

# Logic Design of Neural Networks for High-Throughput and Low-Power Applications

Kangwei Xu<sup>1</sup>, Grace Li Zhang<sup>2</sup>, Ulf Schlichtmann<sup>1</sup>, Bing Li<sup>1</sup>

<sup>1</sup>*Chair of Electronic Design Automation, Technical University of Munich (TUM), Munich, Germany*

<sup>2</sup>*Hardware for Artificial Intelligence Group, Technical University of Darmstadt, Darmstadt, Germany*

Email: kangwei.xu@tum.de, grace.zhang@tu-darmstadt.de, ulf.schlichtmann@tum.de, b.li@tum.de

**Abstract—** Neural networks (NNs) have been successfully deployed in various fields. In NNs, a large number of multiply-accumulate (MAC) operations need to be performed. Most existing digital hardware platforms rely on parallel MAC units to accelerate these MAC operations. However, under a given area constraint, the number of MAC units in such platforms is limited, so MAC units have to be reused to perform MAC operations in a neural network. Accordingly, the throughput in generating classification results is not high, which prevents the application of traditional hardware platforms in extreme-throughput scenarios. Besides, the power consumption of such platforms is also high, mainly due to data movement. To overcome this challenge, in this paper, we propose to flatten and implement all the operations at neurons, e.g., MAC and ReLU, in a neural network with their corresponding logic circuits. To improve the throughput and reduce the power consumption of such logic designs, the weight values are embedded into the MAC units to simplify the logic, which can reduce the delay of the MAC units and the power consumption incurred by weight movement. The retiming technique is further used to improve the throughput of the logic circuits for neural networks. In addition, we propose a hardware-aware training method to reduce the area of logic designs of neural networks. Experimental results demonstrate that the proposed logic designs can achieve high throughput and low power consumption for several high-throughput applications.

## I. INTRODUCTION

Neural networks (NNs) have been successfully applied in various fields, e.g., pattern recognition and natural language processing. In NNs, a large number of multiply-accumulate (MAC) operations need to be performed. Traditional digital hardware platforms such as GPU and TPU use parallel MAC units consisting of multipliers and adders to perform such MAC operations. Due to area constraints, the number of MAC units is limited. Accordingly, MAC units on such platforms have to be reused to implement all the MAC operations in a neural network. Therefore, the throughput of generating classification results on such platforms is not high, which prevents their adoption in extremely high-throughput applications, such as signal compensation in optical fiber communications [1], data collection from physics experiments [2] and malicious packet filtering for network detection [2]. In addition, large power consumption is another issue using such platforms to accelerate NNs, mainly resulting from data movement, e.g., loading weights from external DRAM to MAC units [3]. MAC operations also cause a part of the power consumption. Thus it remains challenging

to perform high-throughput tasks on those resource-constrained platforms requiring low power consumption.

Various methods, from the software to the hardware levels, have been proposed to address the throughput and power consumption issues in traditional digital hardware platforms. On the software level, pruning [4] and quantization [5] of NNs have been explored to reduce the number of MAC operations and the complexity of performing MAC operations, respectively. In addition, different dataflows, e.g., weight stationary [6], output stationary [7], and row stationary [8], have been proposed to reduce data movement and power consumption. On the hardware level, traditional digital hardware platforms using parallel MAC units are modified, e.g., by inserting multiplexers [9] to improve the throughput, and power/clock gating [10] [11] to reduce power consumption.

Another perspective to improve throughput and reduce power consumption in accelerating a neural network is to convert it into a logic circuit or a look-up table (LUT)-based design, where weights are embedded. For example, LogicNets [2] quantizes the inputs and outputs of neurons with low bit widths and implements such neurons with LUTs, which are subsequently deployed on FPGAs. NullaNets [12] train a neural network to produce binary activations and treat the operations at a neuron as multi-input multi-output Boolean functions. Such Boolean functions of neurons are further considered as truth tables and synthesized with logic synthesis tools.

Although converting a neural network into a logic/LUT design is promising to improve throughput and reduce power consumption in accelerating this neural network, the existing methods either incur a large number of LUTs for a high inference accuracy or cannot guarantee the feasibility of logic synthesis of truth tables of neurons due to their complexity. To address this challenge, in this paper, we propose to directly flatten and implement all the operations in a neural network with their corresponding logic circuits and embed weights into such circuits. The key contributions of this paper are summarized as follows.

- All operations including MAC and activation functions in a neural network are flattened and implemented with logic circuits. Flip-flops are inserted at the end of neurons in each layer to improve the throughput of such circuits.
- In implementing MAC operations with logic circuits, pre-determined weights after training are used to simplify the logic of MAC units to reduce their delay, power and area. Since weights are not required to be moved from external

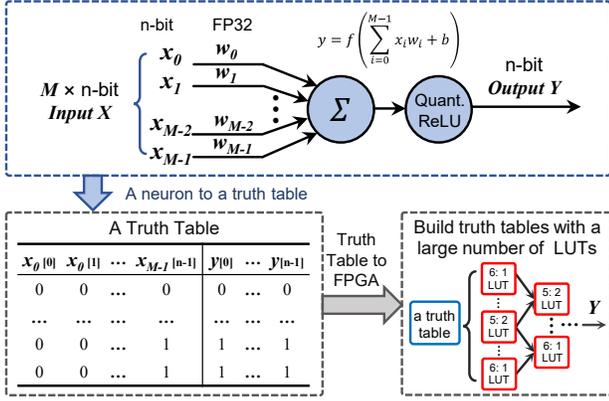


Fig. 1. The concept of implementing a neuron with LUTs in LogicNets [2].

- DRAM, power consumption can be reduced significantly.
- After the MAC units are simplified with weights in a layer, some logic is shared between the layers. Therefore, the whole neural network circuit is further simplified with EDA tools to reduce area and delay.
  - Different weight values affect the logic complexity of the resulting simplified MAC units. Accordingly, the traditional training is adjusted to select those weight values leading to smaller circuit sizes.
  - We present comprehensive evaluations on three high-throughput tasks, demonstrating that the proposed method can achieve high throughput and low power consumption. Furthermore, the proposed method outperforms the state-of-art approaches by achieving an average 2.19% increase in inference accuracy and an average 13.26% increase in throughput while reducing area overheads by an average of 38.76%.

The rest of this paper is organized as follows. Section II gives the background and motivation of this work. Section III explains the details of the proposed method. The experimental results are given in Section IV. Section V concludes the paper.

## II. BACKGROUND AND MOTIVATION

Neural networks have multiple layers of neurons. Synapses connect neurons with individual weights. The weights of a layer form a weight matrix. Computing a layer in an NN requires the multiplication of input data and the weight matrix, incurring a large number of MAC operations. The results of MAC operations are further processed by activation functions, e.g., ReLU.

Traditional digital hardware platforms such as GPU and TPU adopt parallel MAC units to accelerate such MAC operations. Due to the area constraint of such platforms, the number of MAC units is limited, so those MAC units have to be reused to perform all the MAC operations. Therefore, the throughput of generating classification results is low and not well suited for extreme-throughput applications. In addition, such platforms

suffer from large power consumption due to data movement and MAC computation, limiting their application in resource-constrained platforms, e.g., edge devices.

Various techniques have been proposed to address throughput and power consumption issues in accelerating NNs. One of the promising techniques is to convert NNs into LUT-based designs and logic circuits. We will explain their basic concepts as follows.

1) *LogicNets*: [2] proposes to map quantized neurons in a neural network to LUTs and implement such LUTs on FPGAs. The concept of LogicNets is shown in Fig. 1, where the  $M$  inputs and one output of a neuron are quantized to  $n$  bits. To convert this neuron into a LUT, the function of neurons can be expressed by a truth table that enumerates all the combinations of input bits (i.e., the fan-in of this neuron) with their corresponding outputs. For an entry in the truth table, the output can be evaluated according to the input bits and predetermined weights. The weights themselves do not appear in the truth table directly. Once the truth table is established, it can be mapped to LUTs and implemented on FPGAs.

Since neurons in NNs may have many inputs and their quantization bits should be large enough to maintain a high inference accuracy. As the truth table becomes large, the above mapping may lead to a large number of LUTs. This is because the hardware cost of implementing truth tables of neurons with LUTs grows exponentially with neuron fan-in. For example, implementing a 32:1 truth table requires about a hundred million 6:1 LUTs, which is much larger than even the largest FPGAs available today and makes it impractical for direct mapping to LUTs [2]. Although extreme pruning and quantization can limit the neuron fan-in, this degrades the inference accuracy.

2) *NullaNets*: [12] converts the truth table of a neuron into a logic circuit with logic synthesis. To reduce the complexity in logic synthesis, they only evaluate the output values for input combinations in the training dataset, and the output values for the remaining input combinations are set as “don’t care”. The resulting synthesized logic circuit may not generate the correct output for those input data in the test dataset, degrading the inference accuracy. Besides, even though only partial input combinations are considered in logic synthesis, the logic complexity might still be high, preventing the feasibility of logic synthesis of the truth table of a neuron. Contrary to the previous work in converting NNs into LUTs and logic circuits, the proposed method directly flattens all the operations in a neural network into their corresponding logic circuits with weights embedded in such circuits.

## III. LOGIC DESIGN OF NEURAL NETWORKS

In this section, we first introduce the logic implementations of neural networks. Retiming is then used to improve the throughput of the logic designs of NNs. Afterwards, a hardware-aware training technique is proposed to reduce the area overheads of such logic designs.

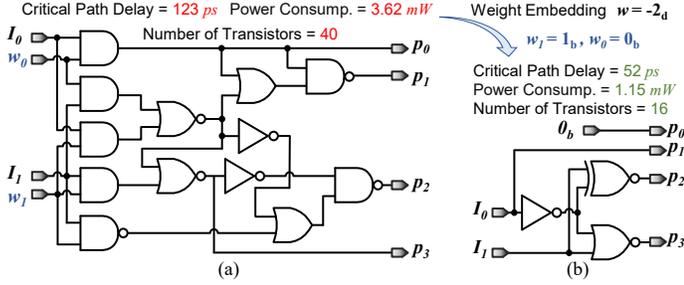


Fig. 2. Logic circuit of 2-bit signed multiplier. (a) The original circuit; (b) The logic circuit simplified with a fixed quantization weight (decimal: -2, binary: 10).

### A. Logic Implementations of NNs with Fixed Weights

To implement all the operations in NNs with logic circuits, we first use quantization-aware training to train a neural network while maintaining the inference accuracy. 8-bit quantization of weights and input activations is used during this process [13]. In addition, unstructured pruning [14] is further used to remove unimportant weights and reduce computational costs. Afterwards, the neural network is fine-tuned to improve the inference accuracy.

After training, the MAC operations in a neural network can be directly implemented with MAC units and the logic realizing the activation functions. The fixed weights after training are used to simplify the MAC operations at neurons. The simplified MAC units are appended with the logic circuit implementing the activation function at a neuron. All the neurons in this network are processed similarly. The resulting logic circuits are concatenated to generate the complete logic of the neural network. Then, flip-flops are inserted at the end of each layer to synchronize data propagation, after which logic redundancy within and across layers is removed by EDA tools. In the following paragraphs, we will introduce each part of the logic design of a neural network.

1) *Multiplier*: To reduce the delay, power and area of multipliers, the fixed weights after training are used to simplify the logic of the multipliers. Fig. 2 illustrates the logic simplification of a 2-bit signed multiplier with the fixed quantized weight  $-2_d$  or  $10_b$ . After this simplification, the delay of the multiplier circuit is reduced by 57.72%, the power consumption is reduced by 68.23%, and the number of transistors is reduced by 60%. All the multipliers at a neuron are processed this way.

2) *Adder*: The simplified multipliers at a neuron will be appended with an adder realizing the addition operation. Due to the weight embedded in multipliers, the resulting logic circuit of the MAC operation can be further simplified. Fig. 3 illustrates the comparison of the MAC unit circuit before and after weight embedding, where the multipliers and the adder are 2-bit and 4-bit, respectively. With weight embedding, the delay of the MAC unit is reduced by 70.07%, the power consumption is reduced by 71.64%, and the number of transistors is reduced by 65%.

In this work, we use an 8-bit neural network to maintain the

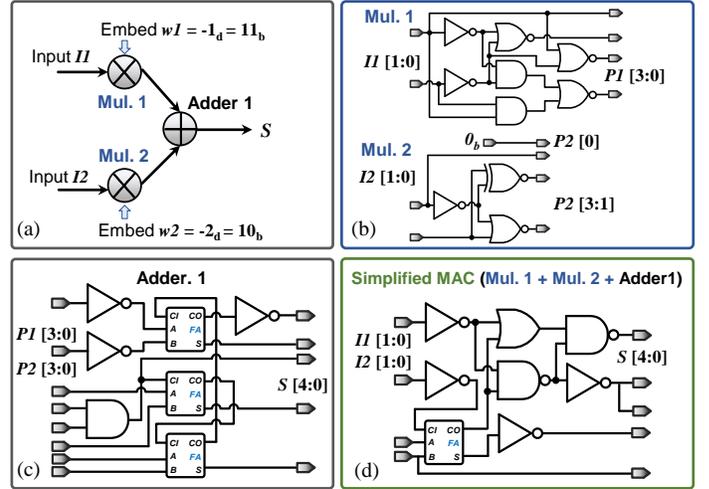


Fig. 3. (a) MAC operations at a neuron; (b) 2-bit signed multipliers simplified with the fixed quantized weights; (c) 4-bit signed adder circuit before logic simplification, where FA is a 1-bit full adder; (d) Circuit of the simplified MAC unit.

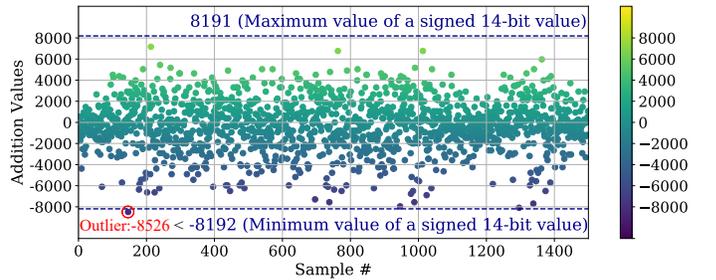


Fig. 4. Distribution of 1500 samples of addition values from the adder output in a neuron for the OFC task. This example illustrates that after removing the outlier  $-8256_d$  (15-bit), the output bit width of the adder needs only 14-bit.

inference accuracy of the hardware implementation. However, the quantization bit of the adder of a MAC unit is larger than 8-bit and increases with the increasing number of inputs at a neuron. In general, the quantization bit of the MAC operation result should be reduced to 8-bit before this result enters the subsequent layer [13]. Two techniques are proposed to maintain the inference accuracy during this process.

First, all the addition results at neurons in NNs are profiled to examine the actual bit width to represent such addition results. Fig. 4 presents an example of the distribution of addition results according to the training dataset. By removing the outliers, the maximum values of the remaining results are used to define the output bit width of the adder. Second, a requantizer is used to convert the high-bit value to 8-bit by multiplying it with a scaling factor [13]. By embedding the scaling factor into the requantizer, the delay, power and area of the requantizer circuit can be reduced.

3) *The Circuit for Implementing Activation Function*: An activation function outputs the corresponding activation state by thresholding the input to determine whether the neuron should

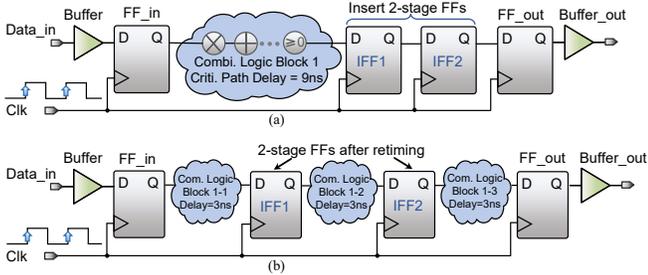


Fig. 5. (a) Insertion of two-stage flip-flops after the original ReLU circuit; (b) Circuit after retiming.

be activated or not activated. The activation function ReLU generates 0 if the input is smaller than 0 and remains the original value if it is larger than 0. The circuit for performing ReLU is generated by describing its function with hardware description languages such as Verilog and synthesizing it by the EDA tool.

The logic circuit of each neuron is generated with the techniques described above. The resulting circuits for all neurons are concatenated together to produce the complete logic design of the neural network. The logic circuits of neurons might share some common logic since they result from the simplification of a MAC unit. To remove the logic redundancy, the complete circuit of a neural network will be optimized by EDA tools.

### B. Retiming with Cascaded Flip-Flops

To synchronize data propagation in each layer, flip-flops are inserted at the end of the logic circuit implementing the activation function. To further improve the throughput of the logic circuit of a neural network, cascaded flip-flops are inserted into the original circuits, and retiming technique is then deployed to reduce the maximum delay between flip-flop stages. As a retiming example at a neuron shown in Fig. 5(a), the combinational logic block 1 has a critical path delay of 9, limiting the whole circuit’s performance. Assume the clock-to-q delay and the setup time of a flip-flop are 3 and 1, respectively. The minimum clock period of this circuit is equal to 13. In this example, two-stage flip-flops (IFF1 and IFF2) are inserted after the circuit implementation of ReLU. Afterwards, retiming is realized by EDA tools automatically. As shown in Fig. 5(b), with the retiming technique, the performance can be optimized by reallocating IFF1 and IFF2 in the combinational logic block 1, resulting in a minimum clock period equal to 7, which reduces the clock period by 46.15%.

### C. Hardware-Aware Training

Different weight values affect the logic complexity of the resulting simplified multiplier. For example, as shown in Fig 6, the quantized 8-bit weight ‘107’ corresponds to a larger multiplier area of 111 units, while the quantized 8-bit weight ‘-16’ leads to an area overhead of only 16 units. To take advantage of this property, we propose to train the neural network with selected weights that lead to a smaller multiplier area. The weight selection and the modified training are explained as follows.

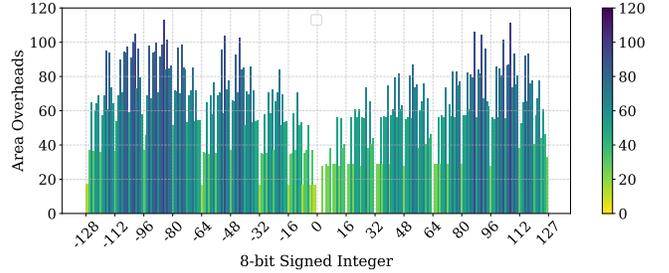


Fig. 6. Area of multipliers simplified with 8-bit quantized weights.

1) *Weight selection*: We first rank the weight values according to the area of the resulting simplified multipliers. Then, we select the top  $n$  weights that lead to the smallest multiplier area. In the experiments,  $n$  was set to 40. The neural network is trained only with such weight values, and the validation accuracy is verified. If the validation accuracy is much lower than that of the original training, more weight values, e.g., 50, are selected and used to train the neural network. In each iteration, 10 more weight values leading to the small area will be added to the previously selected weight value set and used to train the network. The iteration continues until the validation accuracy recovers almost to that of the original training.

2) *Training*: During training, the weights are forced to take the selected values in the forward propagation. In the backward propagation, the straight-through estimator is applied to skip the selection operation. In each epoch, all weights in the model are traversed and replaced with the closest value from the selected weight set. The training continues until the loss function converges or a given number of epochs is finished.

## IV. EXPERIMENTAL RESULTS

In this section, we demonstrate the results of the proposed method in terms of accuracy, throughput, power consumption and area overheads in 3 different extreme-throughput applications:

1) *Optical Fiber Communications (OFC)*: Transmission of optical signals in fibers suffers from chromatic dispersion (CD). Neural networks are used as nonlinear equalizers to compensate CD, where the input is the optical signal affected by CD, and the output is the compensated optical signal. We use the formulation from [1] for OFC as a 21-input 1-output prediction task.

2) *Jet Substructure Classification (JSC)*: In large-scale physics experiments, terabytes of instrumentation data are generated every second. Neural networks with 16 inputs and 5 outputs are employed to filter out the most interesting results [?].

3) *Network Intrusion Detection (NID)*: The neural networks can identify malicious network packets to strengthen network security. For this task, we use the UNSW-NB15 dataset [2], which consists of the packets labeled as either bad (0) or normal (1) with a total of 593 input features.

During training, the Adam optimizer was used and the step decay learning rate schedule was set starting from 0.001. The mini-batch size was set to 1024. The quantization-aware training

TABLE I  
THE PERFORMANCE RESULT WITH THE PROPOSED METHOD.

Network	Neurons / Layer	BER(OFC) or Acc.(JSC/NID) <sup>1</sup>			Max. Frequency (MHz)			Power Consum. (mW)			No. of Weight	
		Float	Pruned	Prop.	Base.	Retim.	Improv.	Base.	Prop.	Drop	Base.	Prop.
OFC-A	21, 40, 25	$4.34e^{-4}$	$5.43e^{-4}$	$5.87e^{-4}$	338	568	68.05%	543	88	83.79%	256	70
OFC-B	21, 50, 25	$3.26e^{-4}$	$3.91e^{-4}$	$4.13e^{-4}$	322	510	58.39%	622	102	83.60%	256	80
OFC-C	21, 50, 50	$2.39e^{-4}$	$2.61e^{-4}$	$3.26e^{-4}$	306	463	51.31%	986	122	87.63%	256	80
JSC-A	16, 64, 16, 16, 8	73.94%	73.69%	73.65%	364	564	54.95%	763	181	76.28%	256	80
JSC-B	16, 64, 32, 32, 32	74.42%	74.20%	74.11%	336	442	31.55%	1270	299	76.46%	256	100
JSC-C	16, 64, 48, 48, 32	74.77%	74.50%	74.44%	305	422	38.36%	1783	432	75.77%	256	110
NID-A	593, 20	90.89%	90.63%	90.34%	458	538	17.47%	349	120	65.62%	256	100
NID-B	593, 20, 20	91.92%	91.85%	91.69%	360	495	37.50%	412	139	66.26%	256	100
NID-C	593, 25, 25	92.07%	91.89%	91.76%	352	481	36.65%	496	164	66.94%	256	120

<sup>1</sup> In general, the Bit Error Rate (BER) indicates the inference performance of the OFC task, while the accuracy (Acc.) indicates the inference performance of the JSC and NID tasks.

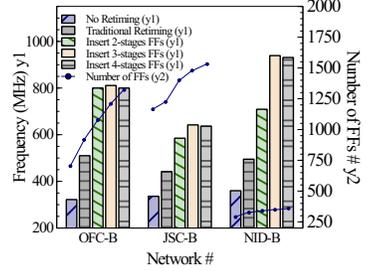
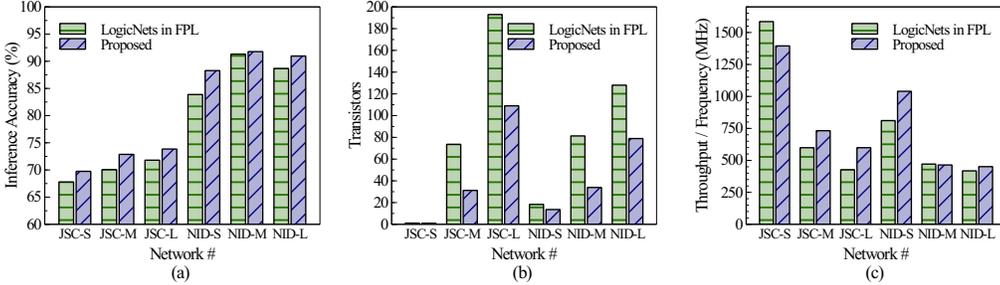


Fig. 7. (a) Comparison of inference accuracy with LogicNets; (b) Comparison in the number of transistors in the hardware implementation with LogicNets; (c) Comparison of maximum frequency with LogicNets.

Fig. 8. Comparison of max frequency before & after FF insertion and retiming.

stops after 300 epochs and the hardware-aware training stops after 100 epochs. Each experiment for a particular task was repeated 10 times and the average is reported. The neural network training was implemented with Pytorch on Nvidia Quadro RTX 6000 GPUs. The logic circuits in implementing NNs were synthesized with Synopsys Design Compiler using the Nangate 45nm open-cell library. The correctness of the circuits is verified by performing simulation and ensuring the same results as the original PyTorch network are returned.

Table I demonstrates the performance results. The first and second columns show the names of the neural networks and their structures. “A”, “B” and “C” represent different versions of neural networks. The third, fourth and fifth columns represent the inference performance of the original floating-point NNs, the traditional quantized and pruned NNs, and our proposed NNs, respectively. According to these columns, the proposed hardware-aware training method can maintain a low bit error rate (BER) in the application of OFC and high inference accuracy in the applications of JSC and NID. Compared with the baseline, where all the operations in a neural network are flattened without embedding fixed weights, the proposed method can achieve higher throughput and lower power consumption, as shown in the eighth and eleventh columns. The last two columns show the numbers of selected weight values in the hardware-aware training. Compared with the original 8-bit of 256 weight

values used for training, we train the NNs with only a small number of weight values (e.g., 70 in OFC-A) to reduce the size of their logic implementations.

To demonstrate the advantages of the proposed method, we compared the proposed method with LogicNets [2] in terms of inference accuracy, the number of transistors and maximum frequency. In this comparison, the network size, the pruning ratio, and the quantization bits of the activations were set as the same with LogicNets. The number of LUTs in LogicNets is equivalent to the number of transistors for area comparison. According to the results illustrated in Fig. 7, the proposed method achieves higher accuracy than LogicNets on all tasks with a much smaller number of transistors. In addition, in most cases, the maximum frequency with the proposed method is higher than that of LogicNets, even though LogicNets uses 16nm technology while the proposed method uses 45nm technology.

To balance the maximum frequency and the number of inserted flip-flops, we iteratively inserted more flip-flops after implementing the circuit ReLU and used retiming to improve the clock frequency. The results are illustrated in Fig. 8, where the histogram denotes the maximum frequency (MHz), and the line represents the number of flip-flops after retiming. According to this figure, the maximum frequency can be improved significantly by inserting more flip-flops initially and later remaining stable even with more flip-flops inserted into the logic circuit.

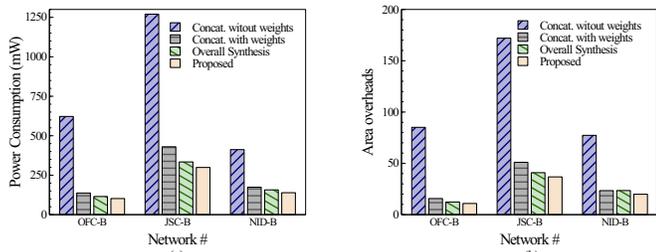


Fig. 9. Comparison of circuit performance. (a) Power consumption; (b) Area overheads.

To demonstrate the reduction in power and area with weight embedding and hardware-aware training, we compared the power consumption and area overheads before and after using such techniques. As shown in Fig. 9 (a), for the OFC task, when weights are embedded in the circuit, the power consumption of the synthesized circuit is 137mW. After considering the simplification between different logic, the power consumption is reduced to 115mW. With the hardware-aware training, the power consumption is further reduced to 102mW. According to this figure, weight embedding and hardware-aware training can significantly reduce power consumption and the corresponding area overheads.

In the proposed hardware-aware training, a certain number of weights resulting in a small multiplier area were selected to train the network. The inference accuracy with the proposed method can still be maintained, as illustrated in Fig. 10 (a). To demonstrate the tradeoff between the number of selected weights and inference performance, we selected different numbers of weights and used them to train the network. The results are illustrated in Fig. 10 (b-d). According to these figures, with the number of selected weights reduced, the BER in OFC increases and the inference accuracy in JSC and NID decreases. Besides, the area overheads reduce with decreasing number of selected weights.

## V. CONCLUSION

In this paper, we have proposed an efficient logic design of neural networks for extremely high-throughput applications. Instead of relying on parallel MAC units to accelerate these MAC operations, all the operations in NNs are implemented through logic circuits with embedded weights. Weight embedding not only reduces the delay of the logic circuit but also reduces power consumption. The retiming technique is used to further improve the circuit throughput. A hardware-aware training method is proposed to reduce the area of logic designs of NNs. Experimental results on three tasks demonstrate the proposed logic design can achieve a high throughput and low power consumption.

## ACKNOWLEDGEMENT

This work is funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) – Project-ID 504518248 and supported by TUM International Graduate School of Science and Engineering (IGSSE).

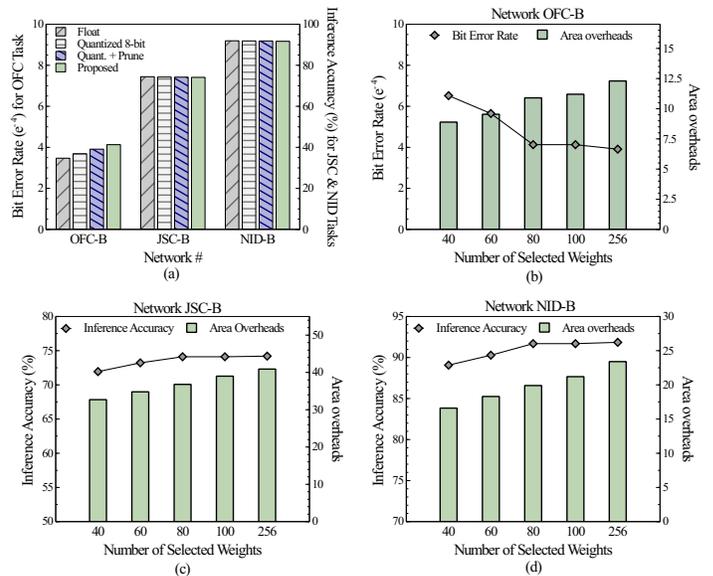


Fig. 10. (a) Accuracy Comparison with traditional quantization and pruning; (b)(c)(d) Tradeoff between accuracy, area overheads and the number of selected weights.

## REFERENCES

- [1] B. Liu, C. Blümm et al., "Area-Efficient Neural Network CD Equalizer for 4x 200Gb/s PAM4 CWD4 Systems," IEEE Optical Fiber Communications Conference and Exhibition (OFC), 2023.
- [2] Y. Umuroglu et al., "LogicNets: Co-Designed Neural Networks and Circuits for Extreme-Throughput Applications," IEEE/ACM International Conference on Field-Programmable Logic and Applications (FPL), 2020.
- [3] T.-J. Yang, Y.-H. Chen et al., "A method to estimate the energy consumption of deep neural networks," IEEE Asilomar Conference on Signals, Systems, and Computers (ACSSC), 2017.
- [4] R. Petri, Grace L. Zhang et al., "PowerPruning: Selecting Weights and Activations for Power-Efficient Neural Network Acceleration," IEEE/ACM Design Automation Conference (DAC), 2023.
- [5] W. Sun, Grace L. Zhang et al., "Class-based Quantization for Neural Networks," IEEE Design Automation and Test in Europe (DATE), 2023.
- [6] N. P. Jouppi, C. Young et al., "In-Datacenter Performance Analysis of a Tensor Processing Unit," IEEE/ACM International Symposium on Computer Architecture (ISCA), 2017.
- [7] T. Chen, Z. Du et al., "DianNao: A small-footprint high-throughput accelerator for ubiquitous machine-learning," ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), 2014.
- [8] Y.-H. Chen, T. Krishna et al., "Eyeriss: An Energy-Efficient Reconfigurable Accelerator for Deep Convolutional Neural Networks," IEEE International Solid-State Circuits Conference (ISSCC), 2016.
- [9] J. Zhang, H. Gu et al., "Hardware-Software Codesign of Weight Reshaping and Systolic Array Multiplexing for Efficient CNNs," IEEE Design Automation and Test in Europe (DATE), 2021.
- [10] N. D. Gundi et al., "EFFORT: Enhancing Energy Efficiency and Error Resilience of a Near-Threshold Tensor Processing Unit," IEEE/ACM Asia and South Pacific Design Automation Conference (ASP-DAC), 2020.
- [11] P. Pandey, N. D. Gundi et al., "UPTPU: Improving Energy Efficiency of a Tensor Processing Unit through Underutilization Based Power-Gating," IEEE/ACM Design Automation Conference (DAC), 2021.
- [12] M. Nazemi, G. Pasandiet al., "Energy-efficient, low-latency realization of neural networks through Boolean logic minimization," IEEE/ACM Asia and South Pacific Design Automation Conference (ASP-DAC), 2019.
- [13] B. Jacob, S. Kligys et al., "Quantization and Training of Neural Networks for Efficient Integer-Arithmetic-Only Inference," IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2018.
- [14] A. Renda, J. Frankle et al., "Comparing Rewinding and Fine-tuning in Neural Network Pruning," International Conference on Learning Representations (ICLR), 2020.