# Low Precision Floating-point Arithmetic for High Performance FPGA-based CNN Acceleration

Chen Wu, Mingyu Wang, Xinyuan Chu, Kun Wang, Lei He

*Abstract*—**Low precision data representation is important to reduce storage size and memory access for convolutional neural networks (CNNs). Yet, existing methods have two major limitations: (1) requiring re-training to maintain accuracy for *deep* CNNs, and (2) needing 16-bit floating-point or 8-bit fixed-point for a good accuracy. In this paper, we propose a low precision (8-bit) floating-point (LPFP) quantization method for FPGA-based acceleration to overcome the above limitations. Without any re-training, LPFP finds an optimal 8-bit data representation with negligible top-1/top-5 accuracy loss (within 0.5%/0.3% in our experiments, respectively, and significantly better than existing methods for *deep* CNNs). Furthermore, we implement one 8-bit LPFP multiplication by one 4-bit multiply-adder (MAC) and one 3-bit adder, and therefore implement *four* 8-bit LPFP multiplications using one DSP slice of Xilinx Kintex 7 family (KC705 in this paper) while one DSP can implement only *two* 8-bit fixed-point multiplications. Experiments on six typical CNNs for inference show that on average, we improve throughput by $64.5\times$ over Intel i9 CPU and by $1.5\times$ over existing FPGA accelerators. Particularly for VGG16 and YOLO, compared to six recent FPGA accelerators, we improve average throughput by 3.5× and 27.5× and improve average throughput per DSP by 4.1× and 5×, respectively. To the best of our knowledge, this is the first in-depth study to simplify one multiplication for CNN inference to one 4-bit MAC and implement four multiplications within one DSP while maintaining comparable accuracy without any re-training.**

*Index Terms*—**low precision floating-point, CNN, deep learning, FPGA processor, FPGA acceleration**

## I. INTRODUCTION

CONVOLUTIONAL neural networks (CNNs) have demonstrated a breakthrough in performance for a broad range of applications including object recognition [1], object detection [2] and speech recognition [3]. However, CNNs often have huge computation complexity. This motivates accelerating CNNs by CPU/GPU clusters [4], FPGAs [5] and ASICs [6]. Customized accelerators/processors on FPGAs have shown more promising throughput and power efficiency than traditional CPU/GPU clusters [7], [8].

Larger and deeper CNNs have been developed to improve performance for a broader range of scenarios. For example, the top-5 error for ImageNet [9] classification decreases from 17% to 2.9%. However, computation complexity and number of parameters increase dramatically as depicted in Figure 1. To be specific, the computation complexity of a feed-forward process of a 224×224 RGB image increases from

C. Wu, M. Wang, X. Chu, K. Wang and L. He are with the Department of Electrical and Computer Engineering, University of California, Los Angeles, CA, 90095, USA (e-mail: chenwu1989@ucla.edu, mingyuw@ucla.edu, xinyuansjtu@gmail.com, wangk@ee.ucla.edu, lhe@ee.ucla.edu).
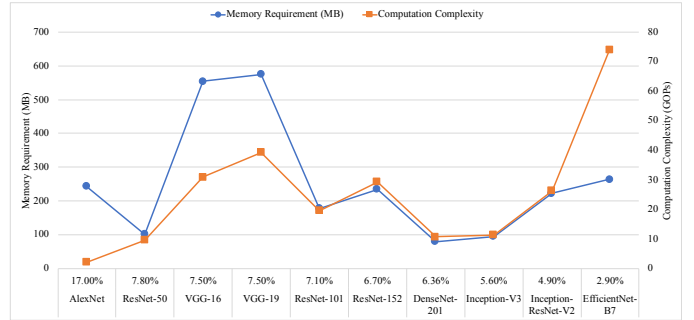
Fig. 1. Computation complexity and memory requirement with respect to different CNNs.

TABLE I
RESOURCE UTILIZATION OF MULTIPLIERS ON FPGA FOR DIFFERENT DATA REPRESENTATIONS. DSP: DIGITAL SIGNAL PROCESSING, LUT: LOOK-UP TABLE, FF: FLIP-FLOP. $M4E3$: 1-BIT SIGN, 4-BIT MANTISSA AND 3-BIT EXPONENT.

| Data Representation | DSP | LUT | FF |
|---|---|---|---|
| **one** 16-bit floating multiplication | 1 | 85 | 167 |
| **one** 16-bit fixed multiplication | 1 | 0 | 0 |
| **two** 8-bit fixed multiplications | 1 | 2 | 0 |
| **four** 8-bit floating ($M4E3$) multiplications | 1 | 20 | 27 |

2.27 GOP of AlexNet [10] in 2012 to 74 GOP of EfficientNet-B7 [11] in 2019. At the same time, the number of parameters stays large at 264 MB. Such great computation complexity makes it harder for general-purpose processor to meet the requirements of real-time applications. On the other hand, the great quantities of parameters lead to a big challenge for communication between off-chip and on-chip memories because of bandwidth constraints.

There are two types of research to reduce computation and parameter complexities for CNN inference. The first one is deep compression including weight pruning, weight quantization and compression storage [12], [13]. However, irregularity caused by deep compression degrades parallelism and hardware performance. Cambricon-S [14] alleviates irregularity in sparse neural networks through a software/hardware co-design approach to improve hardware performance. However, all the above accelerators need time-consuming re-training process to maintain accuracy.

The second type of research is more efficient data representation, also known as quantization for circuit implementation. [15] used 16-bit floating-point in contrast to 32-bit commonly used for computing. However, one 16-bit floating-point multiplier on FPGA needs 1 DSP, 85 LUTs and 167 FFs when using Xilinx floating-point IP [16] as shown in

Table I, leading to a low hardware efficiency. Since one 16-bit or smaller fixed-point multiplier can be fit into one DSP, both 16-bit [17], [18] and 8-bit [6], [19]–[21] fixed-point were employed to gain more hardware efficiency than 16-bit floating-point does. Another 8-bit arithmetic, called block floating-point (BFP), was also applied [22], [23], where a parameter has its own mantissa but shares a same exponent for one data block. [24] proposed a mixed data representation with floating-point for weights and fixed-point for activations (*e.g.,* outputs of a layer). [25] developed an 8-bit floating-point quantization scheme, which needs an extra inference batch to compensate for the quantization error. However, [24] and [25] did not present a circuit design for their approaches. While all aforementioned work has a good accuracy with re-training, more aggressive data representations such as binary [26], ternary [27], and mixed precision (2-bit activations and ternary weights) [28] may suffer from great accuracy loss even with time-consuming re-training.

In this paper, we first propose a low precision floating-point (LPFP) to quantize both weights and activations. During the quantization process, an optimal LPFP data format and the corresponding scale factor are decided for a workload of CNNs. Our proposed quantizer works for *deep* CNNs (more than 100 *convolutional/fully-connected* layers). On average, the top-1 accuracy loss is within 0.5%, while V-Quant [29] that works for such *deep* CNNs has a top-1 accuracy loss about 1% with fine-tuning. Then, we design a LPFP based FPGA processor to further improve the performance for CNN inference. We are able to implement four 8-bit floating-point multiplications within one DSP (see Table I). We experimented for inference of AlexNet, VGG16 [30], ResNet50/101/152 [31] and DenseNet201 [32] via Xilinx KC705. We can achieve an average throughput of 1100.4 GOPS (Giga-Operations Per Second), and it is 1.43 GOPS per DSP. Moreover, the average throughput for these networks is $64.5\times$ and $1.5\times$ over Intel i9 CPU and existing accelerators, respectively. Compared with six existing accelerators for VGG16 and YOLO, on average, our processor improves throughput by $3.5\times$ and $27.5\times$, while improveing per DSP throughput by $4.1\times$ and $5\times$, respectively.

While existing work needs re-training (including calibration [25] and fine-tuning [33]), our quantization method does not need any re-training to compensate for quantization error. Furthermore, to the best of our knowledge, this is the first work that can fit four 8-bit multiplications for inference in one DSP while maintaining comparable accuracy without any re-training.

## II. BACKGROUND AND MOTIVATION

### A. Background

*1) CNNs:* CNNs are used to classify or recognize objects by passing the inputs through multiple types of layers. In each layer, multiple neurons are constructed to process different inputs and pass the outputs to the next layer through connections, and the connections are used to store the weights for the network. Based on different processing procedures, the layers are typically divided into *convolutional, pooling, activation, normalization, fully-connected, residual* and *inception* layers.

TABLE II
CHARACTERISTICS OF CNN BENCHMARKS. GOP IS GIGA-OPERATIONS NEEDED BY ONE 224×224 RGB IMAGE.

| CNN | Type | Operations | Model Weights |
|---|---|---|---|
| AlexNet | *slim* | 2.27 GOP | 249.51 MB |
| VGG16 | *slim* | 30.94 GOP | 553.43 MB |
| ResNet50 | *medium* | 9.74 GOP | 46.05 MB |
| ResNet101 | *medium* | 19.70 GOP | 166.37 MB |
| ResNet152 | *deep* | 29.39 GOP | 229.39 MB |
| DenseNet201 | *deep* | 10.85 GOP | 68.63 MB |

Among them, *convolutional/fully-connected* layers consume most portions of computation while *fully-connected* layers require largest memory to store weights. According to this, we divide the size of CNNs into three categories with respect to the number of *convolutional/fully-connected* layers: 1) *slim* for less than 50 layers, 2) *medium* for 50 to 100 layers, and 3) *deep* for more than 100 layers, as shown in Table II where we report the detailed network information.

*2) Low Precision Floating-point:* Similar to the definition of 32-bit floating-point from the IEEE-754 standard [34], the binary representation of LPFP number comprises *sign, mantissa* and *exponent* in order. The decimal value of LPFP number is then calculated by:

$$V_{dec} = (-1)^S \times 1.M \times 2^{E-E_b}, \qquad (1)$$

where $V_{dec}$ is the value in decimal, $S, M$ and $E$ are all unsigned values and denote the *sign, mantissa* and *exponent*, respectively. For exponent bias $E_b$ in Eq. (1), it is introduced to both positive and negative exponents as

$$E_b = 2^{DW_E - 1} - 1, \qquad (2)$$

where $DW_E$ is the data width of $E$. Different from the IEEE Standard, data widths for $M$ and $E$ in this paper are not fixed. In later sections, we use the term $MaEb$ to indicate different combinations, where $a$ and $b$ indicate the bit width of $M$ and $E$, respectively. For example, $M3E4$ means the mantissa is 3 bits while the exponent is 4 bits.

There are three special definitions in IEEE-754 standard. The first is subnormal numbers when $E = 0$, then Eq. (1) is modified to:

$$V_{dec} = (-1)^S \times 0.M \times 2^{1-E_b}. \qquad (3)$$

Note that Infinity (Inf) and Not a Number (NaN) are the other two special cases, but are not used in our work. This is because our saturation scheme saturates large numbers to the maximal number, as illustrated in detail in Subsection III-A.

### B. Motivation

CNN accelerators with lower data width have significant improvements in terms of memory size, memory bandwidth and power efficiency. Due to the lack of floating-point arithmetic units in FPGA, researchers have used low precision fixed-point instead of floating-point. A 16-bit fixed-point quantization to find the best scale factor for each layer was proposed in [17]. However, this required time-consuming re-training to amend the weights to maintain accuracy. Furthermore, a model was

developed to quantitatively analyze the convolution loops and optimize design objectives such as memory access and latency [18]. However, it had an accuracy loss as large as 2%. A shared drawback for the above two approaches is the low per DSP throughput (0.279 GOPS/DSP for [17] and 0.472 GOPS/DSP for [18]) because of using 16-bit multiplication.

An 8-bit fixed-point accelerator was designed in [35] for embedded FPGAs, with a low per DSP throughput of 0.444 GOPS/DSP. DNNBuilder [36] aimed to automatically build high-performance DNN hardware accelerators for both cloud- and edge-FPGAs with 8-bit fixed-point quantization. It increased the per DSP throughput to 0.771 GOPS by better architecture exploration; however, its quantization method incurred 4.6% top-1 accuracy degradation without fine-tuning. FPGA accelerator with the aforementioned BFP arithmetic [23] had a per DSP throughput of 0.741 GOPS. However, only *slim* and *medium* CNNs were validated in their approach. In short, existing approaches cannot improve the per DSP throughput while maintaining comparable accuracy for all *slim, medium* and *deep* CNNs without using any re-training techniques.

## III. Low Precision Floating-point Quantization

In this section, we present the details of our proposed low precision floating-point (LPFP) quantization method, including the quantization process, data flow in processor and quantization results.

### A. Quantization Process

Our proposed LPFP quantization method is applied to both activations and weights. The quantization function is defined as follows:

$$V_{lfp} = quan(V_{fp32} \times 2^{sf}, MIN_{lfp}, MAX_{lfp}), \quad (4)$$

where $V_{lfp}$ and $V_{fp32}$ denote the decimal values represented by LPFP and traditional single floating-point format, respectively; $MIN_{lfp}$ and $MAX_{lfp}$ indicate the minimal and maximal numbers represented by LPFP, and $sf$ is the scaling factor which is used to better fit the data into the dynamic range of LPFP. The $quan$ function in Eq. (4) rounds the data to the nearest value with saturation considered, formulated as

$$quan(x, MIN, MAX) = \begin{cases} MIN & x <= MIN \\ MAX & x >= MAX \\ round(x) & \text{otherwise} \end{cases}, \quad (5)$$

where $MIN$ and $MAX$ are the minimal and maximal values, respectively.

The mean square error (MSE) of the values before and after quantization is used as the metric to evaluate the quantization error, illustrated as:

$$MSE = \frac{1}{N} \sum_{i=0}^{N} (V_{lfp}/2^{sf} - V_{fp32})^2, \quad (6)$$

where $N$ denotes the amount of data.

As illustrated from Eq. (4) to (6), MSE is influenced by the data format of LPFP and the scaling factor ($sf$). We will
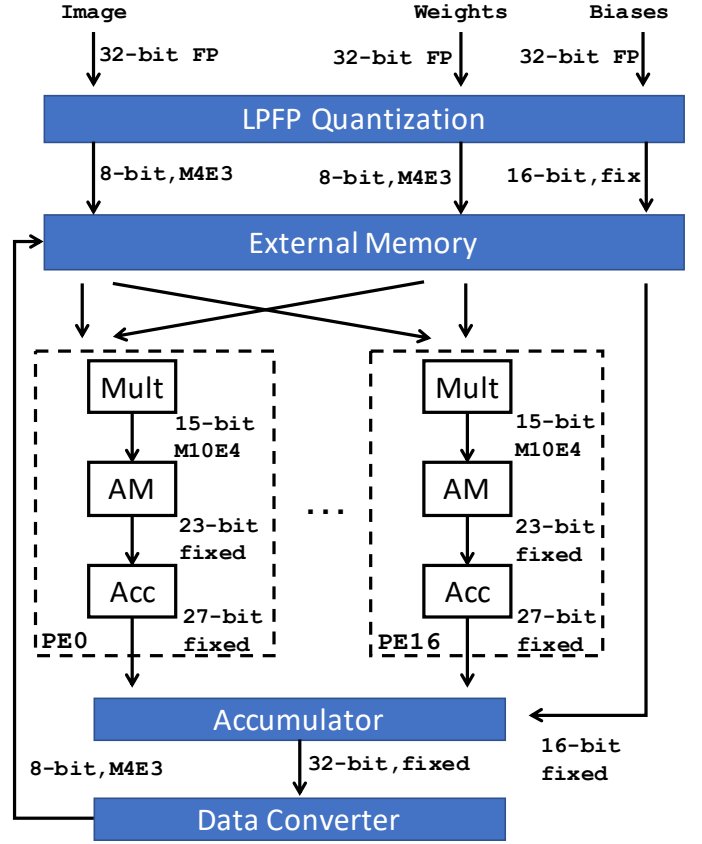


Fig. 2. The data flow in our processor with $M4E3$ data format as an example (FP: floating-point, Mult: LPFP multiplier, AM: alignment module, Acc: accumulator, DC: data converter).

find an optimal combination of LPFP data format and scaling factor for least MSE. In this paper, we assume **the same data format for a CNN and a same scaling factor for each layer**. This assumption can be removed as needed. Furthermore, we choose to use a same optimized data format for all test cases in our experiments, while the problem formulation is to decide a data format for each CNN.

### B. Data Flow in Processor

The data flow to run inference of a quantized network in our processor is shown in Figure 2. In order to explicitly illustrate the data flow, we list the bit width in each step with $M4E3$ data format as an example. All the input image, weights and biases are represented by 32-bit floating-point. In our processor, the raw input image which indicates the input of the first layer and all the weights are quantized with $M4E3$ data format and stored in external memory, while biases are quantized to 16-bit fixed-point to reduce quantization error. Multiplications are performed with the quantized image and weights, and the 15-bit floating-point ($M10E4$) products are converted to 23-bit fixed-point without any precision loss. In this way, all the accumulation can be done in fixed-point accumulators, which consumes fewer resources in FPGA than floating-point accumulators. The final outputs in each output channel are converted to $M4E3$ floating-point again (and
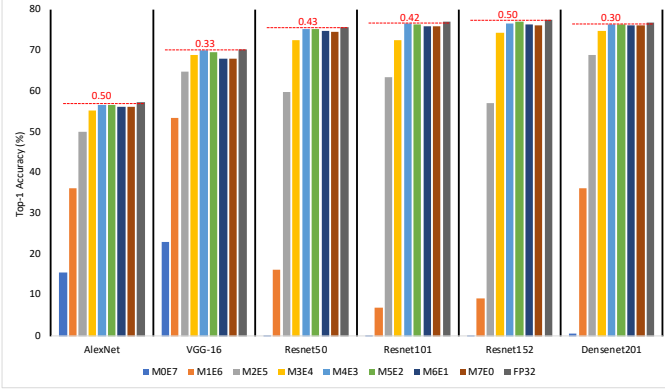
Fig. 3. Top-1 accuracy for different (mantissa, exponent) combinations with respect to different CNNs.



Fig. 4. Top-5 accuracy for different (mantissa, exponent) combinations with respect to different CNNs.

stored in the external memory) before being used by another CNN layer. In the data flow, only the final data conversion step introduces bit truncation and precision loss. However, the precision loss introduced by the final step has little impact on the final accuracy and is validated in Subsection III-C with comprehensive experimental results.

### C. Quantization Results

*1) Experiment Setup:* We implement our LPFP quantization method with C language based on the Darknet framework [37], and the inference process of the quantized network follows the same data flow as that in our processor illustrated in Figure 2. The validation accuracy with single center-crop is then evaluated via the ImageNet validation set (50,000 labelled images) [9]. Our quantization process is run on an Intel (R) Core (TM) i9-7960X CPU working under 2.86GHz, while the evaluation process is run on a Nvidia TITAN Xp GPU. Six representative CNNs including the *slim, medium* and *deep* CNNs are evaluated, as listed in Table II.

*2) 8-bit Quantization:* The detailed validation accuracies on the quantized networks with 8-bit floating-point data format are shown in Figures 3 and 4. We emulate all 8 different (mantissa, exponent) combinations to validate the top-1 and top-5 accuracy of the quantized CNNs, and the 32-bit floating-point results are included as the baseline.

In Figures 3 and 4, the dashed lines illustrate the 32-bit floating-point baseline, while the values above the dashed lines are the accuracy loss compared with the baseline. We can see that our LPFP quantization approach can maintain comparable top-1 and top-5 accuracy to the baseline. On average, the top-1 and top-5 accuracy loss is within 0.5% and 0.3% compared with the full precision results, respectively. Particularly, $M5E2$ always achieves the highest accuracy compared with the other cases. Data formats with more than or equal to 3-bit mantissa all have a low accuracy loss for all the six CNNs, while those with less than 3-bit mantissa can hardly find accurate results. We also compare our proposed approach with the fixed-point situation, marked as $M7E0$ in the figures ($M7E0$ means 1-bit sign, 7-bit mantissa and no exponent, exactly fixed-point). As shown in Figures 3 and 4, $M4E3$ and $M5E2$ outperform the fixed-point for all six benchmarks.
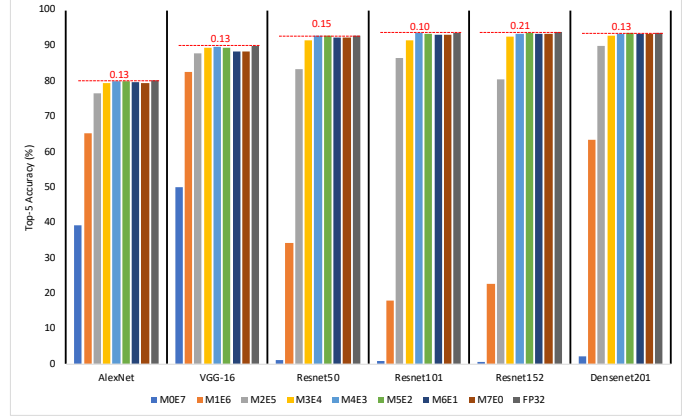
*3) Lower Bit Width Quantization:* We further reduce the bit width from 8-bit to 4-bit and also evaluate the top-1 and top-5 accuracy of the quantized networks. We pick the best (mantissa, exponent) combination for each data format and the results are shown in Figure 5. We can see that both the top-1 and top-5 accuracy decrease when lower bit length is utilized to represent the weights and activations of CNNs. Particularly, the average top-5 accuracy degradations for 7-bit and 6-bit are 0.8% and 4.2%, respectively. However, the accuracy drops dramatically when the bit width decreases to less than 6 bits, which means our LPFP quantization approach can hardly find accurate results without any re-training process.

*4) Comparison with the Prior Quantization Strategies:* $M4E3$ and $M5E2$, which achieve the two best accuracies among all the test cases, are also compared with five typical approaches. We report both the top-1 and top-5 accuracy for all six benchmarks in Table III, where "-" indicates no reported results in the literatures. We use the normalized top-1 accuracy in Table III for the approaches proposed by ARM [24] and Xilinx [25] as reported in their paper. The top-1 and top-5 accuracy in Table III show that our LPFP quantization method without any re-training can outperform the literatures in most cases. Moreover, besides the approach proposed by Nvidia [38], our method is the only one that can reach *deep* networks.

## IV. PROCESSOR ARCHITECTURE

In this section, we discuss in detail the architecture of the processor, which efficiently supports the inference process of quantized networks for various CNNs.

### A. Overview

The overall architecture of the proposed processor is depicted in Figure 6. A floating-point function unit (FPFU), which is composed of multiple processing elements (PEs), is developed to compute the outputs of a layer in parallel. The PE, which is the key component of FPFU, is designed to efficiently perform dot product with LPFP data format. The on-chip memory system (MS) consists of three buffers, *e.g.,* input feature map buffer (IFMB), weight buffer (WB) and output feature map buffer (OFMB). All these three buffers

TABLE III
ACCURACY COMPARISON BETWEEN $M4E3$, $M5E2$, REFERENCES AND FP32. "-" MEANS NO REPORTED RESULTS. RESULTS FOR ARM AND XILINX ARE NORMALIZED TOP-1 ACCURACY AS REPORTED IN THEIR PAPERS WITHOUT CIRCUIT IMPLEMENTATIONS AND WE CONVERT THEM INTO ACTUAL ACCURACY, WHILE OTHERS ARE BASED ON CIRCUIT IMPLEMENTATIONS.

| | Top-1 Accuracy (%) Top-5 Accuracy (%) for each network | | | | | | | | | | | |
| | AlexNet | | VGG16 | | ResNet50 | | ResNet101 | | ResNet152 | | DenseNet201 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Angle-Eye [35] | - | - | 67.72 | 88.06 | - | - | - | - | - | - | - | - |
| Nvidia [38] | 57.05 | 80.06 | 70.84 | - | 73.10 | 91.06 | 74.40 | 91.73 | 74.70 | 91.78 | - | - |
| ARM [24] | 56.71 | - | 70.38 | - | - | - | - | - | - | - | - | - |
| Xilinx [25] | - | - | - | - | 75.80 | - | - | - | - | - | - | - |
| BFP [23] | - | - | 68.32 | - | 72.76 | - | - | - | - | - | - | - |
| FP32 (Baseline) | 57.28 | 80.18 | 70.38 | 89.81 | 75.80 | 92.90 | 77.10 | 93.70 | 77.60 | 93.83 | 76.85 | 93.62 |
| **Ours ($M4E3$)** | **56.69** | **79.99** | **70.05** | **89.68** | **75.25** | **92.75** | **76.68** | **93.60** | **76.79** | **93.44** | **76.40** | **93.43** |
| **Ours ($M5E2$)** | **56.77** | **80.05** | **69.74** | **89.49** | **75.37** | **92.71** | **76.43** | **93.33** | **77.05** | **93.62** | **76.55** | **93.49** |



Fig. 5. Top-1 and top-5 accuracies for different bit width with respect to different CNNs.



Fig. 6. The overall architecture of proposed processor.

are ping-pong architecture to hide the communication time between on-chip and off-chip memories through direct memory access (DMA) module. The central control module (CCM) is designed to arbitrate between different modules. Moreover, the CCM decodes various instructions stored in the instruction RAM (IR) into detailed signals for other modules.

### B. Floating-point Function Unit

FPFU, which is constructed by multiple PEs, is designed to perform convolution in LPFP data format efficiently for performance gain and power reduction. Different parallel computation patterns, including parallel in input feature maps, parallel in output feature maps and parallel in both input and output feature maps, are developed in FPFU and are discussed in the following paragraphs. FPFU receives activations and weights from IFMB and WB, respectively, and distributes the activations and weights to different PEs to perform convolution according to the control signals decoded by CCM.

*1) Architecture of PE:* The PE is designed as a fully pipelined data-flow-based architecture, as shown in Figure 7. Once a PE receives two vectors, it distributes the data to $N_m$ multipliers inside the PE, whose full precision floating-point results are transferred into the alignment module (AM). The full precision floating-point products are aligned and converted to fixed-point numbers without any bit truncation. The aligned products are then fed into four fixed-point adder trees to finalize four dot product processes in parallel, which indicates the feed-forward process of four pixels in two output channels (see
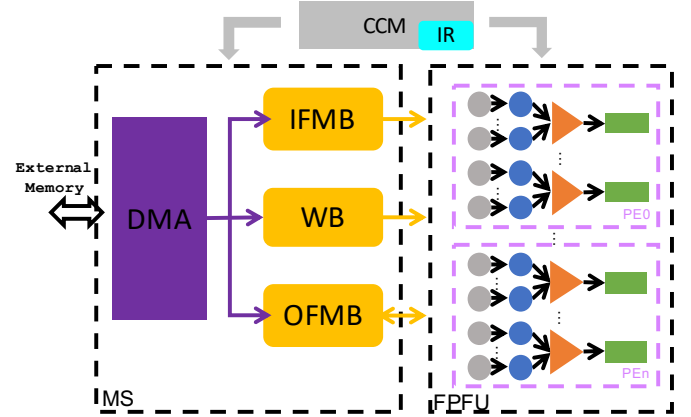
details in Subsubsection IV-B2). The accumulation of partial results (including bias), pooling and activation processes are performed in series inside the post process module (PPM).

The multipliers in each PE are developed for LPFP, which are represented with scientific notation in the sign-and-magnitude format, as illustrated in Eqs. (1) and (3). The multiplication of two LPFP numbers is then divided into three fixed-point components: (1) XOR of the signs; (2) multiplication of mantissas; (3) addition of exponents. Take the $MaEb$ format as an example. An $a$-bit unsigned MAC and a $b$-bit unsigned adder are needed. Although the multiplication of mantissas should be $a+1$-bit considering the first hidden bit of mantissas – "1" for normal numbers and "0" for subnormal numbers – we design the $a$-bit MAC to perform the $a + 1$-bit multiplication to improve per DSP throughput (see details in Subsection V-B). Meanwhile, the exponent bias $E_b$ is not included during addition, because the $E_b$ is the same for all the numbers in one CNN as we assume, and we can address this at the last step to simplify the adders.

*2) Parallel Computation Pattern:* During the convolution process, each pixel in one output channel is calculated as

$$y_i = \sum_{k=0}^{KW \times KH} \sum_{ic=0}^{IC} x_{k,ic} w_{k,ic} + b_i, \quad (7)$$

where $IC$ indicates the number of input channel, $KW$ and $KH$ denotes the width and height of the kernel, and $x, y, w$ and $b$ are input activation, output activation, weight and bias,
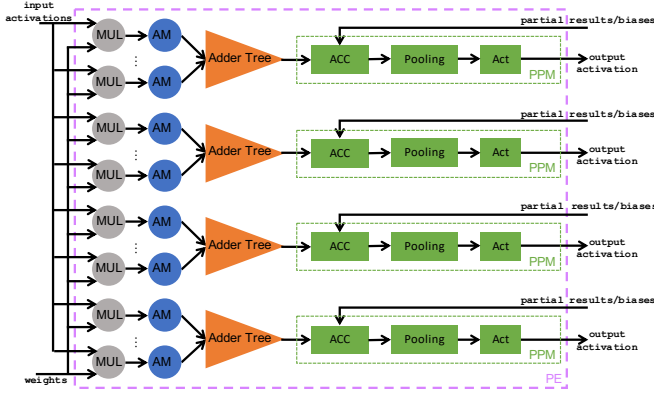
Fig. 7. The architecture of a PE. MUL: LPFP multiplier, AM: alignment module, ACC: accumulator, Act: activation.
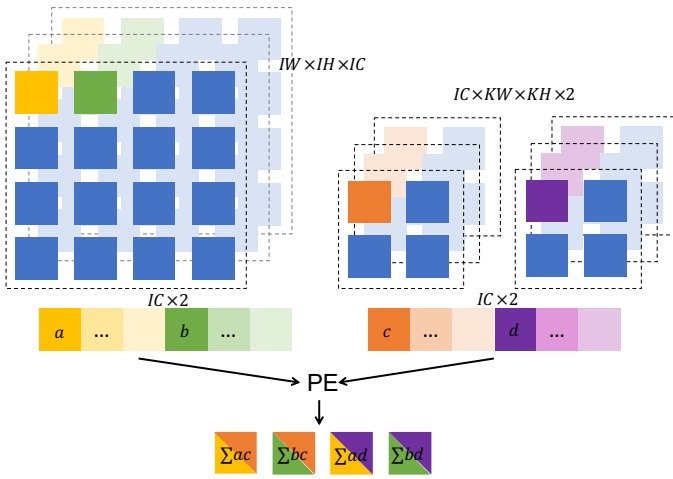


Fig. 8. Parallel computation pattern in one PE.

respectively. In our implementation on FPGA, we implement 4 LPFP multipliers with one DSP slice, which follows the pattern: $(a + b) \times (c + d) = ac + bc + ad + bd$ (see details in Subsection V-B). Therefore, each PE is designed to process convolution in two output channels in parallel, and in each output channel, it will calculate the convolutional results of two pixels at the same time, as shown in Figure 8. To be specific, in the first cycle, the first pixel in $IC$ input channels and the first value in the corresponding kernels are fed into the PE, marked with $a$ and $c$ in Figure 8, respectively. To follow the computation pattern in these four multipliers, the second pixel in $IC$ input channels (marked with $b$), and the corresponding kernels to calculate the pixel in another output channel (marked with $d$) are also fed into the PE. In this way, $a$ and $b$ are reused to produce the pixels in different output channels, while $c$ and $d$ are reused to produce the pixels in different positions of the same output channel. After $KW \times KH$ cycles, four convolution results are produced by one PE.

As illustrated in Subsubsection IV-B1, $N_m$ multipliers are used in each PE, and $IC$ is designed to be $N_m/4$. In this way, $N_m/4$ input channels are calculated in parallel in each PE. With the corresponding weights and biases, 2 pixels in 2 output channels are calculated in parallel. When the number of

input channels is larger than $N_m/4$ and/or when the number of pixels in each output channel is larger than 2 and/or when the number of output channels is larger than 2, multiple rounds of computation are needed in series to finalize the convolution. In order to further increase the parallelism, we use $N_p$ PEs in the FPFU. In different PEs, we can feed in different pixels in input feature maps and weights to perform different parallel computation pattern. For example, the $N_p$ PEs can share the same input feature map and use different weights to parallelize the computation in output channels, or the $N_p$ PEs can share the same weights and use different input feature maps to parallelize the computation in input channels. The $N_m$, $N_p$ and the parallel computation pattern are decided by considering the CNNs, the throughput and the bandwidth requirement. This will be explained with experiments in Subsubsection V-B1.

### C. Memory System

Following the computation pattern in PE, the IFMB and the WB are set to provide $N_m/2$ LPFP input activations and weights to each PE at every cycle, respectively, while the OFMB needs to save 4 output activations from each PE at every cycle. Although each pixel in the output feature map is represented with LPFP data format, we keep the intermediate results with 16-bit precision to reduce accuracy loss. In this way, the bit width of OFMB for each PE is set to 64 bits. As the input activations and/or weights can be shared by different PEs according to different computation patterns, we define $P_{ifm}$ and $P_{ofm}$ ($P_{ifm} \times P_{ofm} = N_p$) to indicate the parallelisms in input feature map and output feature map, respectively. In this definition, $P_{ifm}$ indicates that we have $P_{ifm}$ PE groups where the same weights are shared during calculation, while in each PE group, $P_{ofm}$ PEs share the same input activations. Therefore, the bit width for IFMB, WB and OFMB are $N_m/2 \times P_{ifm} \times BW$, $N_m/2 \times P_{ofm} \times BW$ and $64N_p$, respectively, where $BW$ denotes the bit width of LPFP data format.

The parameters $N_m$, $P_{ifm}$ and $P_{ofm}$ are decided to trade off between the throughput, bandwidth requirement and resource utilization. The sizes of the three buffers also determine the throughput and resource utilization. Previous proposed work applied large enough buffers to store all the activations or weights for one layer [39] to avoid costly off-chip memory access. However, such designs incurred large area and unscalability for larger and deeper CNNs. In our processor, we trade off among the throughput, bandwidth requirement, resource utilization and scalability, and employ the smallest sizes which can hide the DMA communication time. In our implementation on FPGA, we use block RAM to deploy IFMB and OFMB, while we use distributed RAM to deploy WB, as distributed RAM can provide higher bandwidth than block RAM. During inference on our processor, only when all the input feature maps have been processed and reused, or all the weights have been processed and reused, or OFMB is full, will the off-chip memory be accessed for loading new input feature maps, loading new weights or storing output feature maps, respectively.
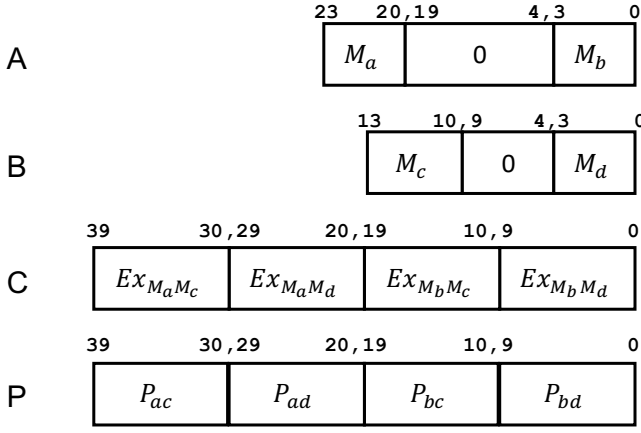
Fig. 9. Data format of the DSP to implement four 4-bit MACs. $M_a$, $M_b$, $M_c$ and $M_d$: the mantissas of LPFP data $a, b, c$ and $d$, respectively; $Ex_{M_a M_c}, Ex_{M_a M_d}, Ex_{M_b M_c}$ and $Ex_{M_b M_d}$: the extra term expressed as $Ex_{M_a M_c} = 1.M_a + 0.M_c$; $P_{ac}, P_{ad}, P_{bc}$ and $P_{bd}$: the mantissas of the product of two LPFP data expressed as $P_{ac} = 1.M_a \times 1.M_c$.

### D. Central Control

The CCM is designed to arbitrate among different modules and control the whole execution process. First, CCM decodes the instructions from IR efficiently and sets the corresponding control registers. Second, different modules are activated according to the control registers and the status of each module is monitored by the control registers as well. Finally, the CCM decides when to fetch the next instruction from the feedback of the control registers. We also design a compiler to generate the block-level instructions.

## V. EVALUATION

In this section, the environment setup of our evaluation is first introduced, then the implementation details and comprehensive experimental results are provided.

### A. Environment Setup

Our processor is implemented on the KC705 evaluation board, which includes a Xilinx Kintex-7 XC7K325T FPGA and a 1GB DDR3 module. First, we explore the parallel computation patterns to find the optimal parameters to best fit the FPGA on KC705. Second, with these parameters, the processor is described in Verilog-HDL, and synthesized and implemented with the Xilinx Vivado 2018.2 Design Suite. Finally, we evaluate the throughput and per DSP throughput of running different networks (shown in Table II) on our processor, and the results are compared with two prior accelerators [18], [33]. The Intel (R) Core (TM) i9-7960X CPU under 2.86GHz working frequency and the Nvidia TITAN Xp GPU with a 12GB DDR5 are also used for comparison. More comprehensive experimental results on VGG16 and YOLO are compared with latest FPGA accelerators [15], [17], [18], [23], [33], [35], [40].

### B. Implementation Details

We use the $M4E3$ data format for FPGA implementation in this paper for two reasons. First, $M4E3$ achieves the top two best validation accuracies among all the LPFP (mantissa, exponent) combinations we tested (see Subsection III-C). Particularly, the average top-1 and top-5 accuracy loss of $M4E3$ compared with 32-bit floating-point are 0.53% and 0.19%, respectively. Second, $M4E3$ only needs a 4-bit fixed-point MAC and a 3-bit fixed-point adder, resulting in fewer resources on FPGA than $M5E2$. To be specific, four 4-bit fixed-point MACs can be implemented inside one DSP48E1 slice in XC7K325T FPGA.

In order to clearly explain the way to implement four MACs with one DSP48E1 slice, we take the multiplication of two normal numbers ($X$ and $Y$) as an example. The mantissa of the product can be explained as:

$$
\begin{aligned}
Prod &= 1.M_x \times 1.M_y \times 2^{(2 - E_x - E_y)} \\
&= (0.M_x \times 0.M_y + (1.M_x + 0.M_y)) \times 2^{(2 - (E_x + E_y))},
\end{aligned}
\tag{8}
$$

where $M_x, M_y, E_x$ and $E_y$ are the mantissas and exponents of $X$ and $Y$, respectively. In Eq. (8), the term $0.M_x \times 0.M_y + (1.M_x + 0.M_y)$ is performed with a 4-bit unsigned fixed-point MAC and the term $E_x + E_y$ is performed with an extra 3-bit unsigned fixed-point adder. As the DSP48E1 slice can be implemented as a MAC followed by $P = A \times B + C$ (where the maximal bit width of $A, B$ and $C$ are 25, 18 and 48, respectively), we add blank bits to the three inputs to fully utilize the functionality of DSP48E1, as shown in Figure 9. During the calculation process, the dot position is kept at the right most position. That is, the terms $0.M_x$ and $0.M_y$ are converted to 4-bit integers, while the extra term $1.M_x + 0.M_y$ is converted to 10-bit integers to make sure that no overlap occurs. In this way, with a few LUTs and FFs to perform additions of the exponents and the extra term $1.M_x + 0.M_y$, four multiplications with $M4E3$ data format can be carried out in on DSP slice (see Table I), thus dramatically increasing the per DSP throughput.

*1) Parallel Exploration:* Since one DSP slice is divided into four 4-bit LPFP MACs in our implementation, the parameters should meet the requirement that $N_m \times N_p = 4 \times \#ofDSP$. Considering the resources of XC7K325T FPGA, we set the targeted number of DSP as 768, which accounts for 91.43% of the available DSPs. We then evaluate the throughput for different CNNs and the bandwidth requirement with respect to different $N_m$ and $N_p$ combinations as shown in Figures 10 and 11, respectively. We also explore different combinations of the parameters $P_{ifm}$ and $P_{ofm}$, and only depict the $P_{ifm}$ and $P_{ofm}$ for achieving the optimal throughput and minimal bandwidth requirement in Figures 10 and 11.

In general, when $N_m$ keeps increasing, the throughput first increases and then decreases when it reaches the peak. The small $N_m$ and large $N_p$ indicate that more output channels are calculated in parallel while large $N_m$ and small $N_p$ mean more input channels are calculated in parallel. When $N_m$ is larger than the total number of input channel (denoted as $IC$), only $IC$ multipliers are used while the rest are wasted, resulting in a low throughput. This is the same for large $N_p$, and the peak throughput comes from balanced $N_m$ and $N_p$. For different CNNs, the peak throughput comes from different $N_m$
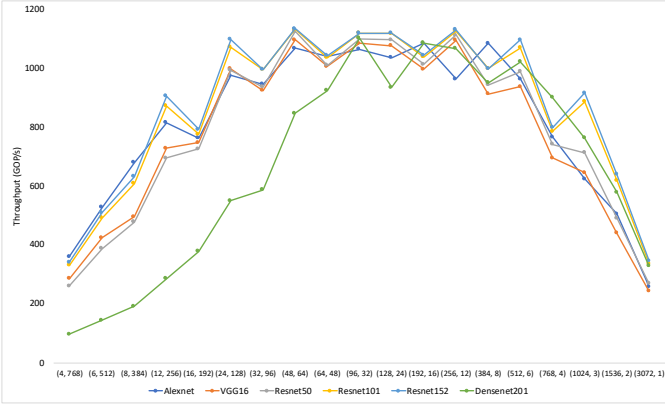
Fig. 10. Throughput for different CNNs with respect to different $N_m$ and $N_p$ combinations.



Fig. 11. Bandwidth requirement with respect to different $N_m$ and $N_p$ combinations.

TABLE IV
RESOURCE UTILIZATION IN XC7K325T.

| Resource | LUT | LUTRAM | FF | BRAM | DSP |
|---|---|---|---|---|---|
| Used | 154625 | 7860 | 180561 | 234.5 | 768 |
| Available | 203800 | 64000 | 407600 | 445 | 840 |
| Utilization | 75.87% | 12.28 | 44.30% | 52.70% | 91.43% |

and $N_p$ combinations due to different network configurations. For example, DenseNet201 has lots of inception layers, which concatenate layers with small output channels (*e.g.,* 32) to form layers with large input channels (*e.g.,* 1568). In this case, larger $N_m$ and smaller $N_p$ incur fewer wasted computations and lead to higher throughput. From Figure 10, we can see that the combination of $N_m = 96$ and $N_p = 32$ results in an optimal throughput for all cases on average.

The bandwidth requirement is extremely high when $N_p$ is large. This is because larger $N_p$ indicates more parallel computations in output channels. Moreover, OFMB is designed to store 16-bit intermediate results, which also lead to higher bandwidth requirement with larger $N_p$. The total bandwidth requirement decreases when $N_p$ decreases, and then increases again since larger $N_m$ needs more bandwidth to load input activations and weights. The smallest bandwidth requirement comes when we have a balanced combination of $N_m$ and $N_p$. As concluded from Figure 11, the optimal combinations are $N_m = 96, N_p = 32$ and $N_m = 128, N_p = 24$. Take the case for optimal throughput, we set $N_m = 96$ and $N_p = 32$ in this implementation.

### C. Experimental Results

*1) Resource Utilization:* Given the parameters that $N_m = 96$ and $N_p = 32$, the detailed post-implementation resource utilization under 200MHz working frequency is listed in Table IV.

*2) Throughput and per DSP throughput for Different CNNs:* Six representative CNNs, including *slim, medium* and *deep* networks (see Table II), are mapped on our processor. When calculating the CNN size, one MAC is counted as two operations. The throughput is measured in GOPS (Giga Operations Per Second), and is reported for different networks on our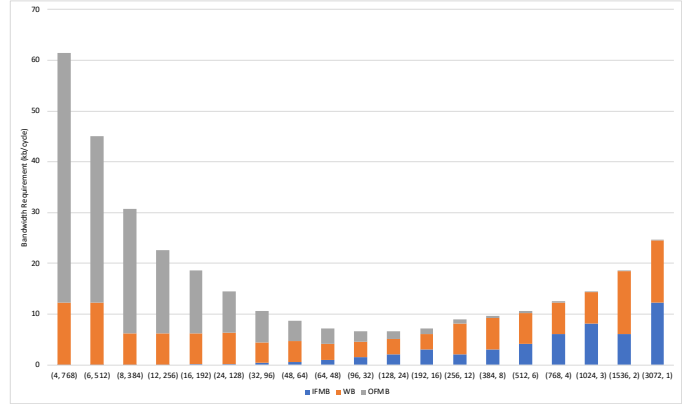 processor, Intel i9 and Nvidia TITAN Xp in Table V. For the evaluation on Intel i9 and Nvidia TITAN Xp, we run the CNNs using the Darknet framework with $batch\_size = 1$ and 32-bit floating-point data format. As the existing studies [18] [33] support multiple networks, we also include their results in Table V.

Compared with the existing accelerators, our processor outperforms them in both throughput and per DSP throughput. Particularly, the average improvement of throughput is 63.5% and 38.6% compared with [18] and [33], respectively. Moreover, the average improvement of per DSP throughput is 2.2× compared with [18]. In the approach proposed in [33], they use LUT to implement multipliers, so we do not compare per DSP throughput with them. Our FPGA processor outperforms Intel i9 by 64.5× in terms of throughput because of the high parallelism in our processor. Although Nvidia TITAN Xp can achieve higher throughput (due to more hardware resources) than our processor does, the average power of Nvidia TITAN Xp is 286W, and the average power efficiency of our processor is 6.9× of the Nvidia TITAN Xp.

*3) Comparison with Previous Accelerators on VGG16:* We run the classification network VGG16 on our processor, and compare the results with six typical studies, as shown in Table VI. We also list the detailed implementation information, such as platform, working frequency and quantization strategy in Table VI. First, our processor, which uses the LPFP quantization scheme, has a negligible top-1 and top-5 accuracy degradation of 0.33% and 0.13%, respectively. Although the work in [15] and [33] can maintain lower accuracy loss than ours, the approach in [15] uses 16-bit floating-point data format, which results in higher bandwidth and memory requirement and lower per DSP throughput, while the approach in [33] needs 144 extra hours for the fine-tuning process. Second, our processor outperforms all the six accelerators in terms of throughput and per DSP throughput. Particularly, the improvements of throughput and per DSP throughput are from 24% to 11.89× and from 92% to 11.14×, respectively. These improvements mainly come from the parallel computation pattern in FPFU and the implementation of four 4-bit MACs within one DSP slice. To the best of our knowledge, this is the first work that can simplify the multiplication to 4-bit and implement four MACs inside one DSP slice while

TABLE V
COMPARISON BETWEEN INTEL I9 CPU, NVIDIA TITAN XP GPU, EXISTING ACCELERATORS AND OUR PROCESSOR WITH RESPECT TO DIFFERENT CNNS. "-" MEANS NO REPORTED RESULTS.

| | Throughput (GOPS)    per DSP throughput (also called DSP efficiency, unit: GOPS/DSP) for each network | | | | | | | | | | |
| | AlexNet | | VGG16 | | ResNet50 | | ResNet101 | | ResNet152 | | DenseNet201 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Intel i9 | 14.1 | - | 19.7 | - | 16.8 | - | 20.1 | - | 19.2 | - | 13.4 | - |
| Nvidia TITAN Xp | 4540.0 | - | 4760.0 | - | 4234.8 | - | 4925.0 | - | 4198.6 | - | 4340.0 | - |
| Ma, et. al. [18] | - | - | 715.9 | 0.47 | 611.4 | 0.40 | - | - | 707.2 | 0.47 | - | - |
| RNA [33] | 687.8 | - | 878.1 | - | 804.3 | - | - | - | - | - | - | - |
| **ours** | **1066.4** | **1.39** | **1086.8** | **1.42** | **1101.9** | **1.43** | **1121.4** | **1.46** | **1121.3** | **1.46** | **1104.7** | **1.44** |

TABLE VI
COMPARISON WITH PRIOR ACCELERATORS ON VGG16. "-" MEANS NO REPORTED RESULTS.

| | Mei, et.al. [15] | Xiao, et. al. [17] | Ma, et. al. [18] | Angle-Eye [35] | RNA [33] | BFP [23] | **ours** |
|---|---|---|---|---|---|---|---|
| Year | 2017 | 2017 | 2018 | 2018 | 2018 | 2019 | **2019** |
| Platform | XC7VX690T | XC7Z045 | Arria 10 GX1150 | XC7Z020 | XC7Z045 | XC7VX690T | **XC7K325T** |
| Frequency (MHz) | 200 | 100 | 200 | 214 | - | 200 | **200** |
| Quantization Strategy | 16-bit floating | 16-bit fixed | 16-bit fixed | 8-bit fixed | 8/4-bit fixed/log | 8-bit block floating | **8-bit floating** |
| Top-1/Top-5 Accuracy (%) | 70.46/89.77 | -/- | -/- | 67.72/88.06 | 70.19/89.81 | 68.31/- | **70.05/89.68** |
| DSP Used | 1728 | 824 | 1518 | 780 | - | 1027 | **768** |
| Throughput (GOPS) | 202.42 | 229.55 | 715.9 | 84.3 (CONV) | 878.11 | 760.83 | **1086.8** |
| per DSP Throughput (GOPS/DSP) | 0.117 | 0.279 | 0.472 | 0.444 | - | 0.741 | **1.42** |
| Power (W) | 10.81 | 9.4 | - | 3.5 | 7.2 | 9.18 | **9.42** |
| Power Efficiency (GOPS/W) | 18.72 | 24.42 | - | 24.1 | 72.8 | 82.88 | **115.4** |

maintaining comparable top-1/top-5 accuracy without any re-training process. Finally, we also show the power efficiency in Table VI, and our processor improves the power efficiency by $39\% - 5.16\times$.

*4) Comparison with Previous Accelerators on YOLO:* We further compare the detection network YOLO [2], [41] with prior accelerators [35], [40], [42], [43] and we use the tiny version of the YOLO network. The comparison results are shown in Table VII, where we also list the mean average precision (mAP) loss of our quantized networks. Compared with the full precision network, the mAP loss of quantized tiny-yolo and tiny-yolo-v2 is 0.3% and 0.1%, respectively. The hardware comparison with prior accelerators shows that our processor is $20.1\times$ and $49.7\times$ higher in terms of throughput for tiny-yolo and tiny-yolo-v2, respectively. Moreover, due to the implementation of four 4-bit MACs within one DSP slice, the per DSP throughput improves by $5\times$ compared with prior accelerators on average.

## VI. RELATED WORK

**Weight and Computation Reduction.** CNNs are typically over-parameterized, and extensive accelerator developers in recent years focus on using CNN approximation algorithms, including weight reduction, computation complexity reduction and quantization to accelerate CNN inference [44]. The accelerator proposed in [45], [46] used Winograd algorithm to reduce the number of multiplication in convolution, thus reducing computation complexity. EIE [39], Cambricon-X [47] and Cambricon-S [14] were the mainstreaming accelerators that benefit from weight and computation complexity reduction techniques. However, the irregularity caused by these

algorithms degrades the parallelism and hardware efficiency [48].

**Quantization.** Accelerators with quantization is another concentration. XNOR_Net [49] applied weights binarization by quantizing weights into {-1, 1} with a scaling factor for AlexNet. The lightweight YOLOv2 [50] was another binarization approach which focused on object detection CNN. Accelerator with ternary representation, which added zero to the binary set, was introduced to help improve the accuracy [51]. Although these accelerators achieve remarkable power and storage saving, they both suffer from significant accuracy loss. Moreover, they all need time-consuming re-training process to compensate for the quantization error. 16-bit quantization oriented accelerators, including floating-point and fixed-point representations, solved the problem of accuracy loss [15], [17], [18], [42]. However, the storage requirement is still huge, and the per DSP throughput is extremely low (less than 0.5GOPS/DSP) because of the usage of 16-bit.

8-bit quantization makes a trade-off between storage and accuracy. The accelerators [33], [35] optimized the computation patterns with 8-bit fixed-point quantization to improve the performance for different CNNs. DNNBuilder [36] was proposed to automatically build DNN accelerators to satisfy the performance and power efficiency demands on embedded and cloud FPGAs, while Cloud-DNN [52] was the framework for mapping DNN models to cloud FPGAs. Block floating-point scheme with 8-bit mantissa was used in [23] to accelerate the inference of CNN while maintaining accuracy. However, all these accelerators need 8-bit MAC to perform convolution, leading to a low per DSP throughput (less than 0.8GOPS/DSP). A more aggressive method quantized the

TABLE VII
COMPARISON WITH PRIOR ACCELERATORS ON YOLO. "-" MEANS NO REPORTED RESULTS.

| | Ma, et.al. [42] | Aristotle [43] | Angle-Eye [35] | Wai, et.al. [40] | ours | |
|---|---|---|---|---|---|---|
| Year | 2017 | 2017 | 2018 | 2018 | **2019** | |
| Platform | XC7V485T | XC7020 | XC7Z020 | Cyclone V | **XC7K325T** | |
| Frequency (MHz) | 143 | 214 | - | 117 | **200** | |
| Quantization Strategy | 16-bit fixed | 8-bit fixed | 8-bit fixed | 8-bit fixed | **8-bit floating** | |
| Network | tiny-yolo | tiny-yolo | tiny-yolo | tiny-yolo-v2 | **tiny-yolo** | **tiny-yolo-v2** |
| mAP loss (%) | - | - | - | - | **0.3** | **0.1** |
| DSP Used | 112 | 198 | - | 122 | **768** | |
| Throughput (GOPS) | 48 | 36.5 | 62.9 | 21.6 | **987.2** | **1095.4** |
| per DSP Throughput (GOPS/DSP) | 0.429 | 0.184 | - | 0.177 | **1.29** | **1.43** |

small values of the weights into 4 bits and keeps the remaining 16 bits as full precision, by dividing the weights into the low-precision and high-precision regions according to the values of the weights [53]. HAQ [54] proposed a mixed precision quantization approach with a trade-off between quantization policy and hardware performance. However, both studies need time-consuming re-training process to compensate for quantization errors.

Different from all the above methods, the proposed LPFP quantization scheme fully exploits the properties of weights and activations, thus obtaining a comparable or better accuracy for *deep* CNNs. Moreover, the LPFP quantization method gets rid of the time-consuming re-training process that needs labelled data and extra computing, because access to labelled data can be difficult in practice as hardware and CNN algorithms are often developed by different parties. Furthermore, with the help of the LPFP quantization method, our processor only needs 4-bit MACs, thus dramatically improving the per DSP throughput. Overall, the proposed processor achieves better performance on FPGA.

## VII. CONCLUSION

We have proposed a low precision floating-point quantization method, called LPFP, to reduce memory size and memory access with negligible accuracy degradation (less than 0.5% for top-1 and 0.3% for top-5 accuracy) for CNN interference. LPFP does not need any re-training. Furthermore, we have reduced the bit width for multiplication to 4-bit with comparable accuracy and implemented four 4-bit MACs within one DSP slice in Xilinx Kintex 7 FPGA family. Experiments using Xilinx KC705 platform and six typical CNN networks show that we achieve an average throughput and per DSP throughput of 1100.4 GOPS and 1.43 GOPS, respectively. Moreover, the average throughput is 64.5× and 1.5× over Intel i9 and existing accelerators, respectively. Particularly for VGG16 and YOLO, we outperform six existing accelerators in terms of average throughput by 3.5× and 27.5×, while improving per DSP throughput by 4.1× and 5×, respectively. To the best of our knowledge, this is the first in-depth work that can simplify the multiplication to 4-bit and accommodate four MACs in one DSP slice while maintaining comparable top-1/top-5 accuracy without any re-training.

## REFERENCES

[1] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.

[2] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 779–788.

[3] D. Amodei, S. Ananthanarayanan, R. Anubhai, J. Bai, E. Battenberg, C. Case, J. Casper, B. Catanzaro, Q. Cheng, G. Chen *et al.*, "Deep speech 2: End-to-end speech recognition in english and mandarin," in *International conference on machine learning*, 2016, pp. 173–182.

[4] J. Dean, G. Corrado, R. Monga, K. Chen, M. Devin, M. Mao, A. Senior, P. Tucker, K. Yang, Q. V. Le *et al.*, "Large scale distributed deep networks," in *Advances in neural information processing systems*, 2012, pp. 1223–1231.

[5] K. Ovtcharov, O. Ruwase, J.-Y. Kim, J. Fowers, K. Strauss, and E. S. Chung, "Accelerating deep convolutional neural networks using specialized hardware," *Microsoft Research Whitepaper*, vol. 2, no. 11, pp. 1–4, 2015.

[6] T. Chen, Z. Du, N. Sun, J. Wang, C. Wu, Y. Chen, and O. Temam, "Diannao: A small-footprint high-throughput accelerator for ubiquitous machine-learning," in *ACM Sigplan Notices*, vol. 49, no. 4, 2014, pp. 269–284.

[7] C. Zhang, G. Sun, Z. Fang, P. Zhou, P. Pan, and J. Cong, "Caffeine: Towards uniformed representation and acceleration for deep convolutional neural networks," in *2016 IEEE/ACM International Conference on Computer-Aided Design*, 2016, pp. 1–8.

[8] X. Wei, C. H. Yu, P. Zhang, Y. Chen, Y. Wang, H. Hu, Y. Liang, and J. Cong, "Automated systolic array architecture synthesis for high throughput cnn inference on fpgas," in *Proceedings of the 54th Annual Design Automation Conference 2017*, 2017, p. 29.

[9] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, "Imagenet large scale visual recognition challenge," *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015.

[10] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.

[11] M. Tan and Q. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," in *International Conference on Machine Learning*, 2019, pp. 6105–6114.

[12] S. Han, J. Pool, J. Tran, and W. Dally, "Learning both weights and connections for efficient neural network," in *Advances in neural information processing systems*, 2015, pp. 1135–1143.

[13] S. Han, H. Mao, and W. J. Dally, "Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding," *arXiv preprint arXiv:1510.00149*, 2015.

[14] X. Zhou, Z. Du, Q. Guo, S. Liu, C. Liu, C. Wang, X. Zhou, L. Li, T. Chen, and Y. Chen, "Cambricon-s: Addressing irregularity in sparse neural networks through a cooperative software/hardware approach," in *2018 51st Annual IEEE/ACM International Symposium on Microarchitecture*, 2018, pp. 15–28.

[15] C. Mei, Z. Liu, Y. Niu, X. Ji, W. Zhou, and D. Wang, "A 200mhz 202.4 gflops@ 10.8 w vgg16 accelerator in xilinx vx690t," in *2017 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*. IEEE, 2017, pp. 784–788.

[16] I. LogiCORE, "Floating-point operator v6. 0," *Xilinx Inc*, 2012.

[17] Q. Xiao, Y. Liang, L. Lu, S. Yan, and Y.-W. Tai, "Exploring heterogeneous algorithms for accelerating deep convolutional neural networks on fpgas," in *2017 54th ACM/EDAC/IEEE Design Automation Conference (DAC)*. IEEE, 2017, pp. 1–6.

[18] Y. Ma, Y. Cao, S. Vrudhula, and J.-s. Seo, "Optimizing the convolution operation to accelerate deep neural networks on fpga," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 26, no. 7, pp. 1354–1367, 2018.

[19] C. Zhang, P. Li, G. Sun, Y. Guan, B. Xiao, and J. Cong, "Optimizing fpga-based accelerator design for deep convolutional neural networks," in *Proceedings of the 2015 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, 2015, pp. 161–170.

[20] N. P. Jouppi, C. Young, N. Patil, D. Patterson, G. Agrawal, R. Bajwa, S. Bates, S. Bhatia, N. Boden, A. Borchers *et al.*, "In-datacenter performance analysis of a tensor processing unit," in *2017 ACM/IEEE 44th Annual International Symposium on Computer Architecture*, 2017, pp. 1–12.

[21] Y. Yu, C. Wu, T. Zhao, K. Wang, and L. He, "Opu: An fpga-based overlay processor for convolutional neural networks," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 2019.

[22] Z. Song, Z. Liu, and D. Wang, "Computation error analysis of block floating point arithmetic oriented convolution neural network accelerator design," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.

[23] X. Lian, Z. Liu, Z. Song, J. Dai, W. Zhou, and X. Ji, "High-performance fpga-based cnn accelerator with block-floating-point arithmetic," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 2019.

[24] L. Lai, N. Suda, and V. Chandra, "Deep convolutional neural network inference with floating-point weights and fixed-point activations," *arXiv preprint arXiv:1703.03073*, 2017.

[25] S. O. Settle, M. Bollavaram, P. D'Alberto, E. Delaye, O. Fernandez, N. Fraser, A. Ng, A. Sirasao, and M. Wu, "Quantizing convolutional neural networks for low-power high-throughput inference engines," *arXiv preprint arXiv:1805.07941*, 2018.

[26] M. Courbariaux, Y. Bengio, and J.-P. David, "Binaryconnect: Training deep neural networks with binary weights during propagations," in *Advances in neural information processing systems*, 2015, pp. 3123–3131.

[27] Z. Lin, M. Courbariaux, R. Memisevic, and Y. Bengio, "Neural networks with few multiplications," *arXiv preprint arXiv:1510.03009*, 2015.

[28] P. Colangelo, N. Nasiri, E. Nurvitadhi, A. Mishra, M. Margala, and K. Nealis, "Exploration of low numeric precision deep learning inference using intel® fpgas," in *2018 IEEE 26th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*. IEEE, 2018, pp. 73–80.

[29] E. Park, S. Yoo, and P. Vajda, "Value-aware quantization for training and inference of neural networks," in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 580–595.

[30] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[31] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[32] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4700–4708.

[33] C. Luo, Y. Wang, W. Cao, P. H. Leong, and L. Wang, "Rna: An accurate residual network accelerator for quantized and reconstructed deep neural networks," in *2018 28th International Conference on Field Programmable Logic and Applications (FPL)*. IEEE, 2018, pp. 60–603.

[34] D. Zuras, M. Cowlishaw, A. Aiken, M. Applegate, D. Bailey, S. Bass, D. Bhandarkar, M. Bhat, D. Bindel, S. Boldo *et al.*, "Ieee standard for floating-point arithmetic," *IEEE Std 754-2008*, pp. 1–70, 2008.

[35] K. Guo, L. Sui, J. Qiu, J. Yu, J. Wang, S. Yao, S. Han, Y. Wang, and H. Yang, "Angel-eye: A complete design flow for mapping cnn onto embedded fpga," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 37, no. 1, pp. 35–47, 2017.

[36] X. Zhang, J. Wang, C. Zhu, Y. Lin, J. Xiong, W.-m. Hwu, and D. Chen, "Dnnbuilder: an automated tool for building high-performance dnn hardware accelerators for fpgas," in *Proceedings of the International Conference on Computer-Aided Design*. ACM, 2018, p. 56.

[37] J. Redmon, "Darknet: Open source neural networks in c," http://pjreddie.com/darknet/, 2013–2016.

[38] S. Migacz, "8-bit inference with tensorrt," in *GPU Technology Conference*, 2017.

[39] S. Han, X. Liu, H. Mao, J. Pu, A. Pedram, M. A. Horowitz, and W. J. Dally, "Eie: efficient inference engine on compressed deep neural network," in *2016 ACM/IEEE 43rd Annual International Symposium on Computer Architecture*, 2016, pp. 243–254.

[40] Y. J. Wai, Z. bin Mohd Yussof, S. I. bin Salim, and L. K. Chuan, "Fixed point implementation of tiny-yolo-v2 using opencl on fpga," *International Journal of Advanced Computer Science and Applications*, vol. 9, no. 10, pp. 506–512, 2018.

[41] J. Redmon and A. Farhadi, "Yolo9000: better, faster, stronger," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 7263–7271.

[42] J. Ma, L. Chen, and Z. Gao, "Hardware implementation and optimization of tiny-yolo network," in *International Forum on Digital TV and Wireless Multimedia Communications*. Springer, 2017, pp. 224–234.

[43] K. Guo, S. Han, S. Yao, Y. Wang, Y. Xie, and H. Yang, "Software-hardware codesign for efficient neural network acceleration," *IEEE Micro*, vol. 37, no. 2, pp. 18–25, 2017.

[44] E. Wang, J. J. Davis, R. Zhao, H.-C. Ng, X. Niu, W. Luk, P. Y. Cheung, and G. A. Constantinides, "Deep neural network approximation for custom hardware: Where we've been, where we're going," *arXiv preprint arXiv:1901.06955*, 2019.

[45] L. Lu, Y. Liang, Q. Xiao, and S. Yan, "Evaluating fast algorithms for convolutional neural networks on fpgas," in *2017 IEEE 25th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*. IEEE, 2017, pp. 101–108.

[46] D. Wu, J. Chen, W. Cao, and L. Wang, "A novel low-communication energy-efficient reconfigurable cnn acceleration architecture," in *2018 28th International Conference on Field Programmable Logic and Applications (FPL)*. IEEE, 2018, pp. 64–643.

[47] S. Zhang, Z. Du, L. Zhang, H. Lan, S. Liu, L. Li, Q. Guo, T. Chen, and Y. Chen, "Cambricon-x: An accelerator for sparse neural networks," in *2016 49th Annual IEEE/ACM International Symposium on Microarchitecture*, 2016, pp. 1–12.

[48] S. Wang, Z. Li, C. Ding, B. Yuan, Q. Qiu, Y. Wang, and Y. Liang, "C-lstm: Enabling efficient lstm using structured compression techniques on fpgas," in *Proceedings of the 2018 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*. ACM, 2018, pp. 11–20.

[49] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi, "Xnor-net: Imagenet classification using binary convolutional neural networks," in *European Conference on Computer Vision*, 2016, pp. 525–542.

[50] H. Nakahara, H. Yonekawa, T. Fujii, and S. Sato, "A lightweight yolov2: A binarized cnn with a parallel support vector regression for an fpga," in *Proceedings of the 2018 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*. ACM, 2018, pp. 31–40.

[51] A. Prost-Boucle, A. Bourge, F. Pétrot, H. Alemdar, N. Caldwell, and V. Leroy, "Scalable high-performance architecture for convolutional ternary neural networks on fpga," in *2017 27th International Conference on Field Programmable Logic and Applications*, 2017, pp. 1–7.

[52] Y. Chen, J. He, X. Zhang, C. Hao, and D. Chen, "Cloud-dnn: An open framework for mapping dnn models to cloud fpgas," in *Proceedings of the 2019 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*. ACM, 2019, pp. 73–82.

[53] E. Park, D. Kim, and S. Yoo, "Energy-efficient neural network accelerator based on outlier-aware low-precision computation," in *2018 ACM/IEEE 45th Annual International Symposium on Computer Architecture*, 2018, pp. 688–698.

[54] K. Wang, Z. Liu, Y. Lin, J. Lin, and S. Han, "Haq: Hardware-aware automated quantization with mixed precision," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8612–8620.