

# BitSys: Bitwise Systolic Array Architecture for Multi-precision Quantized Hardware Accelerators

Yuhao Liu<sup>1,3</sup> , Student Member, IEEE, Salim Ullah<sup>2</sup> , Akash Kumar<sup>2,3</sup> , Senior Member, IEEE

<sup>1</sup>Dresden University of Technology, Germany <sup>2</sup>Ruhr University Bochum, Germany

<sup>3</sup>Center for Scalable Data Analytics and Artificial Intelligence (ScaDS.AI Dresden/Leipzig), Germany

Email: {yuhao.liu1, salim.ullah}@tu-dresden.de, akash.kumar@rub.de

## I. EXTENDED ABSTRACT

*Quantized Neural Networks* (QNN) have been widely applied in hardware accelerator designs for edge. Because lower precision in quantization leads to higher accuracy loss, the mixed-precision scheme has been explored by using different precision in different layers to trade off resource consumption and inference accuracy. Because regular multiplier designs do not support the reconfiguration for multi-precision, we explored a runtime reconfigurable multi-precision bitwise systolic array design, *BitSys*, for mixed-precision multiplication in QNN accelerators. The popular design in previous works, such as [1], divides the inputs of multipliers as two parts for four sub-multipliers and preset left shifting, achieving the reconfigurable multiplication by disabling two of the sub-multipliers. Our design is inspired by the *Bitshifter* architecture from the works of Liu et al. [2, 1]. We convert the  $n \times n$ -bit multiplication as  $A \times B = \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} 2^{i+j} a_i b_j$ . As shown in Figure 1, partial products  $P_{i+j}$  is the sum of sub-partial products  $a_i b_j$  with left shifting value,  $i + j$ . The sub-partial product masks select the desired  $a_i b_j$  to configure the multi-channel according to the precision. One mask square represents one channel. Therefore, we can implement a bitwise systolic array as Figure 2 (left). The sub-partial product computation and mask are fused in one LUT primitive as one processing element in Figure 2 (right up). The sums of elements, considered the sign-bits, in the diagonal with the same left shifting value shown in Figure 2 (left) are the inputs,  $D_k$ , of the output-generate pipeline of Figure 2 (right down). Systolic array sequentially generates the  $D_k$ , and the final multi-precision output is the sum of all  $D_k$ . We implemented our *BitSys* multiplier as a systolic array for mixed-precision QNN acceleration. The comparison with the works of Liu et al. [1] is shown in Table I. Our systolic array accelerator consumes more hardware resources than the three single-layer accelerator instances of Liu et al. [1]. However, because of

TABLE I: Comparison of Resource Consumption and Latency

Design	Precision	LUT	FF	BRAM	Frequency	Latency/ $\mu$ s
Vivado IP [1]	8/8/8/8	24090	22175	135	150MHz	137.654
	1/2/4/8					131.059
Bitshifter [1]	1/2/4/8	42952	22486	138	125MHz	56.658
Multiplier Tree [1]	1/2/4/8	37020	22500	138	100MHz	69.27
<b>BitSys</b>	1/2/4/8	44468	64176	139.5	250MHz	36.741

8 bits x 8 bits Multi-Precision Multiplication										
$P_{i+j}$	$a_i b_j$ ( $i = 0 \rightarrow 7, j = 0 \rightarrow 7$ )								left bitshift	
	1bit	2bit	4bit	8bit	sum					
$P_0$	$a_0 b_0$				0	0	0	0	0	0
$P_1$	$a_1 b_0$	$a_0 b_1$			1	1	1	1	1	1
$P_2$	$a_2 b_0$	$a_1 b_1$	$a_0 b_2$		0+2	2	2	2	2	2
$P_3$	$a_3 b_0$	$a_2 b_1$	$a_1 b_2$	$a_0 b_3$	1+2	3	3	3	3	3
$P_4$	$a_4 b_0$	$a_3 b_1$	$a_2 b_2$	$a_1 b_3$	0+4	0+4	4	4	4	4
$P_5$	$a_5 b_0$	$a_4 b_1$	$a_3 b_2$	$a_2 b_3$	1+4	1+4	5	5	5	5
$P_6$	$a_6 b_0$	$a_5 b_1$	$a_4 b_2$	$a_3 b_3$	0+6	2+4	6	6	6	6
$P_7$	$a_7 b_0$	$a_6 b_1$	$a_5 b_2$	$a_4 b_3$	1+6	3+4	7	7	7	7
$P_8$	$a_7 b_1$	$a_6 b_2$	$a_5 b_3$	$a_4 b_4$	0+8	0+8	0+8	8	8	8
$P_9$	$a_7 b_2$	$a_6 b_3$	$a_5 b_4$	$a_4 b_5$	1+8	1+8	1+8	9	9	9
$P_{10}$	$a_7 b_3$	$a_6 b_4$	$a_5 b_5$	$a_4 b_6$	0+10	2+8	2+8	10	10	10
$P_{11}$	$a_7 b_4$	$a_6 b_5$	$a_5 b_6$	$a_4 b_7$	1+10	3+8	3+8	11	11	11
$P_{12}$	$a_7 b_5$	$a_6 b_6$	$a_5 b_7$		0+12	0+12	4+8	12	12	12
$P_{13}$	$a_7 b_6$	$a_6 b_7$			1+12	1+12	5+8	13	13	13
$P_{14}$	$a_7 b_7$				0+14	2+12	6+8	14	14	14

1-bit Sub-Partial Products Mask								2-bit Sub-Partial Products Mask								4-bit Sub-Partial Products Mask									
$a_7$	$a_6$	$a_5$	$a_4$	$a_3$	$a_2$	$a_1$	$a_0$	$a_7$	$a_6$	$a_5$	$a_4$	$a_3$	$a_2$	$a_1$	$a_0$	$a_7$	$a_6$	$a_5$	$a_4$	$a_3$	$a_2$	$a_1$	$a_0$		
$b_7$	1	0	0	0	0	0	0	$b_7$	1	1	0	0	0	0	0	$b_7$	1	1	1	1	0	0	0	0	
$b_6$	0	1	0	0	0	0	0	$b_6$	1	1	0	0	0	0	0	$b_6$	1	1	1	1	0	0	0	0	
$b_5$	0	0	1	0	0	0	0	$b_5$	0	1	1	0	0	0	0	$b_5$	1	1	1	1	0	0	0	0	
$b_4$	0	0	0	1	0	0	0	$b_4$	0	0	1	1	0	0	0	$b_4$	1	1	1	1	0	0	0	0	
$b_3$	0	0	0	0	1	0	0	$b_3$	0	0	0	0	1	1	0	$b_3$	0	0	0	0	1	1	1	1	
$b_2$	0	0	0	0	0	1	0	$b_2$	0	0	0	0	0	1	1	0	$b_2$	0	0	0	0	1	1	1	1
$b_1$	0	0	0	0	0	0	1	$b_1$	0	0	0	0	0	0	1	1	$b_1$	0	0	0	0	1	1	1	1
$b_0$	0	0	0	0	0	0	0	$b_0$	0	0	0	0	0	0	0	1	$b_0$	0	0	0	0	1	1	1	1

Fig. 1: Tree Structure of Multi-precision Accumulator

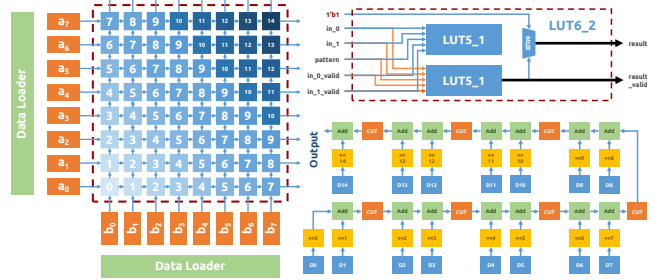


Fig. 2: Tree Structure of Multi-precision Accumulator

our bitwise processing design, *BitSys* instance can support 250MHz clock frequency because of the low critical path delay of our multiplier and achieves 188.5-274.7% speed-up in the evaluation of one four-layer 1/2/4/8-bit mixed-precision quantized MLP. Furthermore, our design does not change the input/out width when configuring to different precision, which can be easily integrated into existing accelerator designs.

## REFERENCES

- Yuhao Liu et al. "High Flexibility Designs of Quantized Runtime Reconfigurable Multi-Precision Multipliers". In: *IEEE Embedded Systems Letters* (2023), pp. 1–1.
- Negar Neda et al. "Multi-Precision Deep Neural Network Acceleration on FPGAs". In: *2022 27th Asia and South Pacific Design Automation Conference (ASP-DAC)*. IEEE, 2022, pp. 454–459.