

# A System-Level Solution for Low-Power Object Detection

Fanrong Li<sup>1,2\*</sup>, Zitao Mo<sup>1,2\*</sup>, Peisong Wang<sup>1</sup>, Zejian Liu<sup>1,2\*</sup>,  
 Jiayun Zhang<sup>1\*</sup>, Gang Li<sup>1,2</sup>, Qinghao Hu<sup>1</sup>, Xiangyu He<sup>1,2</sup>, Cong Leng<sup>1,3</sup>,  
 Yang Zhang<sup>1,3</sup>, Jian Cheng<sup>1,2,3,4</sup>

<sup>1</sup>Institute of Automation, Chinese Academy of Sciences

<sup>2</sup>University of Chinese Academy of Sciences, <sup>3</sup>AiRiA

<sup>4</sup>CAS Center for Excellence in Brain Science and Intelligence Technology

{lifanrong2017, mozitao2017}@ia.ac.cn, {gang.li, jcheng}@nlpr.ia.ac.cn

## Abstract

*Object detection has made impressive progress in recent years with the help of deep learning. However, state-of-the-art algorithms are both computation and memory intensive. Though many lightweight networks are developed for a trade-off between accuracy and efficiency, it is still a challenge to make it practical on an embedded device. In this paper, we present a system-level solution for efficient object detection on a heterogeneous embedded device. The detection network is quantized to low bits and allows efficient implementation with shift operators. In order to make the most of the benefits of low-bit quantization, we design a dedicated accelerator with programmable logic. Inside the accelerator, a hybrid dataflow is exploited according to the heterogeneous property of different convolutional layers. We adopt a straightforward but resource-friendly column-prior tiling strategy to map the computation-intensive convolutional layers to the accelerator that can support arbitrary feature size. Other operations can be performed on the low-power CPU cores, and the entire system is executed in a pipelined manner. As a case study, we evaluate our object detection system on a real-world surveillance video with input size of  $512 \times 512$ , and it turns out that the system can achieve an inference speed of 18 fps at the cost of 6.9W (with display) with an mAP of 66.4 verified on the PASCAL VOC 2012 dataset.*

## 1. Introduction

Since AlexNet [7] won the 2012 large-scale image recognition contest, Deep Convolutional Neural Networks (DCNNs) have shown increasing performance in various computer vision tasks. CNN's impressive performance is mainly due to its high complexity and capacity, in

other words, the great number of parameters and computations. Therefore, high-performance hardwares such as GPUs (clusters) are often utilized for acceleration. However, as for embedded and mobile devices such as drones, security cameras, and smart glasses, GPU-based solutions are not the best choice due to the limitation of volume and power consumption. In addition, modern GPUs that designed for general propose processing are not flexible enough to deal with low-bit integer values less than 8-bit without efforts on tuning the codes. As a result, FPGA-based accelerators are gaining popularity in recent years for both industrial and academic communities.

As for memory efficiency, we find that the advantages of the recent depthwise convolution [3, 5] are apparent. Unlike traditional convolution, in depthwise convolution, each output feature map relies solely on a single input feature map in the previous layer, which dramatically reduces the amount of computations and the demand of on-chip storage. In terms of resource and energy efficiency, recent logarithmic computation [4, 12, 8] has shown its promise. It quantizes the weight as power-of-two in order to efficiently translate multiplication into bit shift operation, which can get rid of the limitation of insufficient on-chip DSP blocks.

Considering the advantages of depthwise convolution and logarithmic computation mentioned above, we put forward an end-to-end hardware-software co-design for low-power object detection on resource-constraint FPGA. Our proposed solution can achieve relatively high performance under extremely low resource budget while retaining considerable accuracy. The contribution of this work can be summarized as follows:

- We propose a dedicated object detection accelerator for customized MobileNet-SSD [9, 5] algorithm through software-hardware co-design. Specifically, we quantize the activations and weights to 4-bit integer and 3-bit power-of-two integer respectively, and

\*These authors contributed equally.

present a fused-layer architecture with shift-based processing elements.

- We adopt a column-prior strategy to map the detection network to the accelerator, which can reduce resource consumption. Besides, a hybrid dataflow is introduced to reuse output or weights according to the heterogeneous property of different layers.
- We highlight the entire pipeline of our heterogeneous system design, including hardware accelerator, host processing and thread management of the main processor, and describe each stage in details.
- We verify the performance of our design on heterogeneous devices Ultra96 SoC that targets to IoT applications. Experiments show that the entire system can reach an inference speed of 18 fps at the cost of around 6.9W.

The rest of the paper is organized as follows. Section 2 describes the quantization algorithm, with which we quantize weights to the power-of-two and enables resource-friendly shift-base multiplications. Section 3 briefly presents the overall system architecture. Section 4 introduces the architecture of the dedicated accelerator, including Processing Elements (PEs), tiling strategy, and dataflow. Section 5 reports the experimental results as well as multithread management on low-power CPUs.

## 2. Quantization

To make the CNN model compatible with our hardware architecture design, we introduce a three-step quantization method, i.e., uniform activation quantization, power-of-two weight quantization as well as scale quantization, as illustrated in Figure 1. It is worth noting that through the proposed three-step quantization, *all computing can be transformed into fixed-point operations, without any floating-point values.*

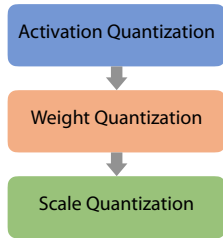


Figure 1. Three-Step Quantization Pipeline.

### 2.1. Uniform Activation Quantization

For  $M$ -bit activation quantization, we want to quantize all the positive activations into the set  $\mathcal{A} = \{0, 1, 2, \dots, 2^M - 1\}$ . As with many other fixed-point

quantization methods, we also introduce a scaling factor  $\alpha$  to lower the quantization error, making the quantization set into

$$\mathcal{A} = \{0, 1, 2, \dots, 2^M - 1\} * \alpha$$

To turn all activations into fixed-point numbers, we can quantize the floating-point activation to the nearest point in the set  $\mathcal{A}$ . The  $2^M - 1$  quantization thresholds can be set to the medians of two successive quantized values:

$$t_i = \frac{(i-1)\alpha + i\alpha}{2} = (i - \frac{1}{2})\alpha \quad (1)$$

for  $i = 1, \dots, 2^M - 1$

Thus the quantization function  $Q_a$  can be formulated as

$$Q_A(x) = \begin{cases} (2^M - 1)\alpha & x > t_{2^M-1}, \\ i\alpha & x \in (t_i, t_{i+1}], \\ 0 & x \leq t_1 \end{cases} \quad (2)$$

### 2.2. Power-of-Two Weight Quantization

For weight, we utilize power-of-two quantization. In this way, the floating-point multiplications within the convolution can be transformed into shifting operations, which can dramatically lower the complexity of CNN and hardware design. The 4-D weight tensor consists of  $n$  kernels of size  $w \times h \times c$ , which are quantized by using different scaling factors. More specifically, the 4-D tensor  $\mathcal{W} \in \mathbb{R}^{w \times h \times c \times n}$  is reshaped into a matrix  $\mathbf{W} \in \mathbb{R}^{(w*h*c) \times n}$ , where each column  $\mathbf{w}_i \in \mathbb{R}^{w*h*c}$  corresponds to a 3-D kernel. To lower the quantization error, a floating-point scaling factor  $\beta_i$  is introduced for each kernel  $\mathbf{w}_i$ , i.e., for  $N$ -bit quantization, the problem is to select weight values from the set

$$\mathcal{B}_i = \{0, \pm 2^0, \pm 2^1, \dots, \pm 2^{2^N-1-2}\} * \beta_i$$

Here we also use the nearest quantization and the  $2^N - 2$  quantization thresholds can also be determined by the medians of two successive quantized values, as in the activation quantization.

### 2.3. Scale Quantization

By activation and weight quantization, the convolution can be performed with only fixed-point operations. However, the whole network still requires floating-point operations due to the introduced scaling factors, bias term of convolution, as well as some other layers like Batch Normalization. To further eliminate the above mentioned floating-point operations, we introduce the scale quantization, which consists of two parts:

**Scale merge:** For the  $l$ -th layer, the input activation  $X$  can be represented by  $X = \alpha \hat{X}$ , where  $\hat{X}$  is the fixed-point version of  $X$  and  $\alpha$  is the scaling factor. Similarly, the  $\mathbf{w} = \beta \hat{\mathbf{w}}$  where  $\hat{\mathbf{w}}$  is one of the fixed-point kernels.

Table 1. Quantization results on ImageNet classification (top-1 accuracy). The #Act., #Wei. and #Sca. represent the number of bits for activations, weights, and scaling factors, respectively.

Model	#Act.	#Wei.	#Sca.	Accuracy
MobileNet	Full	Full	Full	70.1
MobileNet	8	3	8	68.3
MobileNet	4	3	8	68.1

For simplicity, we discard the kernel index. Considering the Batch Normalization term, the convolutional layer can be represented by the following equation:

$$\begin{aligned}
 Y &= \alpha' \hat{Y} = Q_A(BN(\beta \hat{\mathbf{w}} \otimes \alpha \hat{X})) \\
 &= Q_A(\gamma \alpha \beta \hat{\mathbf{w}} \otimes \hat{X} + b) \\
 &= Q_A(a \hat{\mathbf{w}} \otimes \hat{X} + b)
 \end{aligned} \tag{3}$$

where  $Y$  is the output activation,  $\hat{Y}$  is the fixed-point version of output activations, and the  $\alpha'$  is the scaling factor for outputs.  $BN(x) = \gamma x + b$  is the batch normalization layer and  $\otimes$  is the convolution.

To further merge out the output scaling factor, we can divide both sides of Eq. 3 by  $\alpha'$ , resulting in the following equation:

$$\begin{aligned}
 \hat{Y} &= \hat{Q}_A\left(\frac{a}{\alpha'} \hat{\mathbf{w}} \otimes \hat{X} + \frac{b}{\alpha'}\right) \\
 &= \hat{Q}_A(a' \hat{\mathbf{w}} \otimes \hat{X} + b')
 \end{aligned} \tag{4}$$

Note that in the activation quantization function need to be changed accordingly. By defining  $\hat{t}_i = \frac{t_i}{\alpha'}$ , the new quantization function becomes:

$$\hat{Q}_A(x) = \text{round}(\text{clip}(x, 0, 2^M - 1)), \tag{5}$$

where  $\text{round}(x)$  is the rounding operation, and  $\text{clip}(x, u, v)$  clips  $x$  within  $u$  and  $v$ .

**Scale quantization:** In Eq. 4, only the  $a'$  and  $b'$  are floating-points. Note that Eq. 4 only corresponds to one 3-D kernel, for the convolutional layer, there are  $n$  pairs of  $a'$  and  $b'$ , denoted by  $\mathbf{a}'$  and  $\mathbf{b}'$ . In the scale quantization, we need to quantize these values into fixed-point numbers.

During the scale quantization, no scaling factors could be incorporated. However, direct quantizing of  $\mathbf{a}'$  and  $\mathbf{b}'$  will introduce large quantization error. Here we search for the binary point position, resulting in the following set to be quantized into:

$$C = \{0, \pm 1, \pm 2, \dots, \pm 2^{K-1} - 1\} * 2^d,$$

where  $d$  represents the binary point position. More specifically, when  $d$  go through from 0 to -15, we find the best  $d$  that minimize the quantization error for  $\mathbf{a}'$  and  $\mathbf{b}'$ .

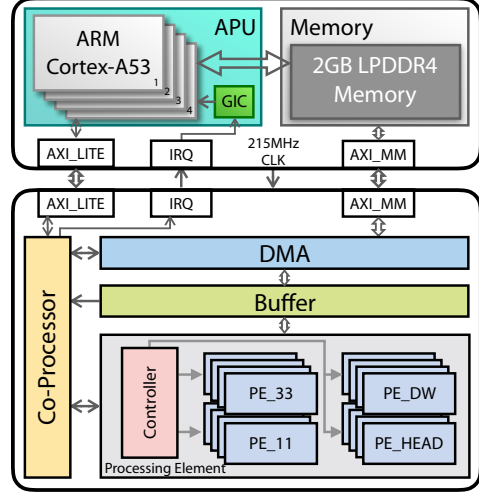


Figure 2. Architecture of the entire system.

## 2.4. Optimization

The optimization problem can be solved efficiently using Lloyd’s algorithm. Take the activation quantization problem of section 2.1 for example, during the assignment step, all activation data points are quantized into the nearest fixed-point values in the set of  $\mathcal{A}$  according to the quantization function  $Q_A(x)$ . In the update step, the new scaling factor can be obtained by solving a one-dimensional optimization problem:

$$\alpha^* = \operatorname{argmin}_{\alpha} \sum_x (x - Q_A(x))^2 \tag{6}$$

By iterative quantization, we could find the optimal scaling factors as well as the quantized values.

After the activation quantization and weight quantization, we need to fine-tune the whole network to restore accuracy.

## 2.5. Performance

The experiments are conducted on the ImageNet classification benchmark, results are shown in Table 1. The results illustrate that the three-step quantization approach has only minimal accuracy drop compared with the floating-point counterpart.

## 3. System Architecture

Our detection network targets to run on the Ultra96 development board, which is a heterogeneous embedded system containing both programmable logic and low-power CPU cores. A 2GB DDR4 is shared by Programmable Logic (PL) and Processing System (PS). Since convolutional layers dominate most of the inference time, we imple-

ment a dedicated CNN accelerator with the Programmable Logic.

The entire system includes the following functional layers. Data forward layer: decode video streams. Encode layer: organize data into the specific pattern for FPGA accelerator. FPGA layer: perform all convolution on the dedicated accelerator. Decode layer: organize extracted features from the accelerator to the storage pattern for CPU. Mbox-conf-reshape layer: reshape bounding boxes. Mbox-conf-softmax layer: softmax layers of the detection. Mbox-conf-flatten layer: reshape data. Detection and visualize layer: generate detecting results and display on the screen. All the layers except for FPGA layer are executed on CPU. All operations before the FPGA layers are referred to as pre-processing, while those operations after the FPGA layer are post-processing.

At the very beginning, images together with the weights and instructions of a specific CNN are stored in DDR. The CPU initiates a calculation request and transfer instructions to the accelerator through AXI. The accelerator receives instructions and completes all convolution computation. Note that the accelerator has its own instruction set, and it can complete the calculations independently unless interrupted by exceptions. Results of the FPGA layer are sent back to CPU for post-processing. Multi-thread technique is exploited to make the most use of 4 low-power ARM cores. The entire system works in a pipelined manner, and the system architecture is shown in Figure 2.

#### 4. Dedicated Accelerator

In this section, we first describe the overall architecture of our accelerator, which exploits multiple PEs for high computing parallelism. Then the design of PE is introduced. After that, the column-prior tiling strategy is presented to support the arbitrary size of input feature maps under limited resources. Finally, a hybrid dataflow is proposed for more efficiency.

##### 4.1. Overall Architecture

Figure 3 shows the overall architecture of our accelerator with different types of PEs inside. The Co-Processor module controls the entire computation flow. It parses instructions to generate control information for the Memory Controller and different kinds of PEs. The addresses of activations and weights are calculated by the Memory Controller, with which all kinds of data can be sent to the proper destinations. Prefetching is enabled since we implement a 4KB instructions cache inside the Co-processor. Note that some cache features are unavailable in this design because they are unnecessary for a specific accelerator without branch and jump instructions. Controllers for different types of PEs generate control signals according to the control information received from Co-Processor. IARAM and

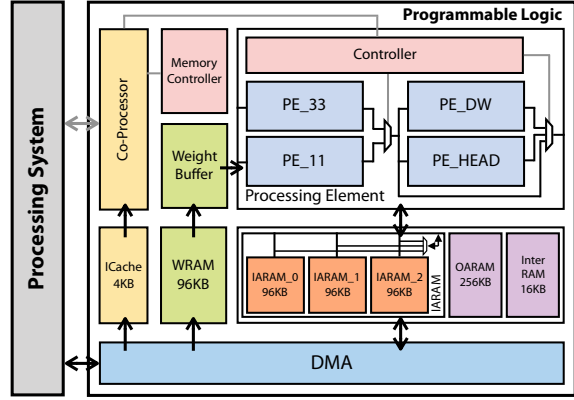


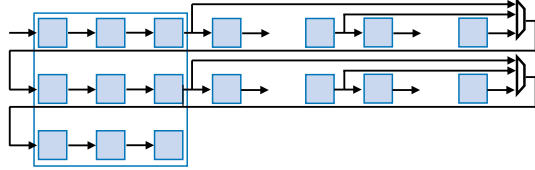
Figure 3. Architecture of dedicated CNN accelerator with only one for each type of PEs.

OARAM are used to store the intermediate feature maps during computation, where IARAM is implemented with three banks, providing sufficient bandwidth to complete the  $3 \times 3$  convolution more efficiently. And the IARAMs and OARAM can be logically swapped between the computation of two adjacent layers. We implement two-level weight caches (Weight buffer and WRAM) with on-chip registers and BRAMs, which can provide sufficient bandwidth for computing.

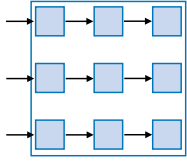
##### 4.2. Processing Elements

Heterogeneous nature of  $1 \times 1$  convolution and depthwise convolution may make the reuse of processing elements costly, so reusing PEs does not necessarily lead to benefits and is contrary to our original intention to design a dedicated low-power accelerator. Therefore, PEs are specialized for different kinds of convolutional layers, i.e.,  $3 \times 3$  convolution (PE\_33),  $1 \times 1$  convolution (PE\_11), and depthwise convolution (PE\_DW) for the consideration of reducing the control complexity and improving hardware efficiency. To efficiently compute the location offsets in the detection algorithm, PE\_HEAD is necessary. Each type of PEs is mainly composed of multipliers and reduction trees, as well as modules that can selectively execute the ReLU and Batch Normalization functions. Each PE processes with only one kernel at a time.

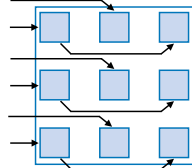
Different from some previous work using line buffer, we implement  $3 \times 3$  convolution in PE\_33 more efficiently, as shown in Figure 4. The input image is divided into three parts according to row number and stored in three IARAMs. During the computation, inputs in three continuous rows can be fetched from different IARAMs simultaneously. Compared to line buffer implementation, it reduces data-preparing time and register consumption. Besides, as for the  $3 \times 3$  convolution with stride=2, each IARAM can provide higher bandwidth to support jump connection for



(a) Line buffer convolution.



(b) PE\_33, stride=1.



(c) PE\_33, stride=2.

Figure 4. The implementation for  $3 \times 3$  convolution with different strategies: (a) Line buffer convolution; (b) Our implementation of  $3 \times 3$  convolution with stride=1; (c) Our implementation of  $3 \times 3$  convolution with stride=2

the registers, as shown in Figure 4(b). Therefore, only the necessary calculations are performed, which can achieve  $4 \times$  speedup than the original convolution based on classic line buffer.

Depthwise convolutional layer can be fused with its adjacent layers in a pipelined manner to speedup computation due to its less data-dependent property. With this insight, in this work, we introduce two types of cascaded PEs to the architecture of our accelerator, which can be summarized as follows.

- **PE\_33, PE\_DW.** The results of  $3 \times 3$  convolution can be sent to PE\_DW directly. Different from PE\_33, PE\_DW are processing with line buffer to accommodate the continuous inflow of data. This manner works in conjunctions with our column-prior tiling strategy to reduce the consumption of registers, which we will present in section 4.3.
- **PE\_11, PE\_DW.** Similarly,  $1 \times 1$  convolution and depthwise convolution can also be processed in a fused manner. During computation, input activations are fetched from one of three input buffers, and the results of  $1 \times 1$  convolution are sent to PE\_DW immediately and processed on the fly. The final results are written back to the corresponding output buffer.

As mentioned in section 2, activations and weights of the network are quantized to low bits. Specifically, the weights are quantized to power-of-two, which enables us to replace multipliers with shift operators. Compared with normal multiplications, it can reduce resource and power consumption. We conduct an experiment to verify the benefits of this shift-based multipliers, which shows that shift-based multi-

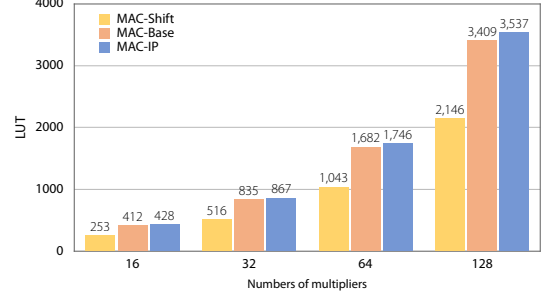


Figure 5. LUT consumption of different implementations of MACs. MAC-Shift (activation 4b/weight 3b) is our implementation of multipliers using shift operations, while MAC-Base (4b/4b) is direct multiplication and MAC-IP (4b/4b) is multiplications using Xilinx IP. Reduction trees are also included in all three cases. Note that if we use multipliers, we have to use 4b/4b inputs in order to represent numbers from -4 to 4.

Table 2. Notation for tiling strategy and dataflow.

Variables	Descriptions
$W_T$	width (column) of a tile of feature maps
$H_T$	height (row) of a tile of feature maps
$K_T$	parallelism on output channel dimension
$C_T$	parallelism on the input channel dimension
$N_k$	number of tiles along the filter dimension
$N_c$	number of tiles along the channel dimension

plication can reduce the usage of LUT by approximately 40%, as shown in Figure 5.

### 4.3. Column-Prior Tiling Strategy

Under the limited on-chip resources, tiling is necessary to map convolutional layers to the accelerator. We adopt a column-prior tiling strategy, as shown in Figure 6, which can reduce both latency and register consumption. We take a feature map with size  $256 \times 256$  as an example, which is expected to be divided into two parts to fit into the limited on-chip buffers. As for the row-prior manner, a tile with size  $128 \times 256$  is generated after  $1 \times 1$  convolution and can be sent to PE\_DW immediately for the processing of depthwise convolution. In this situation, at least  $2 \times 256 + 3 = 515$  registers are required for applying line buffer convolution. However, if the feature maps are divided into the size of  $256 \times 128$  in a column-prior manner, only  $2 \times 128 + 3 = 259$  registers are needed. Thus register consumption can be approximately halved. Similarly, invalid cycles caused by filling registers are also reduced, which will also be beneficial to latency and efficiency.

Since the feature maps are divided into several tiles by column index, overlapping between adjacent tiles are introduced. Suppose that we can obtain output tiles with five valid columns after  $1 \times 1$  and depthwise convolution



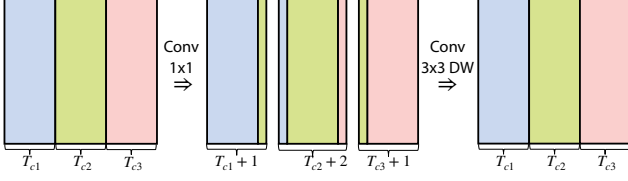


Figure 6. A particular case of column-prior tiling strategy applied to  $1 \times 1$  and depthwise convolution (stride=1). The input feature maps are divided into three tiles and transferred from DDR to on-chip BRAMs sequentially. As shown in the middle of the figure, extra features columns from adjacent tiles are necessary.

---

**Algorithm 1:** Output stationary dataflow for  $1 \times 1$  convolution.

---

```

for  $h = 0 : H_T$  do
  for  $w = 0 : W_T$  do
    for  $n_k = 0 : N_k$  do
      for  $n_c = 0 : N_c$  do
        Parallel for  $k = (n_k-1)K_t : n_k K_t$  do
          Parallel for  $c = (n_c-1)C_T : n_c C_T$  do
             $P = I[c][w][h] * W[k][c];$ 
             $O[k][w][h] = O[k][w][h] + P;$ 
            // Partial sums keep stationary in PE
            until a valid output is obtained;
          // Sent  $O[(n_k-1)K_t : n_k K_t][w][h]$  to Out_buf;
        // Sent  $O[(n_k-1)K_t : n_k K_t][w][h]$  to Out_buf;
      // Sent  $O[(n_k-1)K_t : n_k K_t][w][h]$  to Out_buf;
    // Sent  $O[(n_k-1)K_t : n_k K_t][w][h]$  to Out_buf;
  // Sent  $O[(n_k-1)K_t : n_k K_t][w][h]$  to Out_buf;
// Sent  $O[(n_k-1)K_t : n_k K_t][w][h]$  to Out_buf;

```

---



---

**Algorithm 2:** Weight stationary dataflow for  $1 \times 1$  convolution.

---

```

for  $n_k = 0 : N_k$  do
  for  $n_c = 0 : N_c$  do
    // Fetch weights from weight buffer and keep weights
    stationary in PEs;
    for  $h = 0 : H_T$  do
      for  $w = 0 : W_T$  do
        Parallel for  $k = (n_k-1)K_t : n_k K_t$  do
          Parallel for  $c = (n_c-1)C_T : n_c C_T$  do
             $P = I[c][w][h] * W[k][c];$ 
             $O[k][w][h] = O[k][w][h] + P;$ 
            // Keep partial sums in Inter buffer;
          // Sent  $O[(n_k-1)K_t : n_k K_t][:][:]$  to out_buf;
        // Sent  $O[(n_k-1)K_t : n_k K_t][:][:]$  to out_buf;
      // Sent  $O[(n_k-1)K_t : n_k K_t][:][:]$  to out_buf;
    // Sent  $O[(n_k-1)K_t : n_k K_t][:][:]$  to out_buf;
  // Sent  $O[(n_k-1)K_t : n_k K_t][:][:]$  to out_buf;
// Sent  $O[(n_k-1)K_t : n_k K_t][:][:]$  to out_buf;

```

---

(stride=1), the input tiles should contain seven valid values in each row. During the processing, a column of input features from the last tile is needed.

#### 4.4. Hybrid Dataflow

Although column-prior tiling strategy is utilized for the efficiency of the accelerator, the on-chip buffer requirement and memory accesses depend heavily on the dataflow of computations [1, 2]. The output stationary, as well as the weight stationary, is the most commonly used dataflow in previous designs. Algorithm 1 and 2 illustrate both

dataflows, respectively, where the parameters are shown in Table 2.

- **Output stationary dataflow.** Input activations and weights are fed into the PE array continuously, and the partial sums are held in PEs until the final results are available. These final results are either passed to PE\_DW for the following computation or stored in the IARAMs/OARAM. Since each output is completed after weights in a filter have been calculated, higher bandwidth is required for weight transmission. In addition, because of the implementation of weight buffer, there are more opportunities for weights to be reused.
- **Weight stationary dataflow.** Each PE holds part of weights for reuse until finishing the computation with input activations in the corresponding channels. And the partial sums generated in each PE are stored to the Inter RAM. Only if the kernel group is completed can the final results be sent to IARAMs/OARAM. In this way, weights can be reused as many as possible, but the accelerator requires additional storage, i.e., Inter RAM.

Although our accelerator is specialized for compact detection network, different convolutional layers ( $1 \times 1$  convolution and depthwise convolution) still present heterogeneous property (e.g. width, height, and channel size). The dimensions of feature maps near to the input are relatively large. Thus these layers require more on-chip buffer to store the activations, while weights require less storage. In this case, there are more opportunities for weights to be reused, which is more suitable for output stationary dataflow. However, in the deeper layers, weights become much more intensive in memory, because output stationary dataflow needs to fetch all the weights of a kernel to the PE to calculate each output. If the weight buffer can not accommodate those weights, weights are required to be fetched multiple times during the processing, leading to more energy consumption. In other words, we need a larger weight buffer to reuse weights.

Therefore, we consider a hybrid dataflow that makes a balance between the weight reuse and weight buffer requirements to get the best performance and energy on the resource-limited computing platform. In most of the early layers, we adopt the output stationary dataflow. Thus, all the weights of a kernel group can be reused in weight buffer, and they are fetched from WRAM only once during the processing of a layer. The case becomes different as the network goes deeper, and the weight stationary dataflow is adopted. So the weight buffer requirement can be significantly reduced with only a small Inter RAM overhead.

With the help of Co-Processor, our accelerator is flexible enough to support these two types of dataflow according to the size of kernels.

Table 3. Configurations of each type of PEs and the overall resource utilization on Ultra96 development board.

Parameters	$K_t$	$C_T$	Precision(A/W/O)	Operations
PE_33	8	3	8/3/4 bits	Conv $3 \times 3$
PE_11	16	32	4/3/4 bits	Conv $1 \times 1$
PE_DW	16	16	4/3/4 bits	Conv $3 \times 3$ DW
PE_Head	2	32	4/8/16 bits	Conv $1 \times 1$

Resource	Available	Used	Utilization
LUT	70560	50485	71.55%
FF	141120	74174	52.56%
BRAM	216	178.50	82.86%
DSP	360	83	23.06%

## 5. Experiments

We implement our solution on the Ultra96 development board with Xilinx Zynq UltraScale+ MPSoC. The accelerator runs at a frequency of 215 MHz with clock gating to each type of PE. Power measurement is obtained via a power monitor. We measured the power of approximate 6.9W on the Ultra96 when processing the detection task with the image size of  $512 \times 512$ . The configurations of each type of PE and the overall resource utilization are shown in Table 3, in which we also list the supported precision of activations (A), weights (W), and outputs (O) respectively. It shows that less than 25% of the total on-chip DSPs are used on the FPGA since most of the multiplications are implemented as shift operations using LUTs. Most of the registers are used as weight buffer while BRAMs are mainly used for data buffer and the WRAM. With limited programmable resources on Ultra96 board, the whole system reaches an inference speed of 18 fps. Results are reported when the system is detecting objects from a video. Table 4 shows the specification of the entire system.

Although FPGA undertakes most of the computations in detection algorithm, we find that pre-processing and post-processing on CPUs still account for most of the inference time, as shown in Figure 7 (a). In order to overcome the bottleneck of CPU execution, we adopt a pipelined task management with multi-thread techniques. In this way, the total latency is reduced, and FPGA layers dominate most of the inference time, as shown in Figure 7 (b).

Thread assignments are conducted empirically. Figure 7(c) presents the detailed time breakdown of each layer. The latency can vary greatly depending on the input image because the number of objects within an image varies significantly and thus influence the computational complexity in the post-processing phase. Therefore, time breakdown in Figure 7 is obtained by averaging over a batch of images. As shown in the figure, the softmax layer is the most

Table 4. System specification.

Device	Ultra96 development board
Network	customized MobileNet-SSD
Quantization	activation 4b/weight 3b
Power	6.9Watt
Frame rate	18 fps
Accelerator frame rate	27 fps
mAP on VOC 2012	66.4

time-consuming among all the layers, while the data forward layer and visualization layer account for 34% of the latency. Note that in a real-world application such as ADAS, the detection results are used as part of the control system, in which visualization may not be necessary. In this situation, the latency of CPUs can be further reduced, pushing the system frame rate towards the maximum.

Figure 8 shows a demo of our proposed object detection system. As we can see, the measured power is around 6.9W, and there are slight fluctuations as the detected image changes. Most of the targets are correctly detected (e.g. pedestrian, cars), frame rate for FPGA layers is around 25-30.

Table 5. Comparison with other accelerators

	VGG_ACC[10]	Low-Bit[6]	Synetgy[11]	Ours
<b>Precision (A/W)</b>	16/16 bits	2/1 bits	4/1 bits	4/3 bits
<b>Platform</b>	Zynq XC7Z045	Zynq XC7Z020	Zynq ZU3EG	Zynq ZU3EG
<b>Frequency (MHz)</b>	150	200	250	215
<b>Network</b>	VGG-16	DoReFa	ShuffleNetV2	MobileNet
<b>Classification Top_1 Acc</b>	64.64%	46.10%	68.47%	68.1%
<b>Performance (GOPs)</b>	136.97	410.22	47.09~418	202.76

As shown in Table 5, we also compare our accelerator against previous works. Since the previous works are mainly designed for image classification, we also evaluate the performance of our customized MobileNet on ImageNet classification task. Compared with VGG\_ACC, which is implemented with 16-bits integers, our design can achieve better performance and accuracy even on a smaller FPGA. Low-Bit is implemented with lower bits, which leads to severe accuracy degradation. Synetgy uses shift operations to replace the spatial convolutions. It can achieve high accuracy with lower bits, i.e., 4-bits activations and 1-bit weights. However, our accelerator can achieve more stable performance with comparable accuracy.

## 6. Conclusion

In this paper, we present a system-level solution for object detection on the heterogeneous embedded system. We quantize the compact detection network to low bits, which

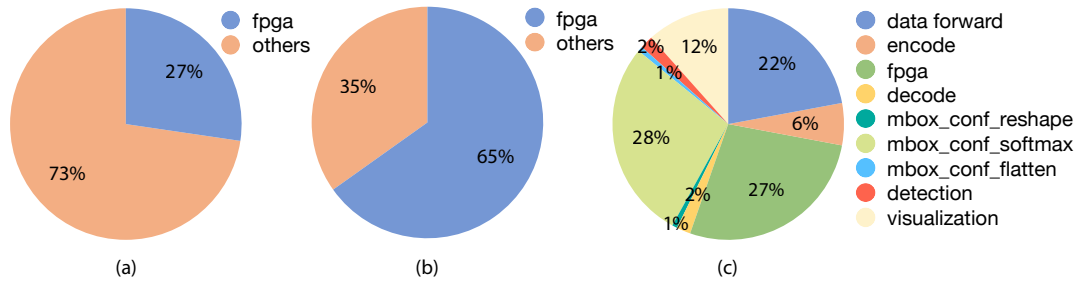


Figure 7. (a) Time breakdown before pipeline. (b) Time breakdown after pipeline. (c) Detailed time breakdown of each layer before pipeline. Figure (a) and (b) demonstrate that with the heterogeneous pipeline, the overall latency is reduced, thus the proportion of FPGA layers becomes larger, dominates most of the latency. Figure (c) shows that data forward layer and mbox-conf-softmax layer are the most time-consuming layer and require more threads to process.



Figure 8. A demo of our proposed object detection system.

allows us to replace multiplications with efficient shift operations. A dedicated CNN accelerator is implemented to carry out convolution computation. In order to support the arbitrary size of input feature maps under limited resources, we adopt a column-prior tiling strategy to map the convolutional layer to the accelerator. Compared to row-prior tiling strategy, it can reduce both register consumption and latency. According to the heterogeneous properties of different layers, we provide a hybrid dataflow, with which we can flexibly reuse the partial sums or filter weights. Multi-thread is also exploited to accelerate the pre-processing and post-processing. We believe that such an efficient and low energy system can play a role in IoT applications.

## Acknowledgment

This work was supported by the Strategic Priority Research Program of Chinese Academy of Sciences (Grant No. XDB32050200) and National Natural Science Foundation of China (Grant No.61972396, 61906193).

## References

- [1] Y.-H. Chen, J. S. Emer, and V. Sze. Eyeriss: A spatial architecture for energy-efficient dataflow for convolutional neural networks. In *ISCA*, 2016.
- [2] Y.-H. Chen, T.-J. Yang, J. S. Emer, and V. Sze. Eyeriss v2: A flexible accelerator for emerging deep neural networks on mobile devices. *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, 9:292–308, 2018.
- [3] F. Chollet. Xception: Deep learning with depthwise separable convolutions. In *CVPR*, 2016.
- [4] D. Gudovskiy and L. Rigazio. ShiftCNN: Generalized low-precision architecture for inference of convolutional neural networks. *arXiv preprint arXiv:1706.02393*, 2017.
- [5] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- [6] L. Jiao, C. Luo, W. Cao, X. Zhou, and L. Wang. Accelerating low bit-width convolutional neural networks with embedded fpga. In *FPL*, 2017.
- [7] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012.
- [8] E. H. Lee, D. Miyashita, E. Chai, B. Murmann, and S. S. Wong. Lognet: Energy-efficient neural networks using logarithmic computation. In *ICASSP*, 2017.
- [9] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. E. Reed, C.-Y. Fu, and A. C. Berg. Ssd: Single shot multibox detector. In *ECCV*, 2016.
- [10] J. Qiu, J. Wang, S. Yao, K. Guo, B. Li, E. Zhou, J. Yu, T. Tang, N. Xu, S. Song, Y. Wang, and H. Yang. Going deeper with embedded fpga platform for convolutional neural network. In *FPGA*, 2016.
- [11] Y. Yang, Q. Huang, B. Wu, T. Zhang, L. Ma, G. Gambardella, M. Blott, L. Lavagno, K. Vissers, J. Wawrzyniek, and K. Keutzer. Synetgy: Algorithm-hardware co-design for convnet accelerators on embedded fpgas. In *FPGA*, 2019.
- [12] A. Zhou, A. Yao, Y. Guo, L. Xu, and Y. Chen. Incremental network quantization: Towards lossless cnns with low-precision weights. *arXiv preprint arXiv:1702.03044*, 2017.