

The slide features a light beige background with four large, semi-transparent brown circles positioned in the corners. A thin brown rectangular border frames the central text area.

Posit-Based Adder and Multiplier Design for MAC Unit



CONTENTS

*INTRODUCTION
OF MAC UNIT*

POSIT ARITHMETIC

WHY POSIT?

*POSIT ADDER
ALGORITHM*

NEXT STEPS

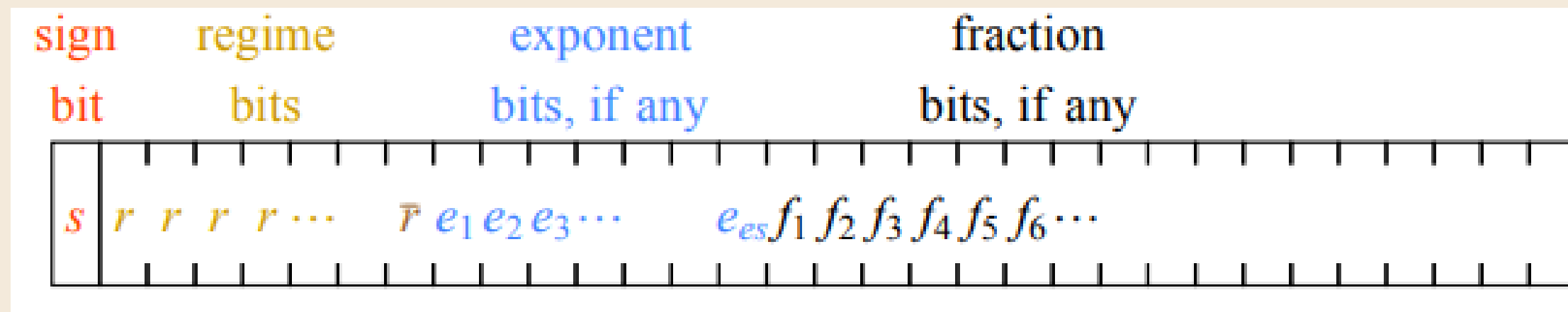
REFERENCES

Introduction of MAC Unit

- MAC - Multiplication & Accumulation Unit
- AI/ML algorithmic architectures have many variants including Deep Neural Networks (DNNs), which are popular due to their diverse applications in various fields viz. scientific, engineering, agriculture, healthcare etc.
- DNN consists of one input layer, one output layer, and multiple hidden layers. Each DNN layer performs matrix (vector) multiplication and accumulation, and a layer's output is input to the subsequent layers within the network.
- The MAC unit is present in each neuron of the hidden layers to perform the multiplication & accumulation operations.
- In general, the arithmetic representations for various numbers in neural network computations are IEEE-754 float, fixed-point, bfloat16 or Posit.

POSIT ARITHMETIC

- Posit is the latest development in Universal Number (unum) system under type-3 unum. Posit is claimed as a possible substitute for floating point (FP) number system.
- The structure of an N-bit Posit representation with ES exponent bits as follows-



- The Sign bit in posit is 0 for positive numbers and 1 for negative numbers. For the case of negative number, first take 2's complement before decoding regime, exponent and mantissa bits.
- There are only two exception cases: zero and infinity.
- Zero - (000...000) representation
- Infinity - (1000...000) representation.

- For all other cases, the value x of a posit is given by -

$$x = (-1)^{MSB} \times useed^k \times 2^{exp} \times \left(1 + \sum_{i=1}^{fn-1} b_{fn-1-i} 2^{-i}\right)$$

- The regime indicates a scale factor of $useed^k$ where $useed = 2^{(2^{ES})}$ and ES is the exponent size. The numerical value of k is determined by the run length of 0 or 1 bits in the string of regime bits.

With $ES=2$: $0_0001_11_1 = + (2^4)^{-3} \times 2^3 \times (1 + \frac{1.0}{2})$, and
 $11101011 \rightarrow 0_001_01_01 = - (2^4)^{-2} \times 2^1 \times (1 + \frac{1.0}{4})$,
 With $ES=3$: $0_110_101_1 = + (2^8)^1 \times 2^5 \times (1 + \frac{1.0}{2})$, and
 $10001111 \rightarrow 0_1110_001 = - (2^8)^2 \times 2^1 \times (1 + 0.0)$.

WHY POSIT?

- MAC Units Designs will differ in terms of Approximation, Architecture, Arithmetic & Algorithms.
- The selection of number representation scheme and the associated numeric computations play a significant role in efficient AI/ML hardware architecture design and the accuracy of machine learning inference.
- For DNN to have better inference accuracy, IEEE-754 floating-point representations are mostly preferred with a precision of 16-bit (half) or 32-bit (single) or bfloat16. However, the hardware design comes with higher resource utilization and more power consumption.
- For an efficient hardware design, fixed-point representation is used with Q4.4 or Q8.8 schemes.
- Recently, a new number system, i.e., Posit, has found a place with an ML inference accuracy level similar to floating-point but requiring fewer hardware resources like fixed-point arithmetic.

POSIT ADDER ALGORITHM (Currently working)

Algorithm 1 Regime, Exponent, Fraction Bits Extraction

```
1: Posit Size:  $N$ ; Exponent Size:  $es$ 
2: OutputTwos  $[N - 2 : 0]$   $\triangleright$  Input Size
3: Regime  $[\log_2(N) : 0]$   $\triangleright$  Input Size
4: FractionBits  $[N + 1 - es : 0]$   $\triangleright$  Output Size
5: ExpBits  $[es - 1 : 0]$   $\triangleright$  Exponent Size
6: Atest  $\leftarrow$  InputA  $[N - 1]$ 
   &  $\sim$  InputA  $[N - 2 : 0]$   $\triangleright$  Checking Special Case
7: OutputTwos  $\leftarrow \{[N]\{\text{InputA}[N - 1]\}\}$ 
   XOR InputA  $[N - 1 : 0]$  + InputA  $[N - 1]$   $\triangleright$  2's
   Complement of Input
8: InvertInput  $\leftarrow \{[N - 1]\{\text{OutputTwos} [N - 2] \}\}$ 
   XOR OutputTwos  $[N - 2 : 0]$ 
9: ZC  $\leftarrow$  Leading Zero Detector (InvertInput)
10: temp  $[N - 4 : 0] \leftarrow$  OutputTwos  $[N - 4 : 0] \ll (ZC - 1)$ 
11: FractionBitsTemp  $\leftarrow \{2'b01, \text{temp}, 3'b0\}$ 
12: ExpBits  $\leftarrow$  temp  $[N - 4 : N - 3 - es]$ 
13: Regime  $\leftarrow$  OutputTwos  $[N - 1] ? (ZC - 1) : -ZC$ 
14: FractionBits  $\leftarrow ((\text{InputA} [N - 1] \text{ XOR}$ 
   InputB  $[N - 1]) \& \text{InputA/B [MSB]})$ 
   ?  $-\text{FractionBitsTemp} : \text{FractionBitsTemp}$ 
```

Algorithm 2 Calculating Scaling Factor

```
1: ShiftRegimeA  $\leftarrow$  RegimeA  $\ll es$ 
2: ShiftRegimeB  $\leftarrow$  RegimeB  $\ll es$ 
3: SFA  $\leftarrow$  ShiftRegimeA + ExpBitsA
4: SFB  $\leftarrow$  ShiftRegimeB + ExpBitsB
5: ValuetoSift  $\leftarrow |SFA - SFB|$ 
6: if InputA > InputB then
7:   GreaterFractionBits  $\leftarrow$  FractionBitsA
8:   SmallerFractionBits  $\leftarrow$  FractionBitsB
9:   GreatestScalingFactor  $\leftarrow$  SFA
10: else
11:   GreaterFractionBits  $\leftarrow$  FractionBitsB
12:   SmallerFractionBits  $\leftarrow$  FractionBitsA
13:   GreatestScalingFactor  $\leftarrow$  SFB
14: end if
```

POSIT ADDER ALGORITHM (Currently working)

Algorithm 3 Adding Significand

```
1: SignificandAdd  $\leftarrow$  (GreaterFractionBits +  
   (SmallerFractionBits  $\gg$  Valuetoshift))  
2: AnsIsZero  $\leftarrow$   $\sim$ | SignificandAdd  
3: AddOperation  $\leftarrow$  InputA[N - 1] OR InputB[N - 1]  
4:  
5: if SignificandAdd[MSB] = 0 & AddOperation=0/1 then  
6:   ToNormalizedFraction  $\leftarrow$  SignificandAdd  
7:   Shift = 0  
8: else if SignificandAdd[MSB] = 1 & AddOperation=0 then  
9:   ToNormalizedFraction  $\leftarrow$  SignificandAdd  $\gg$  1  
10:  Shift = 1  
11: else if SignificandAdd[MSB] = 1 & AddOperation=1 then  
12:   ToNormalizedFraction  $\leftarrow$  -SignificandAdd  
13:   ZC1  $\leftarrow$  Leading Zero Detector(ToNormalizedFraction)  
14:   NormalizedFraction  $\leftarrow$  ToNormalizedFraction  $\ll$  ZC1  
15: end if  
16: Adjusted Scaling factor  $\leftarrow$   
   GreatestScalingFactor+Shift-ZC1
```

Algorithm 4 Decoding Regime and Exponent Value, Encoding and Rounding

```
1: Absolute SF  $\leftarrow$  Adjusted Scaling Factor [MSB]  
   ? -Adjusted Scaling Factor : Adjusted Scaling Factor  
2: if Adjusted Scaling Factor[MSB] = 0 then  
3:   ExpBits  $\leftarrow$  Absolute SF [es - 1 : 0]  
4:   RegimeAns  $\leftarrow$  Absolute SF  $\gg$  es  
5: else  
6:   ExpBits  $\leftarrow$  Adjusted Scaling Factor [es - 1 : 0]  
7:   RegimeAns1  $\leftarrow$  Adjusted Scaling Factor  $\gg$  es  
8:   RegimeAns  $\leftarrow$  -RegimeAns1  
9: end if  
10: TempAns1  $\leftarrow$  {2'b10, ExpBits, NormalizedFraction,  
   {[maxregime]{1'b0}} }  $\triangleright$  Packing Stage 1  
11: TempAns2  $\leftarrow$  {2'b01, ExpBits, NormalizedFraction,  
   {[maxregime]{1'b0}} }  
12: if RegimeAns={log(N){1'b1}} then  
13:   Shiftnegexp  $\leftarrow$  RegimeAns  
14: else  
15:   Shiftnegexp  $\leftarrow$  RegimeAns-1  
16: end if  
17: if Adjusted Scaling Factor[MSB] = 0 then  
18:   TempAns  $\leftarrow$  TempAns1  $\gg$  RegimeAns  
19: else  
20:   TempAns  $\leftarrow$  TempAns2  $\gg$  Shiftnegexp  
21: end if  
22: Guard bit  $\leftarrow$  Extract from TempAns  $\triangleright$  Packing Stage 2  
23: Round  $\leftarrow$  Extract from TempAns  
24: Sticky  $\leftarrow$  Extract from TempAns  
25: LSB  $\leftarrow$  Extract from TempAns  
26: checkround  $\leftarrow$  ((LSB & Guard) | (Guard & Round)) |  
   (Guard & Round | Sticky)  
27: IntermediateAns  $\leftarrow$  {1'b0, TempAns}+checkround  
28: if AnsIsZero=0 then  
29:   FinalAns  $\leftarrow$  0  
30: else if take2sans=1 then  
   FinalAns  $\leftarrow$  -IntermediateAns  
31: else  
   FinalAns  $\leftarrow$  IntermediateAns  
32: end if  
33: AdditionResult  $\leftarrow$  Special ? inf : FinalAns
```


NEXT STEPS

- Need to Understand the algorithm of Posit Multiplier
- Implementation of Posit adder & Posit multiplier using Verilog HDL
- To implement Posit Adder the following modules has to design -
- Top-module which takes N (posit word size) and es (posit exponent size), Posit data extract sub-module, Dynamic right shifter sub-module, Dynamic left shifter sub-module, Leading-One-Detector sub-module, Leading-Zero-Detector sub-module, Two N-bit integer adder sub-module, Two N-bit integer subtract sub-module, Integer Adder sub-module for mantissa overflow addition, Test-bench module.
- Similarly, Posit Multiplier implementation

REFERENCES

- Manish Kumar Jaiswal and Hayden K.-H. So, "Architecture Generator for Type-3 Unum Posit Adder/Subtractor", IEEE International Symposium on Circuits and Systems (ISCAS 2018), pp. 1-5, Florence, Italy, May 2018.
- Manish Kumar Jaiswal and Hayden K.-H. So, "Universal Number Posit Arithmetic Generator on FPGA", Design Automation and Test (DATE 2018), pp. 1159-1162, Dresden, Germany, Mar 2018.
- <https://ieeexplore.ieee.org/document/10089906> - By Gopal Sir
- <https://github.com/manish-kj/Posit-HDL-Arithmetic> - For coding Reference

THANK YOU