponents, shift mantissas

Addition

$$\begin{array}{r} 1.23 \quad \times 10^{10} \\ 0.000345 \times 10^{10} \\ \hline ? \quad \times 10^? \end{array}$$

- Convert to same exponent
 - Compare exponents, shift mantissas

Addition

$$\begin{array}{r} 1.23 \quad \times 10^{10} \\ 0.000345 \times 10^{10} \\ \hline ? \quad \times 10^? \end{array}$$

- Convert to same exponent
 - Compare exponents, shift mantissas
- Add mantissas

Addition

$$\begin{array}{r} 1.23 \quad \times 10^{10} \\ 0.000345 \times 10^{10} \\ \hline 1.230345 \times 10^? \end{array}$$

- Convert to same exponent
 - Compare exponents, shift mantissas
- Add mantissas

Addition

$$\begin{array}{r} 1.23 \quad \times 10^{10} \\ 0.000345 \times 10^{10} \\ \hline 1.230345 \times 10^? \end{array}$$

- Convert to same exponent
 - Compare exponents, shift mantissas
- Add mantissas
- Round

Addition

$$\begin{array}{r} 1.23 \quad \times 10^{10} \\ 0.000345 \times 10^{10} \\ \hline 1.23 \quad \times 10^? \end{array}$$

- Convert to same exponent
 - Compare exponents, shift mantissas
- Add mantissas
- Round

Addition

$$\begin{array}{r} 1.23 \quad \times 10^{10} \\ 0.000345 \times 10^{10} \\ \hline 1.23 \quad \times 10^? \end{array}$$

- Convert to same exponent
 - Compare exponents, shift mantissas
- Add mantissas
- Round
- Update exponent of result

Addition

$$\begin{array}{r} 1.23 \quad \times 10^{10} \\ 0.000345 \times 10^{10} \\ \hline 1.23 \quad \times 10^{10} \end{array}$$

- Convert to same exponent
 - Compare exponents, shift mantissas
- Add mantissas
- Round
- Update exponent of result

Extra steps
compared to
integer / fixed point

Addition - exponent change

$$\begin{array}{r} 1.23 \times 10^{10} \\ 9.99 \times 10^{10} \\ \hline ? \times 10^? \end{array}$$

Addition - exponent change

$$\begin{array}{r} 1.23 \times 10^{10} \\ 9.99 \times 10^{10} \\ \hline ? \times 10^? \end{array}$$

- Already same exponent
- Add mantissas

Addition - exponent change

$$\begin{array}{r} 1.23 \times 10^{10} \\ 9.99 \times 10^{10} \\ \hline 11.22 \times 10^? \end{array}$$

- Already same exponent
- Add mantissas
 - **overflow**

Addition - exponent change

$$\begin{array}{r} 1.23 \times 10^{10} \\ 9.99 \times 10^{10} \\ \hline 1.122 \times 10^? \end{array}$$


- Already same exponent
- Add mantissas
 - overflow
- Shift mantissa

Addition - exponent change

$$\begin{array}{r} 1.23 \times 10^{10} \\ 9.99 \times 10^{10} \\ \hline 1.12 \times 10^? \end{array}$$

- Already same exponent
- Add mantissas
 - overflow
- Shift mantissa
- Round

Addition - exponent change

$$\begin{array}{r} 1.23 \times 10^{10} \\ 9.99 \times 10^{10} \\ \hline 1.12 \times 10^{11} \end{array}$$


- Already same exponent
- Add mantissas
 - overflow
- Shift mantissa
- Round
- Update exponent

Addition - LOD

$$\begin{array}{r} 1.23456 \times 10^{10} \\ -1.23455 \times 10^{10} \\ \hline ? \times 10^? \end{array}$$

Addition - LOD

$$\begin{array}{r} 1.23456 \times 10^{10} \\ -1.23455 \times 10^{10} \\ \hline ? \times 10^? \end{array}$$

- Exponents aligned
- Add mantissas (**note**: subtraction)

Addition - LOD

$$\begin{array}{r} 1.23456 \times 10^{10} \\ -1.23455 \times 10^{10} \\ \hline 0.00001 \times 10^? \end{array}$$

- Exponents aligned
- Add mantissas (**note**: subtraction)

Addition - LOD

$$\begin{array}{r} 1.23456 \times 10^{10} \\ -1.23455 \times 10^{10} \\ \hline 0.00001 \times 10^? \end{array}$$

- Exponents aligned
- Add mantissas (**note**: subtraction)
- **Normalize!!!**

Leading One Detection

Addition - LOD

$$\begin{array}{r} 1.23456 \times 10^{10} \\ -1.23455 \times 10^{10} \\ \hline 1.00000 \times 10^{-5} \times 10^? \end{array}$$

- Exponents aligned
- Add mantissas (**note:** subtraction)
- Normalize!!!

Addition - LOD

$$\begin{array}{r} 1.23456 \times 10^{10} \\ -1.23455 \times 10^{10} \\ \hline 1.00000 \times 10^{-5} \times 10^? \end{array}$$

- Exponents aligned
- Add mantissas (**note**: subtraction)
- Normalize!!!
- Update exponent

Addition - LOD

$$\begin{array}{r} 1.23456 \times 10^{10} \\ -1.23455 \times 10^{10} \\ \hline 1.00000 \times 10^{-5} \times 10^{10} \end{array}$$

- Exponents aligned
- Add mantissas (**note:** subtraction)
- Normalize!!!
- Update exponent

Addition - LOD

$$\begin{array}{r} 1.23456 \times 10^{10} \\ -1.23455 \times 10^{10} \\ \hline 1.00000 \times 10^5 \end{array}$$

- Exponents aligned
- Add mantissas (**note:** subtraction)
- Normalize!!!
- Update exponent
- Final exponent

Variants

- Double precision
 - 1 sign, 11 exponent, 52 mantissa: 64 bits
 - DR: $\approx 10^{600}$

Variants

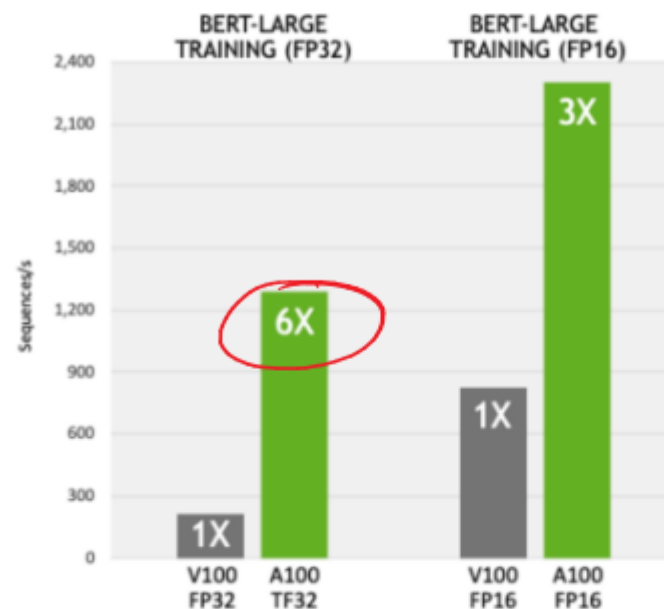
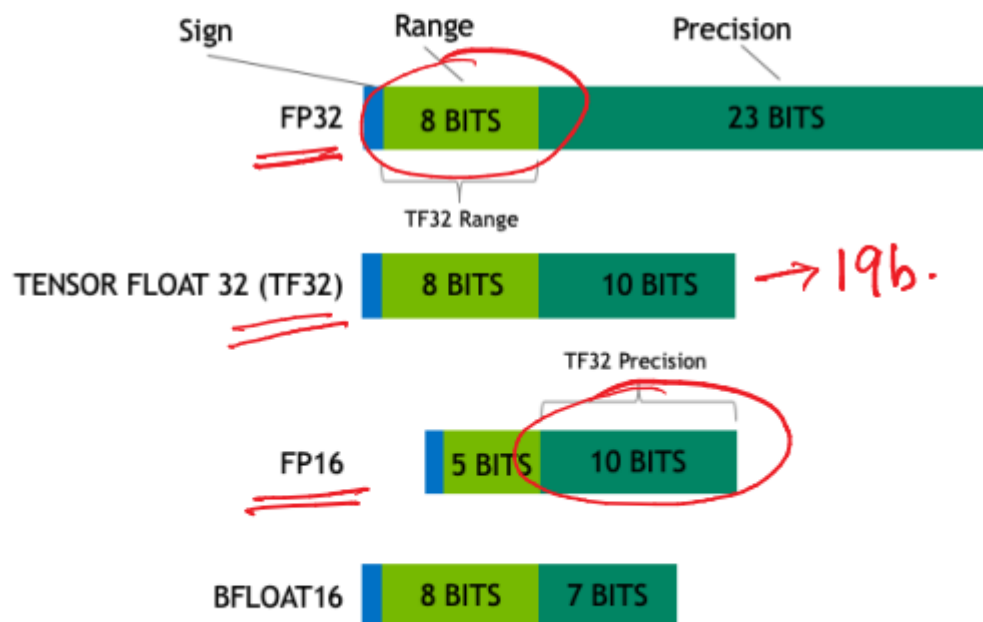
- Double precision
 - 1 sign, 11 exponent, 52 mantissa: 64 bits
 - DR: $\approx 10^{600}$
- Half precision
 - 1 sign, 5 exponent, 10 mantissa: 16 bits
 - DR: $\approx \underline{\underline{10^9}}$ \rightarrow DR of 32b int.

Variants

- Double precision
 - 1 sign, 11 exponent, 52 mantissa: 64 bits
 - DR: $\approx 10^{600}$
- Half precision
 - 1 sign, 5 exponent, 10 mantissa: 16 bits
 - DR: $\approx 10^9$
- Bfloat16
 - 1 sign, 8 exponent, 7 mantissa: 16 bits
 - DR: $\approx 10^{77}$! \rightarrow 32 b fp.
 - Easy conversion from single-precision
 - Much lower precision (eg $\pi \rightarrow$ 3.140625)

Example: TensorFloat-32

Src: nvidia.com -- NVIDIA -- May 2020



BERT Large Training (FP32 & FP16) measures Pre-Training phase, uses PyTorch including (2/3) Phase1 with Seq Len 128 and (1/3) Phase 2 with Seq Len 512, V100 is DGX1 Server with 8xV100, A100 is DGX A100 Server with 8xA100, A100 uses TF32 Tensor Core for FP32 training

Hardware complexity

| Integer | |
|---------|--------|
| Add | |
| 8 bit | 0.03pJ |
| 32 bit | 0.1pJ |
| Mult | |
| 8 bit | 0.2pJ |
| 32 bit | 3.1pJ |

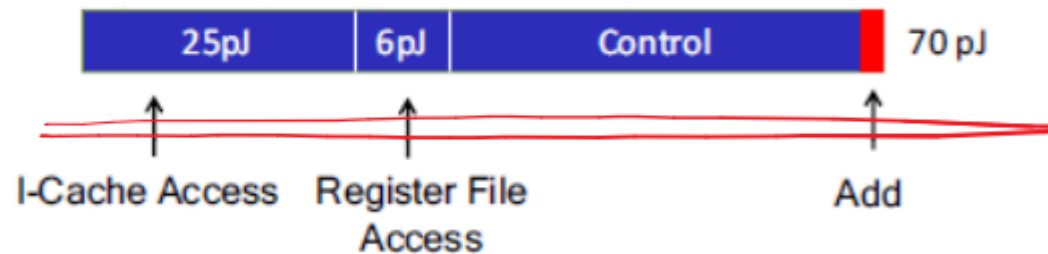
30x

| FP | |
|--------|-------|
| FAdd | |
| 16 bit | 0.4pJ |
| 32 bit | 0.9pJ |
| FMult | |
| 16 bit | 1.1pJ |
| 32 bit | 3.7pJ |

9x

| Memory | |
|---------------|-----------|
| Cache (64bit) | |
| 8KB | 10pJ |
| 32KB | 20pJ |
| 1MB | 100pJ |
| DRAM | 1.3-2.6nJ |

Instruction Energy Breakdown



Any Question...

Thank you