

[DL] A Survey of FPGA Based Neural Network Accelerator

KAIYUAN GUO, SHULIN ZENG, JINCHENG YU, YU WANG AND HUAZHONG YANG, Tsinghua University, China

Recent researches on neural network have shown significant advantage in computer vision over traditional algorithms based on handcrafted features and models. Neural network is now widely adopted in regions like image, speech and video recognition. But the high computation and storage complexity of neural network based algorithms poses great difficulty on its application. CPU platforms are hard to offer enough computation capacity. GPU platforms are the first choice for neural network process because of its high computation capacity and easy to use development frameworks.

On the other hand, FPGA based neural network accelerator is becoming a research topic. With specifically designed hardware, FPGA is the next possible solution to surpass GPU in speed and energy efficiency. Various FPGA based accelerator designs have been proposed with software and hardware optimization techniques to achieve high speed and energy efficiency. In this paper, we give an overview of previous work on neural network accelerators based on FPGA and summarize the main techniques used. An investigation from software to hardware, from circuit level to system level is carried out to complete analysis of FPGA based neural network accelerator design and serves as a guide to future work.

CCS Concepts: • **General and reference** → **Surveys and overviews**; • **Computer systems organization** → *Parallel architectures*;

Additional Key Words and Phrases: FPGA architecture, Neural Network, Parallel Processing

ACM Reference Format:

Kaiyuan Guo, Shulin Zeng, Jincheng Yu, Yu Wang AND Huazhong Yang. 2017. [DL] A Survey of FPGA Based Neural Network Accelerator. *ACM Trans. Reconfig. Technol. Syst.* 9, 4, Article 11 (December 2017), 22 pages. <https://doi.org/0000001.0000001>

1 INTRODUCTION

Recent research on Neural Network (NN) is showing great improvement over traditional algorithms in computer vision. Various network models, like convolutional neural network (CNN), recurrent neural network (RNN), have been proposed for image, video, and speech process. CNN [22] improves the top-5 image classification accuracy on ImageNet [42] dataset from 73.8% to 84.7% and further helps improve object detection [7] with its outstanding ability in feature extraction. RNN [15] achieves state-of-the-art word error rate on speech recognition. In general, NN features a high fitting ability to a wide range of pattern recognition problems. This ability makes NN a promising candidate for many artificial intelligence applications.

But the computation and storage complexity of NN models are high. Current researches on NN are still increasing the size of NN models. Take CNN as an example. The largest CNN model for a 224×224 image classification requires up to 39 billion floating point operations (FLOP) and more than 500MB model parameters [48]. As the computation complexity is proportional to the input image size, processing images with higher resolutions may need more than 100 billion operations.

Author's address: Kaiyuan Guo, Shulin Zeng, Jincheng Yu, Yu Wang AND Huazhong Yang, Tsinghua University, Tsinghua University, Beijing, Beijing, 100084, China, gky15@mails.tsinghua.edu.cn, yu-wang@mail.tsinghua.edu.cn.

ACM acknowledges that this contribution was authored or co-authored by an employee, contractor, or affiliate of the United States government. As such, the United States government retains a nonexclusive, royalty-free right to publish or reproduce this article, or to allow others to do so, for government purposes only.

© 2017 Association for Computing Machinery.

1936-7406/2017/12-ART11 \$15.00

<https://doi.org/0000001.0000001>

Therefore, choosing a proper computation platform for neural-network-based applications is essential. A typical CPU can perform 10-100G FLOP per second, and the power efficiency is usually below 1GOP/J. So CPUs are hard to meet the high performance requirements in cloud applications nor the low power requirements in mobile applications. In contrast, GPUs offer up to 10TOP/s peak performance and are good choices for high performance neural network applications. Development frameworks like Caffe [20] and Tensorflow [1] also offer easy-to-use interfaces which makes GPU the first choice of neural network acceleration.

Besides CPUs and GPUs, FPGAs are becoming a platform candidate to achieve energy efficient neural network processing. With a neural network oriented hardware design, FPGAs can implement high parallelism and make use of the properties of neural network computation to remove additional logic. Algorithm researches also show that an NN model can be simplified in a hardware-friendly way while not hurting the model accuracy. Therefore FPGAs are possible to achieve higher energy efficiency compared with CPU and GPU.

FPGA based accelerator designs are faced with two challenges in performance and flexibility:

- Current FPGAs usually support working frequency at 100-300MHz, which is much less than CPU and GPU. The FPGA's logic overhead for reconfigurability also reduces the overall system performance. A straightforward design on FPGA is hard to achieve high performance and high energy efficiency.
- Implementation of neural networks on FPGAs is much harder than that on CPUs or GPUs. Development framework like Caffe and Tensorflow for CPU and GPU is absent for FPGA.

Many designs addressing the above two problems have been carried out to implement energy efficient and flexible FPGA based neural network accelerators. In this paper, we summarize the techniques proposed in these work from the following aspects:

- We first give a simple model of FPGA based neural network accelerator performance to analyze the methodology in energy efficient design.
- We investigate current technologies for high performance and energy efficient neural network accelerator designs. We introduce the techniques in both software and hardware level and estimate the effect of these techniques.
- We compare state-of-the-art neural network accelerator designs to evaluate the techniques introduced and estimate the achievable performance of FPGA based accelerator design, which is at least 10× better energy efficient than current GPUs.
- We investigate state-of-the-art automatic design methods of FPGA based neural network accelerators.

The rest part of this paper is organized as follows: Section 2 introduces the basic operations of neural networks and the background of FPGA based NN accelerator. In section 3, we analyze the design target of NN accelerators and corresponding methods. Section 4 and section 5 review the techniques on neural network accelerator in NN model and hardware respectively. Section 6 compares existing designs and evaluate the techniques. Section 7 introduces the methods for a flexible accelerator design. Section 8 concludes this paper.

2 PRELIMINARY

Before discussing the system design for neural network acceleration, we first introduce the basic concepts of neural networks and the typical structure of FPGA-based NN accelerator design.

2.1 Neural Network

In this section, we introduce the basic functions in a neural network. In this paper, we only focus on the inference of NN, which means using a trained model to predict or classify new data. The

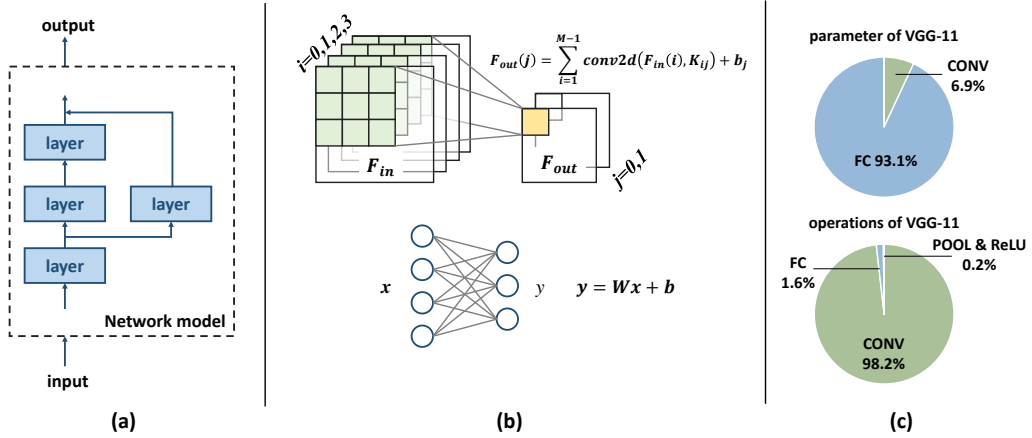


Fig. 1. (a) Computation graph of a neural network model. (b) CONV and FC layers in NN model. (c) CONV and FC layers dominate the computation and parameter of an NN model.

training process of NN is not discussed in this paper. A neural network model can be expressed as a directed graph shown in Figure 1(a). Each vertex of the graph denotes a layer which conducts operations on data from a previous layer or input and generates results to the next layer or output. We refer the parameter of each layer as weights and the input/output of each layer as activations through this paper.

Convolution (CONV) layers and fully connected (FC) layers are two common types of layers in NN models. The functions of these two layers are shown in Figure 1(b). CONV layers conduct 2D convolutions on a set of input feature maps F_{in} and add the results to get output feature maps F_{out} . FC layers receive a feature vector as input and conduct matrix-vector multiplications.

Besides CONV and FC layers, NN layers also have pooling, ReLU [22], concat [50], element-wise [16] and other types of layers. But these layers contributes little to the computation and storage requirement of a neural network model. Figure 1(c) shows the distribution of weights and operations in the VGG-11 model [48]. CONV and FC layers together contribute more than 99% of the network's weights and operations. So most of the neural network acceleration systems focus on these two types of layers.

2.2 FPGA-based Accelerator

In recent years, FPGA is becoming a promising solution for algorithm acceleration. Compared with CPU, GPU, and DSP platforms, for which the software and hardware are designed independently, FPGA enables the developers to implement only the necessary logic in hardware according to the target algorithm. By eliminating the redundancy in general hardware platforms, FPGAs can achieve higher efficiency. Application specific integrated circuits (ASICs) based solutions achieve even higher efficiency but requires much longer development cycle and higher cost.

For FPGA-based neural network accelerator, a typical architecture of the system is shown in Figure 2(a). The system usually consists of a CPU host and an FPGA part. A pure FPGA chip usually works with a host PC/server through PCIe connections. SoC platforms (like the Xilinx Zynq Series) and Intel HARPv2 [12] platform integrate the host and the FPGA in the same chip or package. Both the host and the FPGA can work with their own external memory and access each others'

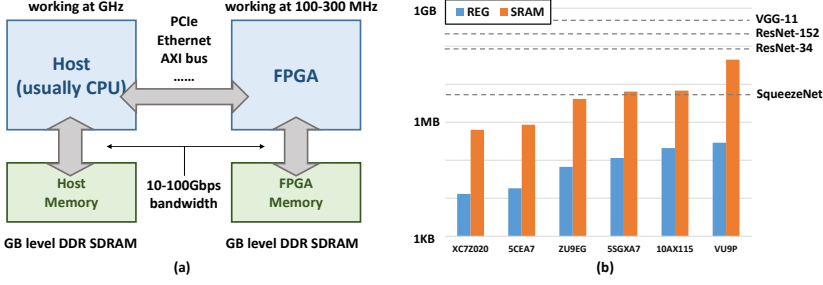


Fig. 2. (a) A typical structure of an FPGA based NN accelerator. (b) Comparison between NN model size and the storage on different FPGA chips.

memory through the connection. Most of the designs implement NN accelerator on the FPGA part and control the accelerator with the software on the host.

Typical FPGA chips consist large on-chip storage units like registers and SRAMs, but still too small compared with NN models as shown in Figure 2(b). Common models implement 100-1000MB parameters while the largest available FPGA chip implements <50MB on-chip SRAM. This gap requires that external memory like DDR SDRAM is needed. The bandwidth and power consumption of DDR limits the system performance.

The computation capacity of FPGA is relatively higher. Common FPGAs implements hundreds to thousands of DSP units, each of which can compute 18×27 or 18×19 , achieving up to 10TFLOP/s (floating point operations per second) on the largest FPGAs. But for low-end FPGAs like Xilinx XC7Z020, this number is reduced to 20GFLOP/s, which is hard to support real-time video processing for applications on mobile platforms.

Even faced with the above challenges, researchers have proposed a series of optimization methods from algorithm to architecture to design high performance NN accelerators on FPGA, which will be discussed in the following sections of this paper.

3 DESIGN METHODOLOGY AND CRITERIA

Before going into the details of the techniques used for neural network accelerators, we first give an overview of the design methodology. In general, the design target of a neural network processing system includes the following two aspects: high speed (high throughput and lower latency), and high energy efficiency.

Speed. The throughput of an NN accelerator can be expressed by equation 1. Peak performance, measured in operations (multiplication or addition) per second, is achieved when all the computation units work every clock cycle. Utilization denotes the average ratio of working cycles of the computation units. The workload measures the number of operations in the target neural network. The on-chip resource for a certain FPGA chip is limited. We can increase the peak performance by reducing the size of each computation unit and increasing the working frequency. Reducing the size of computation units can be achieved by simplifying the basic operations in a neural network model, which may hurt the model accuracy and requires hardware-software co-design. On the other hand, increasing working frequency is pure hardware design work. A high utilization ratio is ensured by reasonable parallelism implementation and efficient memory system. Most of this part is affected by hardware design. But optimization of NN models can also reduce the storage requirements of a neural network model and benefits the memory system.

$$throughput = \frac{actual_performance}{workload} = \frac{peak_performance \times utilization}{workload} \quad (1)$$

Most of the FPGA based NN accelerators compute different inputs one at a time, which means the latency of the system is the reciprocal of throughput. The concurrent process of different inputs is usually realized with batch [28] or layer pipeline [24]. In these cases, the concurrency equals to the batch size or the number of pipeline stages. The latency of the accelerator is expressed as equation 2.

$$latency = \frac{concurrency}{throughput} \quad (2)$$

So in this paper, we focus more on optimizing the throughput. As different accelerators may be evaluated on different NN models, a common criterion of speed is the *actual_performance*, which eliminates the effect of different workload to some extent.

Energy Efficiency. Energy efficiency is the actual performance within a certain power cost, i.e. the average number of operations can be executed within a certain energy budget. The energy efficiency is measured in (OP/s/W) or (OP/J). The workload for the target network is fixed. Increasing the energy efficiency of a neural network accelerator means to reduce the total energy cost, E_{total} to process each input as shown in equation 3. This energy cost comes from 2 parts: computation and memory access, which is expressed in equation 4.

$$energy_efficiency = \frac{E_{total}}{workload} \quad (3)$$

$$E_{total} \approx N_{op} \times E_{op} + N_{SRAM_access} \times E_{SRAM_access} + N_{DRAM_access} \times E_{DRAM_access} + E_{static} \quad (4)$$

The first item in equation 4 is the dynamic energy cost for computation. N_{op} denotes the number of operations processed. E_{op} denotes the average energy cost of each operation. Given a certain network and target FPGA platform, both N_{op} and E_{op} are fixed. For this part, researchers have been focusing on optimizing the NN models by quantization (narrowing the bit-width used for computation) to reduce E_{op} or sparsification (setting more weights to zeros) to skip the multiplications with these zeros to reduce N_{op} .

The second and third item in equation 4 is the dynamic energy cost for memory access. As shown in section 2.2, FPGA based NN accelerator usually works with an external DRAM, we separate the memory access energy into DRAM part and SRAM part. N_{xx_access} denotes the total number of bytes to be accessed from memory. E_{xx_access} denotes the energy cost for accessing each byte. This part can also be reduced by quantization and sparsification, which helps reduce N_{xx_access} . Efficient on-chip memory system and scheduling method also help to reduce N_{DRAM} . E_{xx_access} can hardly be reduced given a certain FPGA platform.

The fourth item E_{static} denotes the static energy cost of the system. This energy cost can hardly be improved given the FPGA chip and the scale of the design.

From the analysis of speed and energy, we see that neural network accelerator involves both optimizations on NN models and hardware. In the following sections, we will introduce previous work in these two aspects respectively.

4 HARDWARE ORIENTED MODEL COMPRESSION

As introduced in section 3, the design of energy efficient and fast neural network accelerator can benefit from the optimization of NN models. A larger NN model usually results in higher model accuracy. This means it is possible to trade the model accuracy for the hardware speed or energy

cost. Neural network researchers are designing more efficient network models from AlexNet [22] to ResNet [16], SqueezeNet [19] and MobileNet [18]. The main differences between these networks are the size of and the connections between each layer. The basic operations are the same and hardly affect the hardware design. For this reason, we will not focus on these techniques in this paper. Other methods try to achieve the tradeoff by compressing existing NN models. They try to reduce the number of weights or reduce the number of bits used for each activation or weight, which help lower down the computation and storage complexity. Corresponding hardware designs can benefit from these NN model compression methods. In this section, we investigate these hardware oriented network model compression methods.

4.1 Data Quantization

One of the most commonly used methods for model compression is the quantization of the weights and activations. The activations and weights of a neural network are usually represented by floating point data in common developing frameworks. Recent work tries to replace this representation with low-bit fixed-point data or even a small set of trained values. On the one hand, using fewer bits for each activation or weight helps reduce the bandwidth and storage requirement of the neural network processing system. On the other hand, using a simplified representation reduce the hardware cost for each operation. The benefit of hardware will be discussed in detail in section 5. Two kinds of quantization methods are discussed in this section: linear quantization and non-linear quantization.

4.1.1 Linear Quantization. Linear quantization finds the nearest fixed-point representation of each weight and activation. The problem with this method is that the dynamic range of floating-point data greatly exceeds that for fixed-point data. Most of the weights and activations will suffer from overflow or underflow. Qiu et al. [41] finds that the dynamic range of the weights and activations in a single layer is much more limited and differs across different layers. Therefore they assign different fractional bit-widths to the weights and activations in different layers. To decide the fractional bit-width of a set of data, i.e. the activations or weights of a layer, the data distribution is first analyzed. A set of possible fractional bit-widths are chosen as candidate solutions. Then the solution with the best model performance on training data set is chosen. In [41], the optimized solution of a network is chosen layer by layer to avoid an exponential design space exploration. Guo et al. [11] further improve this method by fine-tuning the model after the fraction bit-width of all the layers are fixed.

The method of choosing a fractional bit-width equals to scale the data with a scaling factor of 2^k . Li et al. [23] scales the weights with trained parameter W^l for each layer and quantize the weights with 2-bit data, representing W^l , 0 and $-W^l$. The activations in this work are not quantized. So the network still implements 32-bit floating point operations. Zhou et al. [67] further quantize the weights of a layer with only 1 bit to $\pm s$, where $s = E(|w^l|)$ is the expectation of the absolute value of the weights of this layer. Linear quantization is also applied to the activations in this work.

4.1.2 Non-linear Quantization. Compared with linear quantization, non-linear quantization independently assigns values to different binary codes. The translation from a non-linear quantized code to its corresponding value is thus a look-up table. This kind of methods helps further reduce the bit-width used for each activation or weight. Chen et al. [4] assign each of the weight to an item in the look-up table by a pre-defined hash function and train the values in look-up tables. Han et al. [14] assign the values in look-up tables to the weights by clustering the weights of a trained model. Each look-up table value is set as the cluster centre and further fine-tuned with training data set. This method can compress the weights of state-of-the-art CNN models to 4-bit without accuracy loss. Zhu et al. [68] propose the ternary-quantized network where all the weights of a layer

are quantized to three values: W^n , 0, and W^p . Both the quantized value and the correspondence between weights and look-up table are trained. This method sacrifices less than 2% accuracy loss on ImageNet dataset on state-of-the-art network models. The weight bit-width is reduced from 32-bit to 2-bit, which means about 16× model size compression.

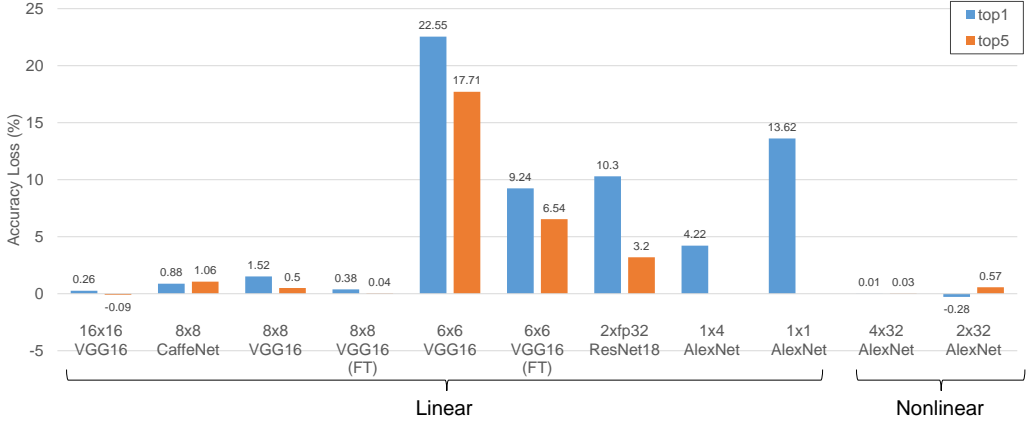


Fig. 3. Comparison between different quantization methods from [11, 14, 23, 41, 67, 68]. The quantization configuration is expressed as (weight bit-width)×(activation bit-width). The "(FT)" denotes that the network is fine-tuned after a linear quantization.

4.1.3 Comparison. We compare some typical quantization methods from [11, 14, 23, 41, 67, 68] in Figure 3. All the quantization results are tested on ImageNet data set and the absolute accuracy loss compared with corresponding baseline floating point models is recorded. Comparing different methods on different models is a little bit unfair. But it still gives some insights. For linear quantization, 8-bit is a clear bound to ensure negligible accuracy loss. With 6 or fewer bits, using fine-tune or even training each weight from the beginning will cause noticeable accuracy degradation. If we require that 1% accuracy loss is within the acceptable range, linear quantization with at least 8×8 configuration and the listed non-linear quantization are available. We will further discuss the performance gain of quantization in section 5.

4.2 Weight Reduction

Besides narrowing the bit-width of activations and weights, another method for model compression is to reduce the number of weights. One kind of method is to approximate the weight matrix with a low-rank representation. Qiu et al. [41] compress the weight matrix W of an FC layer with singular value decomposition. An $m \times n$ weight matrix W is replaced by the multiplication of two matrices $A_{m \times p} B_{p \times n}$. For a sufficiently small p , the total number of weights is reduced. This work compresses the largest FC layer of VGG network to 36% of its original size with 0.04% classification accuracy degradation. Zhang et al. [65] use a similar method for convolution layers and takes the effect of the following non-linear layer into the decomposition optimization process. The proposed method achieves 4× speed up on state-of-the-art CNN model targeting at ImageNet, with only 0.9% accuracy loss.

Pruning is another kind of method to reduce the number of weights. This kind of methods directly remove the zeros in weights or remove those with small absolute values. The challenge in pruning is the tradeoff between the ratio of zero weights and the model accuracy. One solution is

the application of lasso method, which applies L1 normalization to the weights during training. Liu et al. [26] apply the sparse group-lasso method on the AlexNet [22] model. 90% weights are removed after training with less than 1% accuracy loss. Another solution is to prune the zero weights during training. Han et al. [14] directly remove the weights of a network which are zero or have small absolute value. The left weights are then fine-tuned with the training dataset to recover accuracy. Experimental results on AlexNet show that 89% weights can be removed while keeping the model accuracy.

The hardware gain from weight reduction is the reciprocal of the compression ratio. According to the above results, the improvement from weight reduction is up to 10 \times .

5 HARDWARE DESIGN: EFFICIENT ARCHITECTURE

In this section, we investigate the hardware level techniques used in state-of-the-art FPGA based neural network accelerator design to achieve high performance and high energy efficiency. We classify the techniques into three levels: computation unit level, loop unrolling level, and system level.

5.1 Computation Unit Designs

Computation unit level design affects the peak performance of the neural network accelerator. The available resource of an FPGA chip is limited. A smaller computation unit design means more computation units and higher peak performance. A carefully designed computation unit array can also increase the working frequency of the system and thus improve peak performance.

5.1.1 Low Bit-width Unit. Reduce the number of bit-width for computation is a direct way to reduce the size of computation units. The feasibility of using fewer bits comes from the quantization methods as introduced in section 4.1. Most of the state-of-the-art FPGA designs replace the 32-bit floating-point units with fixed-point units. Podili et al. [39] implement 32-bit fixed-point units for the proposed system. 16-bit fixed-point units are widely adopted in [8, 24, 41, 58, 60]. ESE [13] adopts 12-bit fixed-point weight and 16-bit fixed-point neurons design. Guo et al. [11] use 8-bit units for their design on embedded FPGA. Recent work is also focusing on extremely narrow bit-width design. Prost-Boucle et al. [40] implements 2-bit multiplication with 1 LUT for ternary networks. Experiments in [38] show that FPGA implementation of Binarized Neural Network (BNN) outperforms that on CPU and GPU. Though BNN suffers from accuracy loss, many designs explore the benefit of using 1-bit data for computation [21, 25, 33, 35, 36, 51, 66].

The designs mentioned above focus on computation units for linear quantization. For non-linear quantization, translating the data back to full precision for computation still costs many resources. Samragh et al. [43] propose the factorized coefficients based dot product implementation. As the possible values of weights are quite limited for non-linear quantization, the proposed computation unit accumulates the multiplier for each possible weight value and calculate the result as the weighted sum of the values in look-up tables. In this way, the multiplication needed for one output neuron equals to the number of values in look-up table. The multiplications are replaced by random-addressed accumulations.

Most of the designs use one bit-width through the process of a neural network. Qiu et al. [41] finds that neurons and weights in FC layers can use fewer bits compared with CONV layers while the accuracy is maintained. Heterogeneous computation units are used in the designs of [10, 66].

The size of computation units of different bit-widths is compared in Table 1. Three kinds of implementations are tested: separate multiplier and adder with logic resource on Xilinx FPGA, multiply-add function with DSP units on Xilinx FPGA, and multiply-add function with DSP units on Altera FPGA. The resource consumption is the synthesis result by Vivado 2018.1 targeting

Xilinx XCKU060 FPGA and Quartus Prime 16.0 targeting Altera Arria 10 GX1150 FPGA. The pure logic modules and the floating-point multiply and add modules are generated with IP core. The fixed-point multiply and add modules are implemented with $A * B + C$ in Verilog and automatically mapped to DSP by Vivado/Quartus.

We first give an overview of the size of the computation units by logic-only implementations. By compressing the weights and activations from 32-bit floating-point number to 8-bit fixed-point number, the multiplier and the adder are scaled down to about 1/10 and 1/50 respectively. Using 4-bit or smaller operators can bring further advantage but also incur significant accuracy loss as introduced in section 4.1.

Recent FPGAs consist of a large number of DSP units, each of which implements hard multiplier, pre-adder and accumulator core. The basic pattern of NN computation, multiplication and sum, also fits into this design. So we also test the multiply and add function implemented with DSP units. Because of the different DSP architectures, we test on both Xilinx and Altera platforms. Compared with the 32-bit floating-point function, fixed-point functions with narrow bit-width still shows an advantage in resource consumption. But for Altera FPGA, this advantage is not obvious because the DSP units natively support floating-point operations.

Fixed-point functions with 16-or-less-bit fixed-point data are well fit into 1 DSP unit on either Xilinx or Altera FPGA. This shows that quantization hardly benefits the hardware if we use narrower bit-width like 8 or 4 in the aspect of computation. The problem is that the wide multipliers and adders in DSP units are underutilized in these cases. Nguyen et al. [37] propose the design to implement two narrow bit-width fixed-point multiplication with a single wide bit-width fixed-point multiplier. In this design, two multiplications, AB and AC , are executed in the form of $A(B \ll k + C)$. If k is sufficiently large, the bits for AB and AC does not overlap in the multiplication result and can be directly separated. The design in [37] implements two 8-bit multiplications with one 25×18 multiplier, where k is 9. Similar methods can be applied to other bit-width and DSPs.

Table 1. FPGA resource consumption comparison for multiplier and adder with different types of data.

	Xilinx Logic				Xilinx DSP			Altera DSP	
	multiplier		adder		multiply & add			multiply & add	
	LUT	FF	LUT	FF	LUT	FF	DSP	ALM	DSP
fp32	708	858	430	749	800	1284	2	1	1
fp16	221	303	211	337	451	686	1	213	1
fixed32	1112	1143	32	32	111	64	4	64	3
fixed16	289	301	16	16	0	0	1	0	1
fixed8	75	80	8	8	0	0	1	0	1
fixed4	17	20	4	4	0	0	1	0	1

5.1.2 Fast Convolution Unit. For CONV layers, the convolution operations can be accelerated by alternative algorithms. Discrete Fourier Transformation (DFT) based fast convolution is widely adopted in digital signal processing. Zhang et al. [62] propose a 2D DFT based hardware design for efficient CONV layer execution. For an $F \times F$ filter convolved with $K \times K$ filter, DFT converts the $(F - K + 1)^2 K^2$ multiplications in the space domain to F^2 complex multiplications in the frequency domain. For a CONV layer with M input channel and N output channel, MN times of frequency domain multiplications and $(M + N)$ times DFT/IDFT are needed. The conversion of convolution kernels is once for all. So the domain conversion process is of low cost for CONV layers. This technique does not work for CONV layers with stride>1 or 1×1 convolution. Ding et al. [6]

suggest that a block-wise circular constraint can be applied to the weight matrix. In this way, the matrix-vector multiplication in FC layers are converted to a set of 1D convolutions and can be accelerated in the frequency domain. This method can also be applied to CONV layers by treating the $K \times K$ convolution kernels as $K \times K$ matrices and is not limited by K or stride.

Frequency domain methods require complex number multiplication. Another kind of fast convolution involves only real number multiplication [56]. The convolution of a 2D feature map F_{in} with a kernel K using Winograd algorithm is expressed by equation 5.

$$F_{out} = A^T[(GF_{in}G^T) \odot (BF_{in}B^T)]A \quad (5)$$

G , B and A are transformation matrix which only related to the sizes of kernel and feature map. \odot denotes an element-wise multiplication of two matrices. For a 4×4 feature map convolved with a 3×3 kernel, the transformation matrices are described as follows:

$$G = \begin{bmatrix} 1 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} \\ 0 & 0 & 1 \end{bmatrix} \quad B = \begin{bmatrix} 1 & 0 & -1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & -1 & 1 & 0 \\ 0 & 1 & 0 & -1 \end{bmatrix} \quad A = \begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 1 & -1 \\ 0 & -1 \end{bmatrix}$$

Winograd based methods are also limited by the kernel size and stride as DFT based methods. The most commonly used Winograd transformation is for 3×3 convolutions in [28, 58].

The theoretical performance gain from fast convolution depends on the transformation kernel size. As DFT based method uses complex multiplication and costs more hardware, we estimate the theoretical benefit according to the 6×6 Winograd used in [28], which is $4\times$.

5.1.3 Frequency Optimization Methods. All the above techniques introduced targets at increasing the number of computation units within a certain FPGA. Increasing the working frequency of the computation units also improves the peak performance.

Latest FPGAs support 700-900MHz DSP theoretical peak working frequency. But existing designs usually work at 100-300MHz [11, 30, 41, 60]. As claimed in [57], the working frequency is limited by the routing between on-chip SRAM and DSP units. The design in [57] uses different working frequencies for DSP units and surrounding logic. Neighbour slices to each DSP unit are used as local RAMs to separate the clock domain. The prototype design in [57] achieves the peak DSP working frequency at 741MHz and 891MHz on FPGA chips of different speed grades. But this method is not adopted by a complete neural network accelerator design. Existing designs achieve 300MHz without this technique, so we estimate the theoretical hardware performance gain as $2\times$.

5.2 Loop Unrolling Strategies

CONV layers and FC layers contribute to most of the computations and storage requirement of a neural network as introduced in section 2. We express the CONV layer function in Figure 1(b) as nested loops in Algorithm 1. To make the code clear to read, we merge the loops along x and y directions for feature maps and 2-D convolution kernels respectively. An FC layer can be expressed as a CONV layer with feature map and kernel both of size 1×1 . Besides the loops in Algorithm 1, we also call the parallelism of the process of multiple inputs as a batch. This forms the batch loop.

5.2.1 Choosing Unroll Parameters. To parallelize the execution of the loops, we unroll the loops and parallelize the process of a certain number of iterations on hardware. The number of the parallelized iterations on hardware is called the unroll parameter. Inappropriate unroll parameter selection may lead to serious hardware underutilization. Take a single loop as an example. Suppose the trip count of the loop is M and the parallelism is m . The utilization ratio of the hardware is limited by $m/M \lceil M/m \rceil$. If M is not divisible by m , then the utilization ratio is less than 1. For

Algorithm 1 Convolution Layer

Require: feature map F_{in} of size $M \times Y \times X$; convolution kernel Ker of size $N \times M \times K \times K$; bias vector b of size N

Ensure: feature map F_{out}

```

1: function CONV_LAYER( $F_{in}, Ker$ )
2:   Let  $F_{out} \leftarrow$  zero array of size  $N \times (Y - K + 1) \times (X - K + 1)$ 
3:   for  $n = 1; n < N; n++$  do                                ▶ Output channel loop
4:     for  $m = 1; m < M; m++$  do                                ▶ Input channel loop
5:       for each  $(y, x)$  within  $(Y - K + 1, X - K + 1)$  do    ▶ Feature map loop
6:         for each  $(ky, kx)$  within  $(K, K)$  do                ▶ Kernel loop
7:            $F_{out}[n][y][x] += F_{in}[m][y - ky + 1][x - kx + 1] * K[n][m][ky][kx]$ 
8:        $F_{out}[n] += b[n]$ 
9:   return  $F_{out}$ 

```

processing an NN layer, the total utilization ratio will be the product of the utilization ratio on each of the loops.

For a CNN model, the loop dimension varies greatly among different layers. For a typical network used on ImageNet classification like ResNet [16], the channel numbers vary from 3 to 2048; the feature map sizes vary from 224×224 to 7×7 , the convolution kernel sizes vary from 7×7 to 1×1 . Besides the underutilization problem, loop unrolling also affect the datapath and on-chip memory design. Thus loop unrolling strategy is a key feature for a neural network accelerator design.

Various work are proposed focusing on how to choose the unroll parameters. Zhang et al. [61] propose the idea of unrolling the input channel and output channel loops and choose the optimized unroll parameter by design space exploration. Along these two loops, there is no input data cross-dependency between neighbouring iterations. So no multiplexer is needed to route data from the on-chip buffer to computation units. But the parallelism is limited as $7 \times 64 = 448$ multipliers. For larger parallelism, this solution is easy to suffer from the underutilization problem. Ma et al. [30] further extends the design space by allowing parallelism on the feature map loop. The parallelism reaches $1 \times 16 \times 14 \times 14 = 3136$ multipliers. A shift register structure is used to route feature map pixels to the computation units.

The kernel loop is not chosen in the above work because kernel sizes vary greatly. Motamedi et al [34] use kernel unrolling on AlexNet. Even with 3×3 unrolling for the 11×11 and 5×5 kernels, the overall system performance still reaches 97.4% of its peak performance for the convolution layers. For certain networks like VGG [48], only 3×3 convolution kernels are used. Another reason to unroll kernel loop is to achieve acceleration with fast convolution algorithms. Design in [62] implements fully parallelized frequency domain multiplication on 4×4 feature map and 3×3 kernel. Lu et al. [28] implement Winograd algorithm on FPGA with a dedicated pipeline for equation 5. The convolution of a 6×6 feature map with a 3×3 kernel is fully parallelized.

The above solutions are only for a single layer. But there is hardly a one-size-fits-all solution for a whole network, especially when we need high parallelism. Designs in [24, 27] propose fully pipelined structures with each layer a pipe stage. As each layer is executed with an independent part of the hardware and each part is small, loop unrolling method can be easily chosen. This method is memory consuming because ping-pong buffers are needed between adjacent layers for the feature maps. Design in [63] is similar but implemented on FPGA clusters to resolve the scalability problem. Shen et al. [46] group the layers of a CNN by the loops' trip count and map each group onto one hardware module. These solutions can be treated as unrolling the batch loop

because different inputs are processed in parallel on different layer pipeline stages. The design in [28] implements parallelized batch both within a layer and among different layers.

Most of the current designs follow one of the above methods for loop unrolling. A special kind of design is for sparse neural networks. Han et al. [13] propose the ESE architecture for sparse LSTM network acceleration. Unlike processing a dense network, all the computation units will not work synchronously. This causes difficulty in sharing data between different computation units. ESE implements only the output channel (the output neurons of the FC layers in LSTM) loop unrolling within a layer to simplify hardware design and parallelize batch process.

5.2.2 Data Transfer and On-chip Memory Design. Besides the high parallelism, the on-chip memory system should efficiently offer the necessary data to each computation units every cycle. To implement high parallelism, neural network accelerators usually reuse data among a large number of computation units. Simply broadcasting data to different computation units leads to large fan-out and high routing cost and thus reduce the working frequency. Wei et al. [55] use the systolic array structure in their design. The shared data are transferred from one computation unit to the next in a chain mode. So the data is not broadcasted, and only local connections between different computation units are needed. The drawback is the increase in latency. The loop execution order is scheduled accordingly to cover the latency. Similar designs are adopted in [3, 30].

For software implementation on GPU, the `im2col` function is commonly used to map 2D convolution as a matrix-vector multiplication. This method incurs considerable data redundancy and can hardly be applied to the limited on-chip memory of FPGAs. Qiu et al. [41] uses the line buffer design to achieve the 3×3 sliding window function for 2-d convolution with only two lines of duplicated pixels.

5.3 System Design

A typical FPGA based neural network accelerator system is shown in Figure 4. The logic part of the whole system is denoted by the blue boxes. The host CPU issues workload or commands to the FPGA logic part and monitors its working status. On the FPGA logic part, a controller is usually implemented to communicate with the host and generates control signals to all the other modules on FPGA. The controller can be an FSM or an instruction decoder. The on the fly logic part is implemented for certain designs if the data loaded from external memory needs preprocess. This module can be data arrangement module, data shifter [41], FFT module [62], etc. The computation units are as discussed in section 5.1 and section 5.2. As introduced in section 2.2, on-chip SRAM of an FPGA chip is too limited compared with the large NN models. So for common designs, a two-level memory hierarchy is used with DDR and on-chip memory.

From the system level, the performance of a neural network accelerator is limited by two factors: the on-chip computation resource and the off-chip memory bandwidth. Various researches have been proposed to achieve the best performance within a certain off-chip memory bandwidth. Zhang et al. [61] introduce the roofline model in their work to analyze whether a design is memory bounded or computation bounded. An example of a roofline model is shown in Figure 5.

The figure uses the computation to communication (CTC) ratio as the x -axis and hardware performance as the y -axis. CTC is the number of operations that can be executed with a unit size of memory access. Each hardware design can be treated as a point in the figure. So y/x equals to the bandwidth requirement of the design. The available bandwidth of a target platform is limited and can be described as the theoretical bandwidth roof in Figure 5. But the actual bandwidth roof is below the theoretical roof because the achievable bandwidth of DDR depends on the data access pattern. Sequential DDR access achieves much higher bandwidth than random access. The other roof is the computation roof, which is limited by the available resource on FPGA.

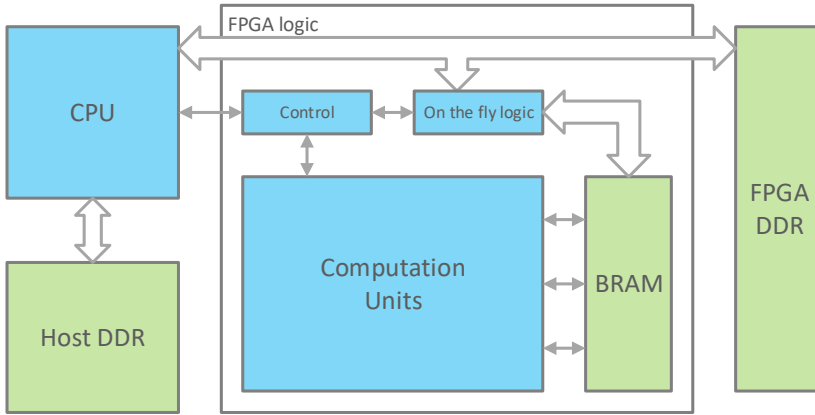


Fig. 4. Block graph of a typical FPGA based neural network accelerator system

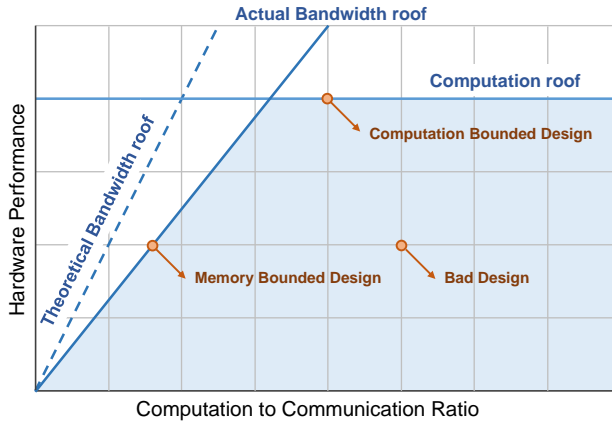


Fig. 5. An example of the roofline model. The shaded part denotes the valid design space given bandwidth and resource limitation.

So a higher CTC ratio means the hardware is more likely to achieve the computation bound. Increasing the CTC ratio also reduce DDR access, which significantly saves energy according to [17]. In section 5.2, we have discussed the loop unrolling strategies to increase the parallelism while reducing the waste of computation for a certain network. When the loop unrolling strategy is decided, the scheduling of the rest part of the loops decides how the hardware can reuse data with on-chip buffer. This involves loop tiling and loop interchange strategy.

Loop tiling is a higher level of loop unrolling. All the input data of a loop tile will be stored on-chip, and the loop unrolling hardware kernel works on these data. A larger loop tile size means that each tile will be loaded from external memory to on-chip memory fewer times. Loop interchange strategy decides the processing order of the loop tiles. External memory access happens when the hardware is moving from one tile to the next. Neighbouring tile may share a part of data. For example in a CONV layer, neighbouring tile can share input feature map or the weights. This is decided by the execution order of the loops.

In [30, 61], design space exploration is done on all the possible loop tiling sizes and loop orders. Many designs also explore the design space with some of the loop unrolling, tiling and loop order is already decided [34, 41]. Shen et al. [47] also discuss the effect of batch parallelism over the CTC for different layers. This is a loop dimension not focused on in previous work.

All the above work give one optimized loop unrolling strategy and loop order for a whole network. Guo et al. [11] implements flexible unrolling and loop order configuration for different layers with an instruction interface. The data arrangement in on-chip buffers is controlled through instructions to fit with different feature map sizes. This means the hardware can always fully utilize the on-chip buffer to use the largest tiling size according to on-chip buffer size. This work also proposes the "back and forth" loop execution order to avoid total on-chip data refresh when an innermost loop finishes.

Alwani et al. [2] address the external memory access problem by fusing two neighbouring layers together to avoid the intermediate result transfer between the two layers. This strategy helps reduce 95% off-chip data transfer with extra 20% on-chip memory cost. Even software program gains $2\times$ speedup with this scheduling strategy. Yu et al. [59] realize this idea on a single-layer accelerator design by modifying the order of execution through an instruction interface.

Besides increasing CTC, increasing the actual bandwidth roof helps improve the achievable performance with a certain CTC ratio. This is achieved by regularizing the DDR access pattern. The common feature map formats in the external memory include *NCHW* or *CHWN*, where *N* means the batch dimension, *C* means the channel dimension, *H* and *W* means the feature map *y* and *x* dimension. Using any of these formats, a feature map tile may be cut into small data blocks stored in discontinuous addresses. Guan [8] suggest that a channel-major storage format should be used for their design. This format avoids data duplication while long DDR access burst is ensured. Qiu et al. [41] propose a feature map storage format that arranges the $H \times W$ feature map into (HW/rc) tile blocks of size $r \times c$. So the write burst size can be increased from $c/2$ to $rc/2$.

6 EVALUATION

In this section, we compare the performance of state-of-the-art neural network accelerator designs and try to evaluate the techniques mentioned in section 4 and section 5. We mainly reviewed the FPGA-based designs published in the top FPGA conferences (FPGA, FCCM, FPL, FPT), EDA conferences (DAC, ASPDAC, DATE, ICCAD), architecture conferences (MICRO, HPCA, ISCA, ASPLOS) since 2015. Because of the diversity in the adopted techniques, target FPGA chips, and experiments, we need a trade-off between the fairness of comparison and the number of designs we can use. In this paper, we pick the designs with 1) whole system implementation; 2) experiments on real NN models with reported speed, power, and energy efficiency.

The designs used for comparison are listed in Table 2. For data format, the "INT A/B" means that activations are A-bit fixed-point data and weights are B-bit fixed-point data. We also investigate the resource utilization and draw advice to both accelerator designers and FPGA manufacturers.

Each of the designs in Table 2 drawn as a point in Figure 6, using $\log_{10}(\text{power})$ as *x* coordinate and $\log_{10}(\text{speed})$ as *y*-axis. Therefore, $y - x = \log_{10}(\text{energy_efficiency})$. Besides the FPGA based designs, we also plot the GPU experimental results used in [11, 13] as standards to measure the FPGA designs' performance.

Bit-width Reduction. Among all the designs, 1-2 bit based designs [21, 33, 35] show outstanding speed and energy efficiency. This shows that extremely low bit-width is a promising solution for high-performance design. As introduced in section 4.1, linear quantized 1-2 bit network models suffer from great accuracy loss. Further developing related accelerator will be of little use. More

Table 2. Performance and resource utilization of state-of-the-art neural network accelerator designs

	Data Format	Speed (GOP/s)	Power (W)	Efficiency (GOP/J)	Resource(%)			FPGA chip
					DSP	logic	BRAM	
[3]	FP16	1382	45	30.7	97	58	92	Arria 10 GX1150
[13]	INT16/12	2520	41	61.5	54.4	88.6	87.7	XCKU060
[52]	INT16	12.73	1.75	7.27	94.54	66.64	6.07	XC7Z020
[62]	FP32	123.5	13.18	9.37	87.5	85.4	64	Stratix V
[64]	INT16	1790	37.46	47.8	91	43	53	GX1150
	FP32	866	41.73	20.75	87	-	46	
[30]	INT16/8	645.25	21.2	30.43	100	38	70	GX1150
[41]	INT16	136.97	9.63	14.22	89.2	83.5	86.7	XC7Z045
[49]	INT16/8	117.8	19.1	6.2	12.5	22	65.2	5SGSD8
[61]	FP32	61.62	18.61	3.3	80	61.3	50	XC7VX485T
[8]	INT16	364.4	25	14.6	65	25	46	5SGSMD5
[28]	INT16	2940.7	23.6	124.6	-	-	-	ZCU102
[39]	INT32	229	8.04	28.5	100	83.7	17.6	Stratix V
[35]	1bit	329.47	2.3	143.2	0.5	34.4	11.4	Zynq XC7Z020
[21]	2bit	410.22	2.26	181.51	40.5	82.7	37.7	Zynq XC7Z020
[33]	1bit	40770	48	849.38	-	-	-	GX1155
[24]	INT16	565.94	30.2	22.15	59.56	63.21	65.07	XC7VX690T
[27]	INT16/8	222.1	24.8	8.96	39.9	26.6	39.7	XC7VX690T
[63]	INT16	1280.3	160	8	-	-	-	XC7Z020+ XC7VX690T×6
[11]	INT8	84.3	3.5	24.1	87	84	89	XC7Z020
[58]	INT16	229.5	9.4	24.42	91.9	71	83.2	XC7Z045
[9]	FP32	7.26	19.63	0.37	42	65.31	52.04	XC7VX485T
[60]	INT16	354	26	13.6	78	81	42	XC7VX690T
[45]	INT16	431	25	17.1	42	56	52	XC7VX690T
		785	26	30.2	53	8.3	30	XCVU440

efforts should be put on the models. Even trading speed with accuracy can be acceptable considering the current hardware performance.

Besides the 1/2bit designs, the rest of the designs adopts 32-bit floating-point data or linear quantization with 8 or more bits. According to the results in section 4.1, within 1% accuracy loss can be achieved. So we think the comparison between these designs is fair in accuracy. INT16/8 and INT16 are commonly adopted. But the difference between these designs is not obvious. This is because the underutilization of DSPs discussed in section 5.1.1. The advantage of INT16 over FP32 is obvious except for [64], where the hard-core floating-point DSPs are utilized. To a certain extent, this shows the importance of fully utilizing the DSPs on-chip.

Fast Convolution Algorithm. Among all the 16-bit designs, [28] achieves the best energy efficiency and the highest speed with the help of the 6×6 Winograd fast convolution, which is $1.7\times$ faster and $2.6\times$ energy efficient than the 16-bit design in [64]. The design in [62] achieves $2\times$ speedup and $3\times$ energy efficiency compared with [61] where both designs use 32-bit floating-point data. Overall, the improvement does not match the estimation but can still reach $2\text{-}3\times$.

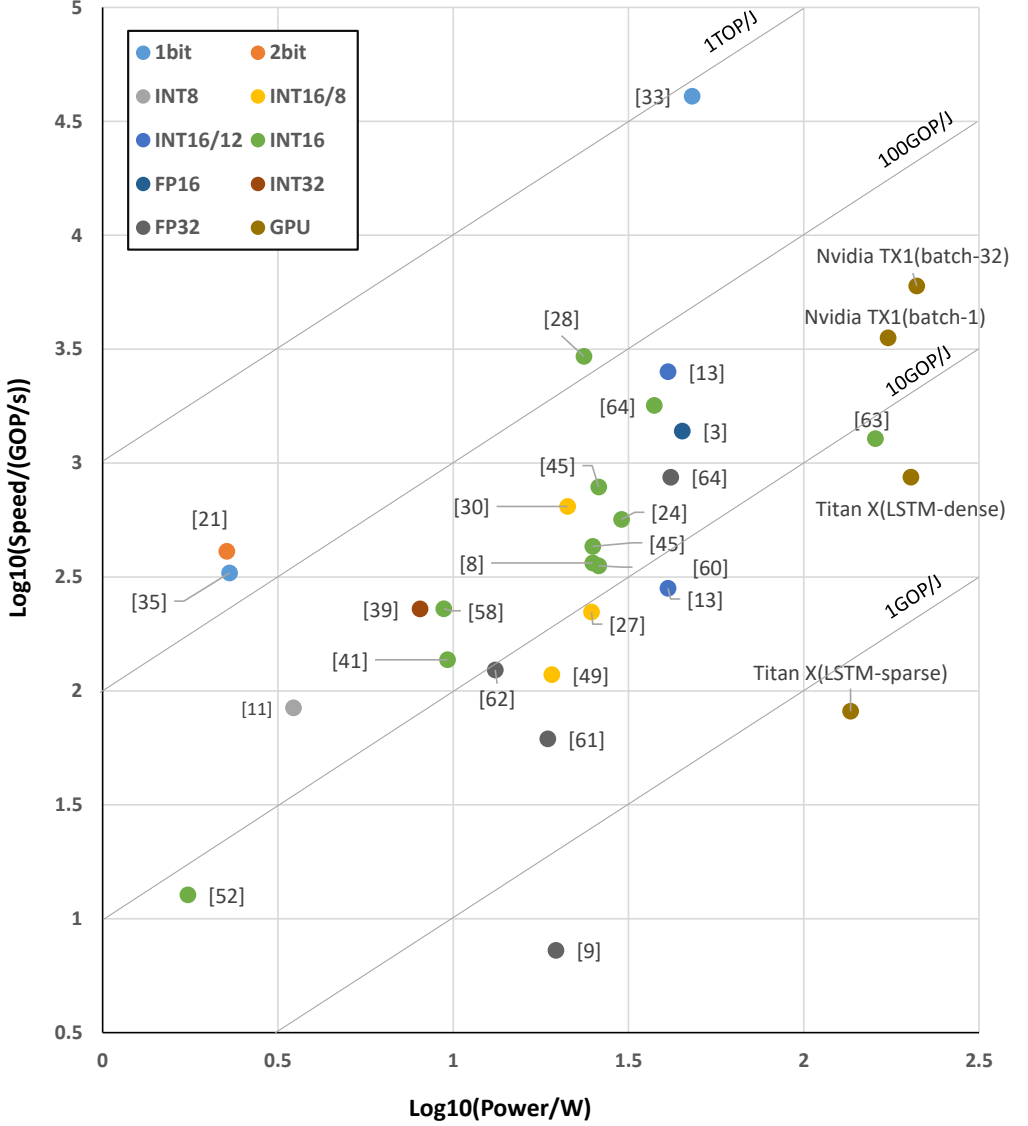


Fig. 6. A comparison between different designs on a logarithm coordinate of power and performance.

System Level Optimization. The overall system optimization is not well addressed in most of the work. As this is also related to the HDL design quality, we can roughly evaluate the effect. Here we compare three designs[24, 27, 60] on the same XC7VX690T platform and try to evaluate the effect. All the three designs implement 16-bit fixed-point data format except that [27] uses 8-bit for weights. No fast convolution or sparsity is utilized in any of the work. Even though, [24] achieves $2.5\times$ the energy efficiency of [27]. It shows that a system level optimization has a strong effect even comparable to the use of fast convolution algorithm.

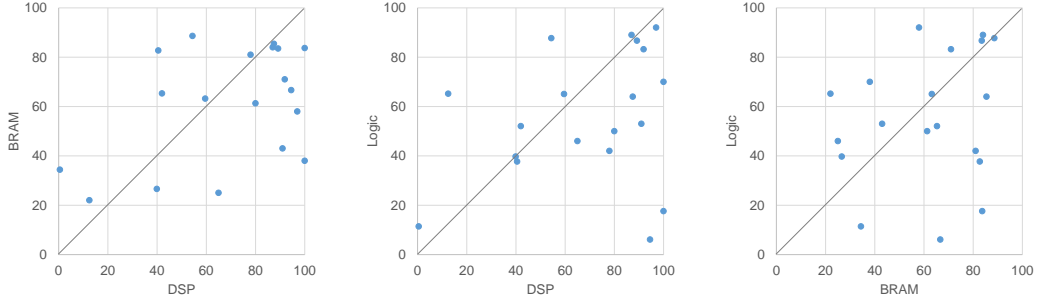


Fig. 7. Resource utilization ratio of different accelerator designs.

We also investigate the resource utilization of the designs in Table 2. Three kinds of resources (DSP, BRAM, and logic) are considered. We plot the designs in Figure 7 using two of the utilization ratio as x and y coordinate. We draw the diagonal line of each figure to show the designs' preference on hardware resource. The BRAM-DSP figure shows an obvious preference on DSP over BRAM. A similar preference appears on DSP over logic. This indicates that current FPGA designs are more likely computation bounded. FPGA manufacturers targeting neural network applications can adjust the resource allocation accordingly.

Comparison with GPU. In general, FPGA based designs have achieved comparable energy efficiency to GPU with 10-100GOP/J. But the speed of GPUs still surpasses FPGAs. Scaling up the FPGA based design is still a problem. Zhang et al. [63] propose the FPGA-cluster-based solution using 16-bit fixed-point computation. But the energy efficiency is worse than the other 16-bit fixed-point designs.

Here we estimate the achievable speed of an ideal design. We use the 16-bit fixed-point design in [28] as a baseline, which is the best 16-bit fixed-point design with both the highest speed and energy efficiency. 8-bit linear quantization can be adopted according to the analysis in section 4.1, which achieves another $2\times$ speedup and better energy efficiency by utilizing 1 DSP as 2 multipliers. The double frequency optimization further improves the system speed by $2\times$. Consider a sparse model which is similar to the one in [13] with 10% non-zero values. We can estimate a similar $6\times$ improvement as [13]. In general about $24\times$ speedup and $12\times$ better energy efficiency can be achieved, which means 72TOP/s speed with about 50W. This shows that it is possible to achieve over $10\times$ higher energy efficiency on FPGA over 32-bit floating-point process on GPU.

The left problem is: does all the techniques: double MAC, sparsification, quantization, fast convolution, and the double frequency design work well together? Pruning a single element in a 2D convolution kernel is of no use for fast convolution because the 2D kernel is always processed as a whole. Directly pruning 2D kernels as a whole may help. But the reported accuracy of this method is lower [31] than a fine-grained pruning. The irregular data access pattern for processing sparse network and the increase in parallelism also brings challenges to the design of memory system and scheduling strategy.

7 DESIGN AUTOMATION AND FLEXIBILITY

In some application scenarios, various NN models are to be supported with the FPGA accelerator. Whether the accelerator can respond to the change in network model promptly and keeps high performance becomes a key feature. To address this problem, various researches have been carried

out to map an NN model to an FPGA automatically. Mainly two kinds of methods are used: hardware design automation and software design automation. Hardware design automation generates different hardware designs according to different NN models. Software design automation keeps the same accelerator and generates different inputs to the accelerator.

7.1 Hardware Design Automation

Hardware design automation is widely adopted in FPGA based accelerators because of the reconfigurability of FPGAs [5, 29, 32, 44, 52–54]. These proposed techniques focus on automatically generate the HDL design based on the network parameter. Difference between these methods is the selection of an intermediate level description of the network to cover the gap between high-level network description and low-level hardware design.

A straightforward way is no intermediate description. The design flow in [29] searches the optimized parameter for a handcrafted Verilog template with the input network description and platform constraint. This method is similar to the optimization methods mentioned in section 5. DiCecco et al. [5] use a similar idea based on OpenCL model. This enables that the development tool be integrated with Caffe and one network can be executed on different platforms.

Venireis, et al. [53] describes the network model as a DFG in their design tool. Then the network computation process is translated to hardware design with DFG mapping method.

DnnWeaver [44] use a virtual instruction set to describe a network. The network model is first translated into an instruction sequence. Then the sequence is mapped as hardware FSM states but not executed like traditional CPU instructions.

Hardware design automation directly modifies the hardware design to support different networks. This means the hardware can always achieve the best performance on the target platform. This is suitable for FPGA because of its reconfigurability. It works in situations where network switching is not frequent and the reconfiguration overhead does not care. For example, for a large scale-cloud service, the change in network models can be covered by switching between different FPGA chips. So the FPGAs do not need to be reconfigured frequently.

7.2 Software Design Automation

Software design automation tries to run different networks on the same hardware accelerator by simply changing the input, in most cases, an instruction sequence. The difference between these work is the granularity of instruction. At a lower level, Guo, et al. [11] propose the instruction set with only three kinds of instructions: LOAD, CALC, and SAVE. The granularity of the LOAD and SAVE instructions are the same as the data tiling size. Each CONV executes a set of 2-D convolutions given the feature map size encoded in the instruction. The channel number is fixed as the hardware unrolling parameter. At this level, the software compiler can carry out static scheduling and dynamic data reuse strategy accordingly for each layer.

Zhang et al. [60] use a layer level instruction set. The control of a CNN layer is designed as a configurable hardware FSM. Compared with [11], this reduces the memory access for instruction while increasing the hardware cost on the configurable FSM.

Instruction based methods do not modify hardware and thus enables that the accelerator can switch between networks at run-time. An example of the application scenario is the real-time video processing system on a mobile platform. The process of a single frame can involve different networks if the task is complex enough. Reconfigure the hardware causes unacceptable overhead while instruction based methods can solve the problem if all the instructions of all the networks are prepared in memory.

7.3 Mixed Method

Wang, et al. [54] propose a design automation framework mixing the above two by both optimizing hardware design and compile software instructions. The hardware is first assembled with pre-defined HDL templates using the optimized hardware parameter. The data control flow of the computation process is controlled by software binaries, which is compiled according to the network description. It is possible that the hardware can be used for a new network by simply changing the software binaries.

8 CONCLUSION

In this paper, we review state-of-the-art neural network accelerator designs and summarize the techniques used. According to the evaluation result in section 6, with software hardware co-design, FPGA can achieve $13\times$ better energy efficiency than state-of-the-art GPU while using 30% power with conservative estimation. This shows that FPGA is a promising candidate for neural network acceleration. We also review the methods used for flexible accelerator design, which shows that current development flow can achieve both high performance and run-time network switch.

But there is still a gap between current designs and the estimation. Combining all the techniques requires software-hardware co-design. Using quantization and weight reduction together while maintaining the performance is challenging. Scaling up the design is another problem. Future work should focus on solving these challenges. There is still $10\times$ potential performance gain from using $1/2$ bits for neuron and weight representation. Future work should try to improve the model accuracy in these cases.

ACKNOWLEDGEMENT

This work was supported by National Natural Science Foundation of China (No. 61373026, 61622403, 61621091), National Key R&D Program of China 2017YFA0207600, and Joint fund of Equipment pre-Research and Ministry of Education (No. 6141A02022608).

REFERENCES

- [1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, et al. 2016. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467* (2016).
- [2] Manoj Alwani, Han Chen, Michael Ferdman, and Peter Milder. 2016. Fused-layer CNN accelerators. In *Microarchitecture (MICRO), 2016 49th Annual IEEE/ACM International Symposium on*. IEEE, 1–12.
- [3] Utku Aydonat, Shane O’Connell, Davor Capalija, Andrew C Ling, and Gordon R Chiu. 2017. An OpenCL (TM) Deep Learning Accelerator on Arria 10. *arXiv preprint arXiv:1701.03534* (2017).
- [4] Wenlin Chen, James Wilson, Stephen Tyree, Kilian Weinberger, and Yixin Chen. 2015. Compressing neural networks with the hashing trick. In *International Conference on Machine Learning*. 2285–2294.
- [5] Roberto DiCecco, Griffin Lacey, Jasmina Vasiljevic, Paul Chow, Graham Taylor, and Shawki Areibi. 2016. Caffeinated FPGAs: FPGA Framework For Convolutional Neural Networks. In *Field-Programmable Technology (FPT), 2016 International Conference on*. IEEE, 265–268.
- [6] Caiwen Ding, Siyu Liao, Yanzhi Wang, Zhe Li, Ning Liu, Youwei Zhuo, Chao Wang, Xuehai Qian, Yu Bai, Geng Yuan, et al. 2017. CirCNN: accelerating and compressing deep neural networks using block-circulant weight matrices. In *Proceedings of the 50th Annual IEEE/ACM International Symposium on Microarchitecture*. ACM, 395–408.
- [7] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. 2014. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 580–587.
- [8] Yijin Guan, Hao Liang, Ningyi Xu, Wenqiang Wang, Shaoshuai Shi, Xi Chen, Guangyu Sun, Wei Zhang, and Jason Cong. 2017. FP-DNN: An Automated Framework for Mapping Deep Neural Networks onto FPGAs with RTL-HLS Hybrid Templates. (2017), 152–159.
- [9] Yijin Guan, Zhihang Yuan, Guangyu Sun, and Jason Cong. 2017. FPGA-based accelerator for long short-term memory recurrent neural networks. In *Design Automation Conference (ASP-DAC), 2017 22nd Asia and South Pacific*. IEEE,

- 629–634.
- [10] Jianxin Guo, Shouyi Yin, Peng Ouyang, Leibo Liu, and Shaojun Wei. 2017. Bit-Width Based Resource Partitioning for CNN Acceleration on FPGA. In *Field-Programmable Custom Computing Machines (FCCM), 2017 IEEE 25th Annual International Symposium on*. IEEE, 31–31.
 - [11] Kaiyuan Guo, Lingzhi Sui, Jiantao Qiu, Jincheng Yu, Junbin Wang, Song Yao, Song Han, Yu Wang, and Huazhong Yang. 2017. Angel-Eye: A Complete Design Flow for Mapping CNN onto Embedded FPGA. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* (2017).
 - [12] PK Gupta. 2016. Accelerating datacenter workloads. In *26th International Conference on Field Programmable Logic and Applications (FPL)*.
 - [13] Song Han, Junlong Kang, Huizi Mao, Yiming Hu, Xin Li, Yubin Li, Dongliang Xie, Hong Luo, Song Yao, Yu Wang, et al. 2017. ESE: Efficient Speech Recognition Engine with Sparse LSTM on FPGA.. In *FPGA*. 75–84.
 - [14] Song Han, Huizi Mao, and William J Dally. 2015. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. *arXiv preprint arXiv:1510.00149* (2015).
 - [15] Awni Hannun, Carl Case, Jared Casper, Bryan Catanzaro, Greg Diamos, Erich Elsen, Ryan Prenger, Sanjeev Satheesh, Shubho Sengupta, Adam Coates, et al. 2014. Deep speech: Scaling up end-to-end speech recognition. *arXiv preprint arXiv:1412.5567* (2014).
 - [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
 - [17] M. Horowitz. [n. d.]. Energy table for 45nm process, Stanford VLSI wiki.[Online]. <https://sites.google.com/site/seecproject/>. ([n. d.]).
 - [18] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. 2017. MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications. (2017).
 - [19] Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer. 2016. SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and< 0.5 MB model size. *arXiv preprint arXiv:1602.07360* (2016).
 - [20] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. 2014. Caffe: Convolutional Architecture for Fast Feature Embedding. *arXiv preprint arXiv:1408.5093* (2014).
 - [21] Li Jiao, Cheng Luo, Wei Cao, Xuegong Zhou, and Lingli Wang. 2017. Accelerating low bit-width convolutional neural networks with embedded FPGA. In *Field Programmable Logic and Applications (FPL), 2017 27th International Conference on*. IEEE, 1–4.
 - [22] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*. 1097–1105.
 - [23] Fengfu Li, Bo Zhang, and Bin Liu. 2016. Ternary weight networks. *arXiv preprint arXiv:1605.04711* (2016).
 - [24] Huimin Li, Xitian Fan, Li Jiao, Wei Cao, Xuegong Zhou, and Lingli Wang. 2016. A high performance FPGA-based accelerator for large-scale convolutional neural networks. In *Field Programmable Logic and Applications (FPL), 2016 26th International Conference on*. IEEE, 1–9.
 - [25] Yixing Li, Zichuan Liu, Kai Xu, Hao Yu, and Fengbo Ren. 2017. A 7.663-TOPS 8.2-W Energy-efficient FPGA Accelerator for Binary Convolutional Neural Networks. *arXiv preprint arXiv:1702.06392* (2017).
 - [26] Baoyuan Liu, Min Wang, Hassan Foroosh, Marshall Tappen, and Marianna Pensky. 2015. Sparse convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 806–814.
 - [27] Zhiqiang Liu, Yong Dou, Jingfei Jiang, and Jinwei Xu. 2016. Automatic code generation of convolutional neural networks in FPGA implementation. In *Field-Programmable Technology (FPT), 2016 International Conference on*. IEEE, 61–68.
 - [28] Liqiang Lu, Yun Liang, Qingcheng Xiao, and Shengen Yan. 2017. Evaluating fast algorithms for convolutional neural networks on fpgas. In *Field-Programmable Custom Computing Machines (FCCM), 2017 IEEE 25th Annual International Symposium on*. IEEE, 101–108.
 - [29] Yufei Ma, Yu Cao, Sarma Vrudhula, and Jae-sun Seo. 2017. An automatic RTL compiler for high-throughput FPGA implementation of diverse deep convolutional neural networks. In *Field Programmable Logic and Applications (FPL), 2017 27th International Conference on*. IEEE, 1–8.
 - [30] Yufei Ma, Yu Cao, Sarma Vrudhula, and Jae-sun Seo. 2017. Optimizing Loop Operation and Dataflow in FPGA Acceleration of Deep Convolutional Neural Networks. In *Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*. ACM, 45–54.
 - [31] Huizi Mao, Song Han, Jeff Pool, Wenshuo Li, Xingyu Liu, Yu Wang, and William J. Dally. 2017. Exploring the Granularity of Sparsity in Convolutional Neural Networks. In *Computer Vision and Pattern Recognition Workshops*. 1927–1934.

- [32] Raghdh Marcel, Haitham Akkary, Hazem Hajj, Mazen Saghir, Anil Keshavamurthy, Rahul Khanna, and Hassan Artail. 2017. Minimalist Design for Accelerating Convolutional Neural Networks for Low-End FPGA Platforms. In *Field-Programmable Custom Computing Machines (FCCM), 2017 IEEE 25th Annual International Symposium on*. IEEE, 196–196.
- [33] Duncan JM Moss, Eriko Nurvitadhi, Jaewoong Sim, Asit Mishra, Debbie Marr, Suchit Subhaschandra, and Philip HW Leong. 2017. High performance binary neural networks on the Xeon+ FPGADc platform. In *Field Programmable Logic and Applications (FPL), 2017 27th International Conference on*. IEEE, 1–4.
- [34] Mohammad Motamedi, Philipp Gysel, Venkatesh Akella, and Soheil Ghiasi. 2016. Design space exploration of fpga-based deep convolutional neural networks. In *Design Automation Conference (ASP-DAC), 2016 21st Asia and South Pacific*. IEEE, 575–580.
- [35] Hiroki Nakahara, Tomoya Fujii, and Shimpei Sato. 2017. A fully connected layer elimination for a binarized convolutional neural network on an FPGA. In *Field Programmable Logic and Applications (FPL), 2017 27th International Conference on*. IEEE, 1–4.
- [36] Hiroki Nakahara, Haruyoshi Yonekawa, Hisashi Iwamoto, and Masato Motomura. 2017. A Batch Normalization Free Binarized Convolutional Deep Neural Network on an FPGA. In *Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*. ACM, 290–290.
- [37] Dong Nguyen, Daewoo Kim, and Jongeun Lee. 2017. Double MAC: Doubling the performance of convolutional neural networks on modern FPGAs. In *2017 Design, Automation & Test in Europe Conference & Exhibition (DATE)*. IEEE, 890–893.
- [38] Eriko Nurvitadhi, David Sheffield, Jaewoong Sim, Asit Mishra, Ganesh Venkatesh, and Debbie Marr. 2016. Accelerating Binarized Neural Networks: Comparison of FPGA, CPU, GPU, and ASIC. In *Field-Programmable Technology (FPT), 2016 International Conference on*. IEEE, 77–84.
- [39] Abhinav Podili, Chi Zhang, and Viktor Prasanna. 2017. Fast and efficient implementation of Convolutional Neural Networks on FPGA. In *Application-specific Systems, Architectures and Processors (ASAP), 2017 IEEE 28th International Conference on*. IEEE, 11–18.
- [40] Adrien Prost-Boucle, Alban Bourge, Frédéric Pétrot, Hande Alemdar, Nicholas Caldwell, and Vincent Leroy. 2017. Scalable high-performance architecture for convolutional ternary neural networks on FPGA. In *Field Programmable Logic and Applications (FPL), 2017 27th International Conference on*. IEEE, 1–7.
- [41] Jiantao Qiu, Jie Wang, Song Yao, Kaiyuan Guo, Boxun Li, Erjin Zhou, Jincheng Yu, Tianqi Tang, Ningyi Xu, Sen Song, et al. 2016. Going deeper with embedded fpga platform for convolutional neural network. In *Proceedings of the 2016 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*. ACM, 26–35.
- [42] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. 2015. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)* 115, 3 (2015), 211–252. <https://doi.org/10.1007/s11263-015-0816-y>
- [43] Mohammad Samragh, Mohammad Ghasemzadeh, and Farinaz Koushanfar. 2017. Customizing neural networks for efficient fpga implementation. In *Field-Programmable Custom Computing Machines (FCCM), 2017 IEEE 25th Annual International Symposium on*. IEEE, 85–92.
- [44] Hardik Sharma, Jongse Park, Divya Mahajan, Emmanuel Amaro, Joon Kyung Kim, Chenkai Shao, Asit Mishra, and Hadi Esmaeilzadeh. 2016. From high-level deep neural models to FPGAs. In *Microarchitecture (MICRO), 2016 49th Annual IEEE/ACM International Symposium on*. IEEE, 1–12.
- [45] Junzhong Shen, You Huang, Zelong Wang, Yuran Qiao, Mei Wen, and Chunyuan Zhang. 2018. Towards a Uniform Template-based Architecture for Accelerating 2D and 3D CNNs on FPGA. In *Acm/sigda International Symposium*. 97–106.
- [46] Yongming Shen, Michael Ferdman, and Peter Milder. 2016. Overcoming resource underutilization in spatial CNN accelerators. In *Field Programmable Logic and Applications (FPL), 2016 26th International Conference on*. IEEE, 1–4.
- [47] Yongming Shen, Michael Ferdman, and Peter Milder. 2017. Escher: A CNN Accelerator with Flexible Buffering to Minimize Off-Chip Transfer. In *Proceedings of the 25th IEEE International Symposium on Field-Programmable Custom Computing Machines (FCCM2017)*. IEEE Computer Society, Los Alamitos, CA, USA.
- [48] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).
- [49] Naveen Suda, Vikas Chandra, Ganesh Dasika, Abinash Mohanty, Yufei Ma, Sarma Vrudhula, Jae-sun Seo, and Yu Cao. 2016. Throughput-optimized OpenCL-based FPGA accelerator for large-scale convolutional neural networks. In *Proceedings of the 2016 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*. ACM, 16–25.
- [50] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich, et al. 2015. Going deeper with convolutions. *Cvpr*.
- [51] Yaman Umuroglu, Nicholas J Fraser, Giulio Gambardella, Michaela Blott, Philip Leong, Magnus Jahre, and Kees Visser. 2017. Finn: A framework for fast, scalable binarized neural network inference. In *Proceedings of the 2017 ACM/SIGDA*

- International Symposium on Field-Programmable Gate Arrays*. ACM, 65–74.
- [52] Stylianos I Venieris and Christos-Savvas Bouganis. 2017. fpgaConvNet: Automated Mapping of Convolutional Neural Networks on FPGAs. In *Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*. ACM, 291–292.
 - [53] Stylianos I Venieris and Christos-Savvas Bouganis. 2017. Latency-driven design for FPGA-based convolutional neural networks. In *Field Programmable Logic and Applications (FPL), 2017 27th International Conference on*. IEEE, 1–8.
 - [54] Ying Wang, Jie Xu, Yinhe Han, Huawei Li, and Xiaowei Li. 2016. DeepBurning: automatic generation of FPGA-based learning accelerators for the neural network family. In *Design Automation Conference (DAC), 2016 53rd ACM/EDAC/IEEE*. IEEE, 1–6.
 - [55] Xuechao Wei, Cody Hao Yu, Peng Zhang, Youxiang Chen, Yuxin Wang, Han Hu, Yun Liang, and Jason Cong. 2017. Automated Systolic Array Architecture Synthesis for High Throughput CNN Inference on FPGAs. In *Proceedings of the 54th Annual Design Automation Conference 2017*. ACM, 29.
 - [56] Shmuel Winograd. 1980. *Arithmetic complexity of computations*. Vol. 33. Siam.
 - [57] Ephrem Wu, Xiaoqian Zhang, David Berman, and Inkeun Cho. 2017. A high-throughput reconfigurable processing array for neural networks. In *Field Programmable Logic and Applications (FPL), 2017 27th International Conference on*. IEEE, 1–4.
 - [58] Qingcheng Xiao, Yun Liang, Liqiang Lu, Shengen Yan, and Yu-Wing Tai. 2017. Exploring Heterogeneous Algorithms for Accelerating Deep Convolutional Neural Networks on FPGAs. In *Proceedings of the 54th Annual Design Automation Conference 2017*. ACM, 62.
 - [59] Jincheng Yu, Yiming Hu, Xuefei Ning, Jiantao Qiu, Kaiyuan Guo, Yu Wang, and Huazhong Yang. 2017. Instruction driven cross-layer CNN accelerator with winograd transformation on FPGA. In *International Conference on Field Programmable Technology*. 227–230.
 - [60] Chen Zhang, Zhenman Fang, Peipei Zhou, Peichen Pan, and Jason Cong. 2016. Caffeine: Towards uniformed representation and acceleration for deep convolutional neural networks. In *Computer-Aided Design (ICCAD), 2016 IEEE/ACM International Conference on*. IEEE, 1–8.
 - [61] Chen Zhang, Peng Li, Guangyu Sun, Yijin Guan, Bingjun Xiao, and Jason Cong. 2015. Optimizing fpga-based accelerator design for deep convolutional neural networks. In *Proceedings of the 2015 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*. ACM, 161–170.
 - [62] Chi Zhang and Viktor Prasanna. 2017. Frequency domain acceleration of convolutional neural networks on CPU-FPGA shared memory system. In *Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*. ACM, 35–44.
 - [63] Chen Zhang, Di Wu, Jiayu Sun, Guangyu Sun, Guojie Luo, and Jason Cong. 2016. Energy-Efficient CNN Implementation on a Deeply Pipelined FPGA Cluster. In *Proceedings of the 2016 International Symposium on Low Power Electronics and Design*. ACM, 326–331.
 - [64] Jialiang Zhang and Jing Li. 2017. Improving the Performance of OpenCL-based FPGA Accelerator for Convolutional Neural Network.. In *FPGA*. 25–34.
 - [65] Xiangyu Zhang, Jianhua Zou, Xiang Ming, Kaiming He, and Jian Sun. 2015. Efficient and accurate approximations of nonlinear convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1984–1992.
 - [66] Ritchie Zhao, Weinan Song, Wentao Zhang, Tianwei Xing, Jeng-Hau Lin, Mani B Srivastava, Rajesh Gupta, and Zhiru Zhang. 2017. Accelerating Binarized Convolutional Neural Networks with Software-Programmable FPGAs.. In *FPGA*. 15–24.
 - [67] Shuchang Zhou, Yuxin Wu, Zekun Ni, Xinyu Zhou, He Wen, and Yuheng Zou. 2016. DoReFa-Net: Training low bitwidth convolutional neural networks with low bitwidth gradients. *arXiv preprint arXiv:1606.06160* (2016).
 - [68] Chenzhuo Zhu, Song Han, Huizi Mao, and William J Dally. 2016. Trained ternary quantization. *arXiv preprint arXiv:1612.01064* (2016).