

Learning on Hardware: A Tutorial on Neural Network Accelerators and Co-Processors

LUKAS BAISCHER, MATTHIAS WESS, and NIMA TAHERINEJAD, TU Wien, Institute of Computer Technology, Austria

Deep neural networks (DNNs) have the advantage that they can take into account a large number of parameters, which enables them to solve complex tasks. In computer vision and speech recognition, they have a better accuracy than common algorithms, and in some tasks, they boast an even higher accuracy than human experts. With the progress of DNNs in recent years, many other fields of application such as diagnosis of diseases and autonomous driving are taking advantage of them. The trend at DNNs is clear: The network size is growing exponentially, which leads to an exponential increase in computational effort and required memory size. For this reason, optimized hardware accelerators are used to increase the performance of the inference of neuronal networks. However, there are various neural network hardware accelerator platforms, such as graphics processing units (GPUs), application specific integrated circuits (ASICs) and field programmable gate arrays (FPGAs). Each of these platforms offer certain advantages and disadvantages. Also, there are various methods for reducing the computational effort of DNNs, which are differently suitable for each hardware accelerator. In this article an overview of existing neural network hardware accelerators and acceleration methods is given. Their strengths and weaknesses are shown and a recommendation of suitable applications is given. In particular, we focus on acceleration of the inference of convolutional neural networks (CNNs) used for image recognition tasks. Given that there exist many different hardware architectures. FPGA-based implementations are well-suited to show the effect of DNN optimization methods on accuracy and throughput. For this reason, the focus of this work is more on FPGA-based implementations.

Additional Key Words and Phrases: neural network, hardware accelerator, deep learning, CNN, FPGA, ASIC, GPU, dataflow processing, energy efficient accelerators, performance gap

ACM Reference Format:

Lukas Baischer, Matthias Wess, and Nima TaheriNejad. 2021. Learning on Hardware: A Tutorial on Neural Network Accelerators and Co-Processors. *arXiv* 0, 0, Article 0 (April 2021), 29 pages. <https://doi.org/x>

1 INTRODUCTION

Deep learning is currently one of the most prominent machine learning approaches for solving complex tasks that could only be solved by human beforehand [60]. In applications such as computer vision or speech recognition, DNNs achieve higher accuracy compared to non-learning algorithms and in some cases even higher than human experts [45, 57, 60]. The higher accuracy of DNNs compared to non-learning algorithms comes from the ability to extract high-level features from the input data after using statistical learning over a high number of training data. Statistical learning leads to an efficient representation of the input space and a good generalization [45].

However, this capability requires high computational effort [45]. It has shown that, by increasing the number of parameters, the accuracy of a network can be increased [57]. Consequently, the trend

Authors' address: Lukas Baischer, lukas.baischer@student.tuwien.ac.at; Matthias Wess, matthias.wess@student.tuwien.ac.at; Nima TaheriNejad, nima.taherinejad@tuwien.ac.at, TU Wien, Institute of Computer Technology, Gusshausstraße 27-29/384, Vienna, Vienna, Austria, 1040.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

© 2021 Copyright held by the owner/author(s).

XXXX-XXXX/2021/4-ART0

<https://doi.org/x>

in DNNs is clearly that the network size is growing exponentially. Which leads to an exponentially increasing computational effort and required memory size [57]. Therefore, central processing units (CPUs) can hardly handle the necessary computations. Hence, structurally optimized hardware accelerators are used to increase the inference performance of neuronal networks. For inference of a neural network running on edge devices, energy efficiency is an important factor that has to be considered, in addition to throughput [57].

In this paper, we focus on optimization methods and hardware accelerators used for inference of CNNs. FPGA-based implementations are discussed in more detail, since they are well suited to show different architecture approaches and how optimization methods affect throughput and accuracy. However, in addition to hardware accelerators mentioned in this paper, other promising approaches such as spiking neural network (SNN) or in memory computation (IMC) exist. Nevertheless, these approaches require different optimization methods or different math, which is beyond the scope of this paper. Below is an outline of the paper.

Section 2 gives a general introduction to the field of neural networks. The intention is to explain the necessary background to understand how neural networks work and to show the basic mathematical operations. Section 3 highlights the need for neural network hardware accelerators and shows why simply using general-purpose GPUs cannot fulfill future requirements in terms of performance density and energy efficiency. Based on the knowledge that more efficient neural network structures are required to meet future requirements, Section 4 shows methods for decreasing the computation effort, required for deep neural networks. In Section 5 algorithmic optimization methods are explained, which increase the algorithmic efficiency and therefore, increase the throughput of hardware accelerators. Section 6 explains the parallelization strategies used by neural network hardware accelerators. Section 7, Section 8 and Section 9 present the most commonly used neural network hardware accelerator platforms, namely GPUs, ASICs and FPGAs. Additionally, an overview of existing implementations is introduced. Section 10 compares the platforms in terms of speed, accuracy, power, and usability. Finally, we draw our conclusions in section 11.

2 NEURAL NETWORKS

Neural networks are a part of the broad field of artificial intelligence (AI). AI generally deals with building intelligent machines, which can solve tasks similar to human [45]. Neural networks are inspired by the human brain. Figure 1 visualizes the principle of a single neuron in a neural network and the similarities to the human brain. Additionally, the fundamental Equation (1) for computing the output activation of a single neuron is depicted. x_i , w_i , $f(\cdot)$, b and y_j are the input activations, weights, non-linear activation function, bias and output activation [45].

$$y_j = f\left(\sum_i w_i x_i + b\right) \quad (1)$$

In principle, any non-linear function can be allowed as an activation function for a neural network. Historically the Sigmoid function was the most popular activation function. However in modern DNNs mostly rectified linear unit (ReLU) is used [45]. It has shown that ReLU achieves good results and provides the benefit of less computational effort since only negative values are set to zero. All others stay untouched [26]. However, cutting off all negative neuron outputs leads to a loss of information. For this reason leaky ReLU is used in order to additionally consider negative neuron outputs [4].

In a neuronal network, single neurons are arranged in layers. How neurons are arranged in

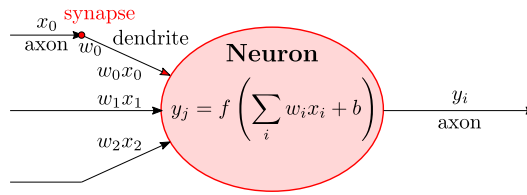


Fig. 1. A single neuron and the connection to human brain. x_i , w_i , $f(\cdot)$, b and y_j are the input activations, weights, non-linear activation function, bias and output activation. Adapted from [45, Fig. 2].

a layer and how they are connected to the successor layer indicates the type of layer or network. Three important neural network types are:

- Fully connected (FC) or feed-forward neural networks
- Recurrent neural networks (RNNs)
- Convolutional neural networks (CNNs)

Fully connected neural networks are historically the first form. In modern networks, they are often used as part of CNNs or RNNs. In some literature FC-layers are called a dense layer, which comes from the high parameter density of fully connected neural networks. CNNs are mainly used in the field of computer vision since CNNs consider geometric information. Additionally, CNNs require less memory space, since weights in a filter stay constant for a single layer. RNNs are mainly used in tasks such as speech recognition or predictions in the financial world because they provide the possibility to consider the information of former events. This is done by using feedback from a layer to a previous one. DNN, on the other hand, do not necessarily designate the structure of a network. It is related to the keyword deep learning, which should indicate that particularly deep networks, i.e. networks with many layers, are used [45].

2.1 Fully connected network

Figure 2 visualizes the principle and terminology of a fully connected network. The first layer is called the input layer, which receives the data that contains the information to be analyzed. That can be, for example, an image or a sampled audio sequence. The layers which are processing the data are called hidden layers. The output layer is the last layer of a neural network and provides the results of a neural network. Usually, the output of a neural network is represented in the form of a probability of specific possible results. Each layer of a fully connected network is described by Equation (1), whereas the weight matrix w_i contains an entry for each neuron connection between two layers.

2.2 CNN

Convolutional neural networks (CNNs) are state of the art for images classification and detection tasks. Compared to other neural network structures, CNNs can consider geometric information included in each pixel, which enables them to surpass the accuracy of fully connected networks using fewer parameters [26]. The 2-D convolution operation is fundamental for CNNs. A filter with size $k \times k$ moves along an image or feature map producing a single output image at each step. This is visualized by Figure 3a. It shows that at each step an element-wise multiplication of the filter with the overlapping image tile is performed. The results of the element-wise multiplication are accumulated. Equation (2) shows how to compute a single output pixel using a 2-D convolution [1, 43, 45].

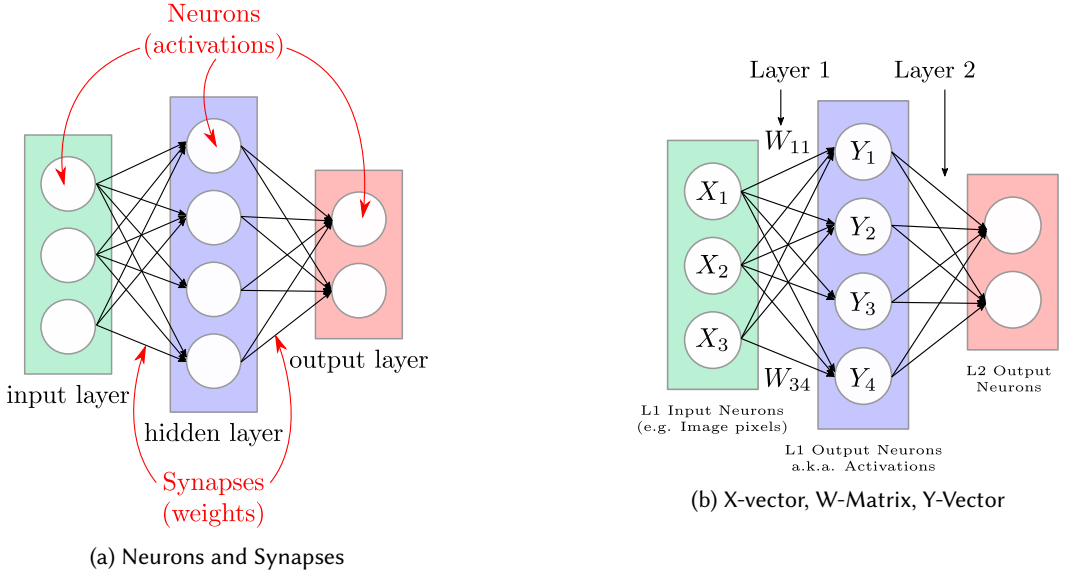


Fig. 2. Simple neural network example and terminology. (a) shows a simple example of a fully connected neural network. (b) gives an overview how the X-Vector, W-matrix, Y-vector are interpreted. Adapted from [45, Fig. 3].

$$\text{conv2D}(X^{\text{conv}}[c], \Theta[n, c]) = \sum X^{\text{conv}}[c] \odot \Theta[n, c] \quad (2)$$

Multiple filters are applied in parallel to create complex high-level feature maps for object classification. Figure 3b visualizes that C_o filter applied on a single input image create C_o output feature maps. Parallel feature maps in a single layer are called *channels*. For computing a single pixel of a feature map, the results of the 2-D convolution of all C_i input channels are accumulated. After the summation a non-linear function $f(\cdot)$ is applied. The weights of each filter are different for each channel. That leads to $C_i \times C_o$ different filter for a single layer. Equation (3) shows how to compute a single output pixel of a multidimensional convolutional layer [1].

$$Y^{\text{conv}}[n] = f \left(\sum_{c=0}^C \text{conv2D}(X^{\text{conv}}[c], \Theta[n, c]) \right) \quad (3)$$

It visualizes that for computing an output feature map of a single CNN layer, $C_o \cdot C_i \cdot W \cdot H \cdot K \cdot K$ multiply and accumulate (MAC) operations are required. The huge computational effort, which is necessary for CNNs, leads to the need of dedicated and structural optimized hardware acceleration platforms such as GPUs, FPGAs and ASICs [1].

2.3 RNN

In contrast to FC-neural networks and CNNs, recurrent neural networks (RNNs) use intralayer recurrent connection. These recurrent connections allow RNNs to consider past results or states in addition to the current input data [13, 45]. Due to this ability, RNNs are particularly suitable for tasks like speech recognition, translation, and financial predictions.

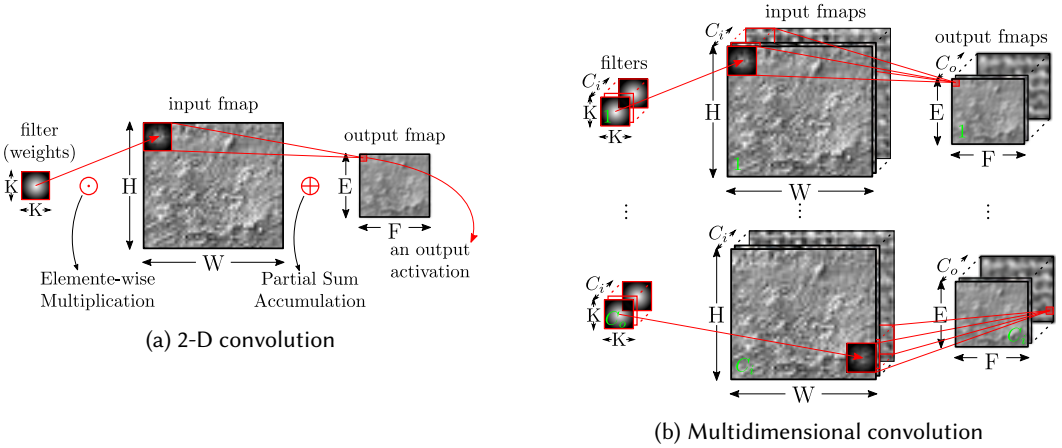


Fig. 3. Convolutional neural network principle. (a) shows the principle of a 2-D convolution used in image processing. (b) visualizes the principle of convolutions in high dimensional CNNs. Adapted from [45, Fig. 9].

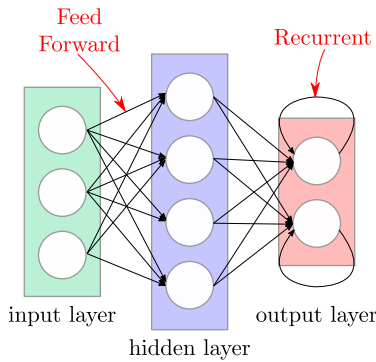
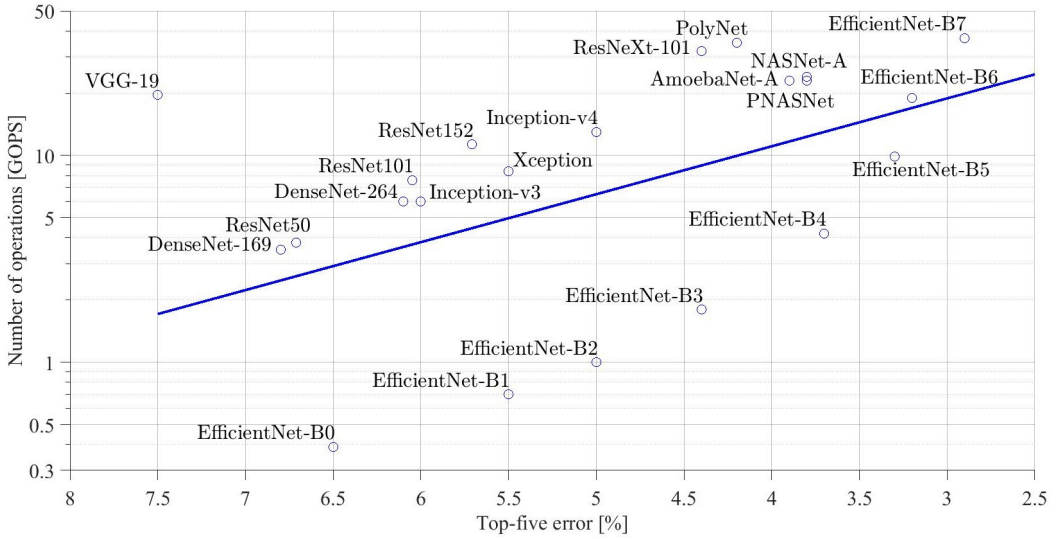


Fig. 4. Simplified structure of a recurrent neural network [45, Fig. 8].

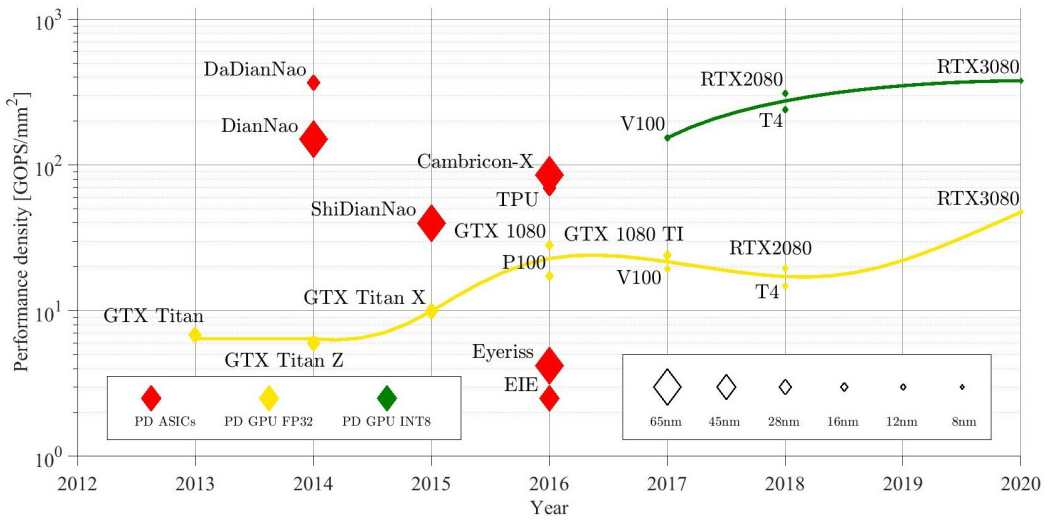
Figure 4 shows a simplified structure of a RNN. It visualizes that the output of the networks is traced back to the input of the output layer. That enables a RNN to additionally consider the former output of the network. Recurrent layer are mostly implemented using long short term memory (LSTM) or gated recurrent unit (GRU) layer [13].

3 GAPS BETWEEN ACCURACY REQUIREMENT AND HARDWARE PERFORMANCE

Figure 5a shows the required number of operations for structural optimized neural networks versus the achievable top-5 error using the ImageNet dataset. The blue line visualizes the for a linear increase of the accuracy leads to an exponential increase of the computational effort. Figure 5b visualizes the improvement of the performance density of neural network hardware accelerators. Comparing Figure 5a and Figure 5b concludes that the number of operation required to compute a neural network increases faster than the hardware performance can follow. That leads to a gap between the required hardware performance and the actual available hardware performance [57]. Figure 5b shows that the performance density of GPUs using 32bit floating-point operations depends



(a)



(b)

Fig. 5. Gap between required number of operations and performance density. (a) Number of operations versus top-five error rate for leading DNN designs from ImageNet classification competition. Data from: [21, 48]. (b) Performance density (PD) of leading GPU and ASIC platforms. To catch up with the required number of operations, simply increasing the chip area is not feasible. Data from [6–8, 18, 25, 37–39, 57, 61], Adopted from [57, Fig. 3]

mainly on the feature size¹ of the manufactured die. The feature size is visualized by the size of the

¹In this context, feature size refers to the technology used to manufacture the integrated circuit (IC) and not on the number of feature maps in a CNN

diamonds in Figure 5b. It is also shown that a simple way to increase the performance density of a hardware accelerator is to reduce the bit width of arithmetic operations. That is the main reason why ASICs and FPGAs can reach a higher performance density compared to general-purpose GPUs [57]. For this reason, modern GPUs implement specialized fixed-point arithmetic units, which increases the performance density and power efficiency of GPUs [38, 39]. However, GPUs still require additional 32-bit floating-point units for use as graphics cards. Nevertheless, GPUs exhibit a higher throughput compared to FPGAs and most ASICs [45, 57].

The performance and energy efficiency of hardware accelerators is tied to the feature size of the respective manufacturing technology. However, the feature size of a die can only be decreased to a certain point due to Moore's law end. Therefore, increasing the performance density at the same speed as the network parameter number is not possible [57].

4 DNN COMPRESSION METHODS

To narrow the gap between accuracy requirement and hardware performance, substantially two approaches are followed. On the one hand, compact models and pruning create more efficient neural network structures. On the other hand, decreasing the bit width of weights and activations leads to more efficient hardware [45, 53, 57]. General-purpose GPUs are usually designed for computation of large floating-point vectors. Since sparsity reduces the vector size and quantization focus on fixed-point operations, general-purpose GPUs have little profit from these methods. Due to their greater freedom of design, ASIC and FPGA based hardware accelerators achieve higher efficiency for such structures [45, 53, 57].

SNNs use discreet spikes that propagate through the network instead of activations represented as a number. They are a more biologically realistic replica of the human brain since biological neurons use a similar technique. SNNs bring the advantage that they are more power-efficient and hardware friendly compared to common artificial neural networks (ANNs). In terms of accuracy SNNs are still behind common ANNs, but the gap is decreasing [50]. The methods used for SNNs are very different from the methods used in conventional DNNs and are therefore considered separately from other approaches. For this reason, SNNs will not be discussed further here.

4.1 Pruning

Pruning pursues the approach of finding those weights that have little or no influence on the result and can be eliminated, therefore. That reduces network complexity and over-fitting [19, 40]. Reducing the network complexity has two advantages: The disk storage of the network reduces and the computational effort reduces [9, 57]. A pruned network is also called a sparse network [19]. Whether pruning is permitted or not can be determined with different methods. The absolute value of the weights is found to be a proper measure [17, 20]. Retraining is required after pruning to avoid accuracy loss [17, 20]. An issue with retraining pruned networks can be that the gradient vanishes, which leads to slow convergence. The reason for this is simply that a large proportion of the connections is pruned. That effect can be reduced by pruning convolutional and fully connected layers separately. Additionally, slower pruning helps to increase the convergence speed during training [17, 20].

Pruning can be done in various granularity: fine-grained pruning, vector-level pruning, kernel-level pruning, group-level pruning, and filter-level pruning [9]

4.1.1 Fine-grained pruning. Fine-grained pruning has no constraints. Any weight which finds to be unimportant can be pruned[9]. Han et al. introduced a Fine-grained pruning method called Deep Compression. In the first step, the network is trained as usual. In the second step, all weights

below a certain threshold are pruned. Finally, the remaining network is retrained. That method reduces, the number of parameters required by AlexNet by $9\times$ and by VGG-16 by $13\times$ without any drop in accuracy [19]. Incorrect pruning leads to severe accuracy loss. Therefore, caution is advised using pruning, as connections cannot be restored [17]. For this reason, Guo et al. introduced splicing. Splicing reconnects pruned connections, which have shown to be important. Splicing is inspired by the human brain since pruning and splicing are analogical to the synthesis of excitatory and inhibitory neurotransmitter in the human nervous system [17]. Based on these methods the dynamic network surgery algorithm iteratively compresses the network. Using dynamic network surgery leads to a $17.7\times$ compression rate using AlexNet, with only a minor accuracy loss of 0.33%.

4.1.2 Vector-level and kernel-level pruning. In vector-level and kernel-level pruning, pruning is constrained to vectors or kernels of convolutional layer [9]. Compared to fine-grained pruning vector-level and kernel-level pruning is more hardware friendly, since hardware accelerators are optimized for vector operations and therefore only benefit slightly from fine-grained pruning. Additionally, vector-level and kernel-level pruning require fewer indices due to de vector structure, which reduces the required storage [9].

4.1.3 Group-level pruning. If a filter in a convolutional layer operates on the same input feature map and provide the same pattern, the filters can be merged. Group level pruning bases on this method [9]. Lebedev et al. introduce the group-wise brain damage algorithm for group-level pruning. It achieves a $3.2\times$ speed up for convolutional layers of AlexNet [27].

4.1.4 Filter-level pruning. Filter-level pruning focus on reducing the dimension of a layer. The feature map of the following layer indicates which channel can be pruned [9].

4.2 Tensor Decomposition

As mentioned before, tensor (including matrix) computations are the basic operations required by neural networks. For this reason, tensor decomposition methods are well suited to compress and accelerate neural network models [13]. Tensor decomposition has two advantages. First, it decomposes a large tensor into two smaller tensors that fit better into local memory. Second, if the original tensor doesn't have full rank, tensor decomposition reduces the computational effort, because it eliminates redundancies.

4.2.1 Tensor Decomposition for CNNs. In tensor decomposition, it is assumed that the output produced by $\mathbf{y} = W\mathbf{x}$ can be approximately replicated by a low-rank subspace [63]. This low-rank subspace is created using low-rank decomposition, which results in $\mathbf{y} = PW'\mathbf{x} + b$. That brings the benefit that the computational complexity can be reduced from $O(dk^2c)$ to $O(d'k^2c) + O(dd')$. With c is the input channel size, d is the original output channel size and d' is the intermediate channel size [63]. Figure 6 visualizes the principle of tensor decomposition. Zhang et al. showed that tensor decomposition can speed up the computation by a factor of 2.9 with only a minor accuracy loss of 0.3% [63].

4.3 Compact Models

In contrast to creating sparse models using pruning, compact models can be used for the targeted creation of sparse neural network models. Compact models aim to replace large kernels of convolutional layers, with smaller kernels. For example a 5×5 convolution can be replaced with two 3×3 convolutions requiring 28% fewer parameter [13]. However, 3×3 convolutions are still computationally expensive. For this reason, depthwise separable convolutions were introduced, which consists of 1×1 convolutions followed by depthwise separated 3×3 convolutions [10]. The 1×1 convolutions observe the correlation of the different channels, while the 3×3 convolutions

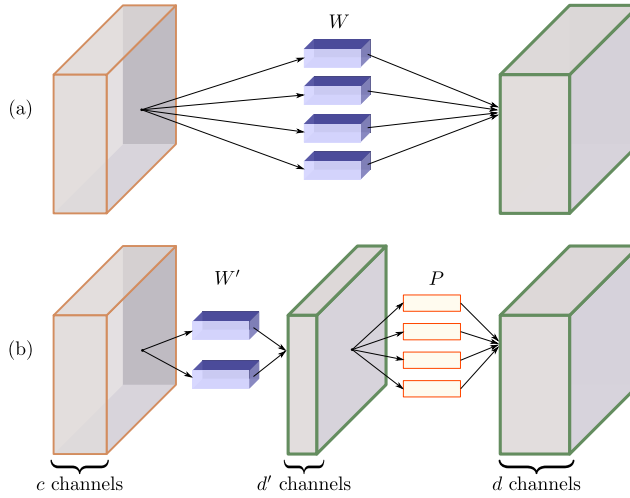


Fig. 6. Decomposition of weight tensors in CNNs.(a) original layer with complexity $O(dk^2c)$. (b) approximated layer with complexity reduced to $O(d'k^2c) + O(dd')$. Figure adopted from [63, Fig. 1].

act as a conventional filter in the x and y-direction. This approach is similar to that of tensor decomposition in CNNs (See Section 4.2), except that this structure is specified before the training and is not generated from the trained weights.

Another method introduced by ResNet is to reuse the information of features generated by a former layer. That is done by adding it to the current feature map. Such networks are called residual network [21]. Residual networks allow the creation of very deep networks, which yields an accuracy gain [13].

Figure 7 shows how these approaches are typically realized in compact models. Figure 7a shows a residual connection. Figure 7b shows a group convolution, which splits the channels into group-wise separated channels. Figure 7c shows a densely connected block. Figure 7d shows a bottleneck layer, which downscales the channel number by using a less computational intensive 1×1 convolution layer. After the downscale, the computational costly 3×3 convolution layer performs on fewer input and output channels. After that, the channel number expands again using a 1×1 convolution layer. The inverse bottleneck layer visualized in Figure 7e, does the opposite, but it uses depthwise separated 3×3 convolutions. Figure 7f shows the difference between depthwise separated convolutions and pointwise convolutions. It shows that depthwise convolutions only consider a single input channel for computing a new output channel. In most modern compact models these approaches are combined [24, 41, 46, 48].

In MobileNetV2 an inverted residual block is introduced, which is a combination of an inverse bottleneck layer and a residual connection [42]. Using these blocks Tan et al. show how to find the optimum network model for given hardware resources [48]. The goal of this approach is to achieve the highest possible accuracy with the lowest possible parameter number. However, the lower number of required operations of EfficientNet does not necessarily lead to a faster computation time of a single image. Bochkovski et al. show that the less complex network structure used by Darknet-53 leads to a lower computation time² compared to EfficientNet, even though nearly 5 times more computation operations are required [4].

²Using an Nvidia RTX 2070 GPU in single batch mode.

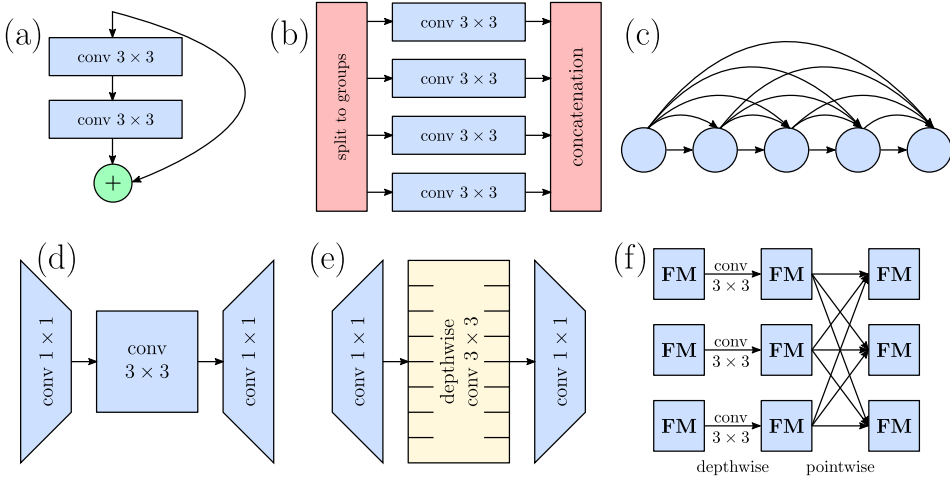


Fig. 7. Compact model design elements. (a) Residual connection. (b) Group convolution. (c) Densely connected block. (d) Bottleneck layer. (e) Inverse bottleneck layer. (f) Depthwise separable convolution. Adopted from [13, Fig. 6].

4.4 Knowledge Distillation

The idea of knowledge distillation bases on the assumption that averaging the predictions of different models gives higher accuracy than a single model can achieve. Since that would lead to an enormous additional computational effort, knowledge distillation combines the knowledge of different networks into a new more simple network. The original models are called the teacher. The new model is called the student. Based on the output of the teacher networks, the weights and biases of the student network, are adjusted using stochastic gradient descent. Using knowledge distillation, a simple student model can achieve an accuracy that would be unachievable using direct training [45].

4.5 Quantization

Quantize data is essential in digital computations since the bit width of arithmetic and logic units (ALUs) is limited in all digital devices. Basically, the higher the bit width and thus the higher the computing accuracy, the greater the hardware costs [45]. Floating-point quantization allows a higher precision compared to fixed-point quantization, due to the arbitrary exponent value. However, for floating-point more complex arithmetic operations are required, compared to fixed-point operations. This lead to higher power consumption, resource costs, or latency [53]. For this reason fixed-point quantization is heavily used in ASIC- and in FPGA-based hardware accelerators [6, 8, 15, 18, 23, 25, 33, 40, 53, 55, 59, 61].

The goal in quantizing neural networks is to minimize the error between the quantized and the original network while reducing hardware costs [45]. It is possible to use distinctive types of conditioning for different layers, channels, and filter [45]. Both, the weights and the activations can be quantized. A change in the bit width of the weights has a smaller effect on the accuracy than with the activations [45].

Data can be quantized linear and non-linear. Linear quantization brings the benefit that well-known calculation methods can be used, while specialized computation methods for non-linear quantized

data are necessary. The advantage of non-linear quantization is that numerical ranges that occur frequently have a higher resolution than those that occur less frequently. These numerical ranges can either be found by observing the distribution of weights and activations in a neural network or by learning them e.g. using k-means clustering [45]. Non-linear computations typically use lookup tables (LUTs), since they offer the possibility to output a specific bit-pattern depending on an input pattern.

Binarized Quantization. Binarized quantization is an extreme version of fixed-point quantization. The basic idea is to constrain both weights and activations to $\{1, -1\}$ [12]. That brings the benefit that hardware-friendly XNOR operations can replace costly multiplications [29]. However, binarized quantization of both weights and activations leads to a heavy loss in accuracy [45]. For this reason, retraining is necessary for binarized neural network (BNN) to regain some accuracy. Besides, the original is modified to increase accuracy. For example, the precision of the activations can be increased, or the first and/or the last layer are computed with full precision [33, 45].

Quantization of BNN can be done using stochastic or deterministic methods. However, typically the $Sign(x)$ function is used [12, 29]:

$$x_b = Sign(x) = \begin{cases} +1 & x \geq 0 \\ -1 & x < 0 \end{cases} \quad (4)$$

The problem using the $Sign(x)$ function for training, is that the derivative is a Dirac $\delta(x)$ function, which is not suitable for calculating the direction of the gradient in order to update the weight. For this reason $tanh(x)$ or $Htanh(x)$ functions are used instead of the $Sign(x)$ [29].

Ternary Quantization. Indifference to binarized quantization, ternary quantization adds a bit in order to be able to represent a weight of 0 [22]. That enables ternary neural networks (TNNs) to prune connections in the network. Which can increase the accuracy and significantly decrease the inference latency, since zero weights can be skipped [22]. Originally ternary quantization uses $\{-1, 0, +1\}$. However, some other constant values different from 1 are used in some related works to increase accuracy. For example, $\{-E, 0, +E\}$ can be used, with E is the mean absolute value of the learned weights. Another possibility is to use learned full-precision weights $\{-W_l^n, 0, +W_l^p\}$ which are constant for each layer [64].

Compared to full-precision models, ternary quantization reduces the model size by a factor of 16 [64]. Additionally if $\{-1, 0, +1\}$ quantization is used, no costly multiplication operations are required [22, 64]. In contrast $\{-E, 0, +E\}$ and $\{-W_l^n, 0, +W_l^p\}$ quantization still requires multiplications. Nevertheless, constant multipliers can be used. In both cases, a more efficient hardware accelerator implementation in terms of speed, power efficiency, and area consumption is possible, with only a moderate loss in accuracy [22, 64]. However, compared to BNNs, TNNs require double the model size, but provide a gain of accuracy. Therefore, TNNs represent a trade-off between model size and accuracy [64].

Nevertheless, retraining is required to achieve a competitive accuracy using ternary quantization, which can surpass the full-precision accuracy in some cases [22, 64].

Zhu et al. outperformed full the precision model using ResNet-32,44,56 on CIFAR-10 dataset by 0.04%, 0.16%, 0.36%, and AlexNet on ImageNet by 0.3% using ternary $\{-W_l^n, 0, +W_l^p\}$ quantization [64].

He et al. achieved an accuracy loss of $\sim 3.9\%$, 2.52% , 2.16% in comparison to full precision using a $\{-1, 0, +1\}$ ternarized ResNet-18/34/50 on ImageNet dataset.

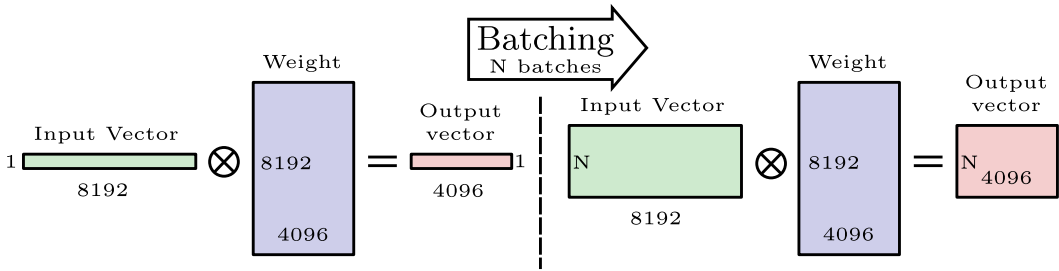


Fig. 8. Batching to convert vector matrix multiplication into matrix-matrix multiplication. It is demonstrated that the reuse of the weight matrix increases by a factor of N if batching is used. Figure adopted from [31, Fig. 7].

Log2 Quantization. Log2-quantization is a subclass of non-linear quantization, which allows transforming multiplications into shifts. The optimal quantization strategy for a single layer depends on the distribution of the weights. Typically, training uses an L_2 regularization, which forces gradients to favor weights close to 0. That causes the weights in a layer to be normal distributed with the mean value 0 [28]. For this reason, an alternative for arbitrary non-linear quantization is to use logarithmic quantization, since it provides a higher resolution for small weights. In a digital system, base 2 is optimal, since it allows to transform multiplications into shifts [28, 45]. Log2 quantization shows to achieve higher accuracy compared to linear quantization, for low-resolution weights of 4 bits and less. Additionally, area intensive multipliers change to shifts, which allows a higher performance density [28].

5 ALGORITHMIC OPTIMIZATION

Algorithmic optimization applies computational transformation or vectorization of data to reduce the number of performed arithmetic operations and memory accesses. Mainly CPUs and GPUs use these techniques. However, various FPGAs based neural network implementations are using algorithmic optimization [1, 31, 45].

5.1 General Matrix Multiplication

General matrix multiplication (GEMM) intendeds to minimize the necessary memory accesses for vector-matrix multiplications by combining multiplications with the same weight into a vector. The benefit of this method is that a single weight is fetched only once. For this reason, GEMM implementations are most efficient for large matrices. That is the case when an entire batch of input data is computed, which is called batching. Figure 8 visualizes that the reuse of a weight matrix increases by a factor of N for N batches. Especially for FC-layers, batching increase the utilization of single instruction multiple data (SIMD) processors [31]. However, if real-time applications use batching, it has to be ensured that the latency requirements are not violated, since batching may increase the latency [31].

GEMM is widely used by GPUs, since matrix multiplications are processed more efficiently by SIMD and single instruction multiple threads (SIMT) architectures, if GEMM is used [1, 34, 45]. Also for high level synthesis (HLS) based hardware accelerator designed for FPGAs some frameworks for using GEMM based on open computer language (OpenCL) already exists [49].

5.2 Winograd Transformation

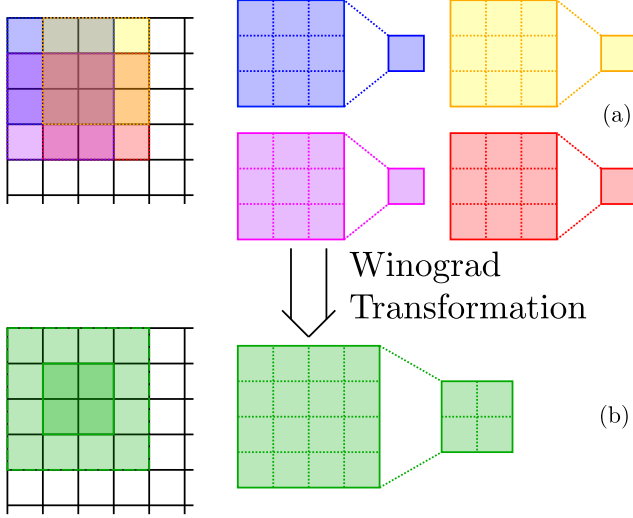


Fig. 9. Comparison of computations required for computing 2×2 feature map tile (a) and computations required computing the same 2×2 feature map tile using Winograd-Transformation (b).

It is a method found by Winograd. To goal is to reduce the number of multiplication required for implementing finite impulse response (FIR) filter. Winograd transformation transforms overlapping filter kernels into non-overlapping kernels [56]. The original method, developed for FIR filter, can be adopted for CNNs since the convolutional layer uses similar moving filter operations. Figure 9 visualizes that by transforming overlapping kernels into non-overlapping kernels, the arithmetic complexity of CNN is reduced by a factor up to 4 [2]. It shows that computing a 2×2 output matrix requires 36 multiplications, compared to 16 multiplications needed for the same 2×2 output matrix using Winograd-Transformation. The efficiency increases with the kernel size since more elements are overlapping [2].

The standard algorithm introduced by Winograd for a $F(2, 3)$ uses $2 \cdot 3 = 6$ multiplications [2, 56]:

$$F(2, 3) = \begin{bmatrix} d_0 & d_1 & d_2 \\ d_1 & d_2 & d_3 \end{bmatrix} \begin{bmatrix} g_0 \\ g_1 \\ g_2 \end{bmatrix} = \begin{bmatrix} m_1 + m_2 + m_3 \\ m_2 - m_3 - m_4 \end{bmatrix} \quad (5)$$

with:

$$\begin{aligned} m_1 &= (d_0 - d_2)g_0 \\ m_2 &= (d_1 + d_2) \frac{g_0 + g_1 + g_2}{2} \\ m_3 &= (d_2 - d_1) \frac{g_0 - g_1 + g_2}{2} \\ m_4 &= (d_1 - d_3)g_2 \end{aligned} \quad (6)$$

Equation (6) shows that the Winograd algorithm uses only 4 multiplications, therefore it is minimal for $\mu(F(2, 3)) = 2 + 3 - 1 = 4$. Additionally, 4 additions involving the data d , plus 3 additions, and 2 multiplications involving the filter parameter g , are required. However, these computations can be

pre-calculated. This algorithm can also be written in matrix form [2]:

$$Y = A^T [(Gg) \odot (B^T d)] \quad (7)$$

with \odot indicating an element-wise product. The matrices for $F(2, 3)$ are as follows [2]:

$$\begin{aligned} B^T &= \begin{bmatrix} 1 & 0 & -1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & -1 & 1 & 0 \\ 0 & 1 & 0 & -1 \end{bmatrix} \\ G &= \begin{bmatrix} 1 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} \\ 0 & 0 & 1 \end{bmatrix} \\ A^T &= \begin{bmatrix} 1 & 1 & 1 & 0 \\ 0 & 1 & -1 & -1 \end{bmatrix} \\ g &= [g_0 \quad g_1 \quad g_2]^T \\ d &= [d_0 \quad d_1 \quad d_2 \quad d_3]^T \end{aligned} \quad (8)$$

The 1D Equation (7) for $F(m, r)$ can be transformed into 2D for $F(m \times m, r \times r)$ as follows [2]:

$$Y = A^T [(GgG^T) \odot (B^T dB)]A \quad (9)$$

with g is a 3×3 filter and d is a $(m + r - 1) \times (m + r - 1) = 4 \times 4$ image tile [2]:

$$\begin{aligned} g &= \begin{bmatrix} g_{11} & g_{12} & g_{13} \\ g_{21} & g_{22} & g_{23} \\ g_{31} & g_{32} & g_{33} \end{bmatrix} \\ d &= \begin{bmatrix} d_{11} & d_{12} & d_{13} & d_{14} \\ d_{21} & d_{22} & d_{23} & d_{24} \\ d_{31} & d_{32} & d_{33} & d_{34} \\ d_{41} & d_{42} & d_{43} & d_{44} \end{bmatrix} \end{aligned} \quad (10)$$

The additions and subtractions to calculate the new 4×4 weight matrix (GgG^T) perform only once since they are constant for each layer. However, the new 4×4 matrix $(B^T dB)$ have to be computed before the multiplications. That leads to $16 \cdot 3 = 48$ additions or subtraction required during run time for each 2×2 output matrix.

A reverse transformation is done to get the results m . That leads to the following additions and subtractions with $e = (GgG^T) \odot (B^T dB)$:

$$\begin{aligned} m_{0,0} &= e_{0,0} + e_{0,1} + e_{0,2} + e_{1,0} + e_{1,1} + e_{1,2} + e_{2,0} + e_{2,1} + e_{2,2} \\ m_{0,1} &= e_{0,1} - e_{0,2} - e_{0,3} + e_{1,1} - e_{1,2} - e_{1,3} + e_{2,1} - e_{2,2} - e_{2,3} \\ m_{1,0} &= e_{1,0} + e_{1,1} + e_{1,2} - e_{2,0} - e_{2,1} - e_{2,2} - e_{3,0} - e_{3,1} - e_{3,2} \\ m_{1,1} &= e_{1,1} - e_{1,2} - e_{1,3} - e_{2,1} + e_{2,2} + e_{2,3} - e_{3,1} + e_{3,2} + e_{3,3} \end{aligned} \quad (11)$$

Equation (11) shows that each output element requires 9 additions or subtractions. That is equal to the original approach and doesn't influence the performance.

This can be summarized as follows: $F(2 \times 2, 3 \times 3)$ uses $4 \times 4 = 16$ multiplications and 16 additions, compared to that Equation (2) uses $2 \cdot 2 \cdot 3 \cdot 3 = 36$ multiplications, which reduces the amount of required multiplications by a factor of 2.25 [2].

However, the comparison is not that simple when using application-specific hardware such as

FPGAs and ASICs, considering that loading the data from memory is often the bottleneck. However, since less local data reuse is possible/necessary using Winograd transformation, it may not be possible to achieve a performance gain. Nevertheless, an advantage of the Winograd-Transformation is that, by using the local 2×2 output matrix, pooling can be done immediately. Since the pooling layer often follows a convolutional layer, Winograd-Transformation can reduce the required data movement.

Winograd transformation is used in GPUs, FPGAs and ASICs. It increases throughput as well as power efficiency. For the use of Winograd transformation in FPGA and ASIC implementations, it is recommended to consider it already in the hardware design phase, to optimally design the data flow and the processing elements.

5.3 Fast Fourier Transformation

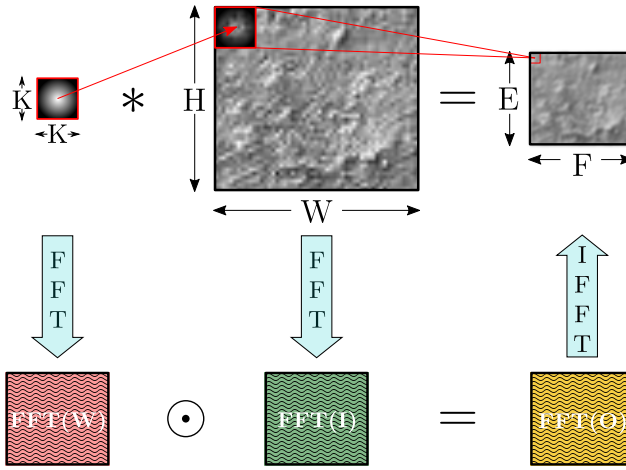


Fig. 10. Fast Fourier Transformation used to accelerate CNN. By applying FFT to the filter and to the input image the convolution operation can be transformed into an element-wise multiplication. Figure adopted from [45, Fig. 20].

Fast fourier transformation (FFT) is an algorithm to transform signals from image to time domain and vice versa. Convolutions in the time domain correspond to multiplications in the image domain and vice versa. Figure 10 visualizes the transformation of a convolution operation in a convolutional layer into an element-wise multiplication using FFT [1, 45]. A 2D FFT is computed by calculating first a 1D FFT of every row followed by a 1D FFT of every column [31]. FFT reduces the arithmetic complexity from $O(N_K^2 N_O^2) \rightarrow O(N_O^2 \log_2(N_O))$ with $N_K \times N_K$ is the kernel size and $N_O \times N_O$ is the output size [1, 45]. The benefits of FFT decreases with the kernel size [45]. It is also difficult to combine FFT with sparsity, which often offers higher profits [45].

Since modern CNNs mainly use small kernel sizes FFT has lost importance in modern hardware accelerators.

6 ACCELERATING NESTED LOOPS

The idea of neural network hardware accelerators is to parallelize operations. From an algorithmic point of view, matrix computations can be seen as nested loops. There are 3 approaches to accelerate the computation of these loops: Loop recording, unrolling, and pipelining. Loop tiling, on the other hand, deals with efficient memory allocation [31].

6.1 Loop recording

It is also known as loop interchange [9]. The goal is to rearrange the processing of nested loops, to avoid redundant memory accesses, and to maximize cache usage efficiency [31]. That can lead to a significant performance gain [9].

6.2 Loop unrolling

Independent loops can execute simultaneously using parallel hardware. The goal is to decrease the computation time by increasing the hardware utilization [9, 31]. Loop unrolling defines the optimal processing element (PE) structure [30].

6.3 Loop pipelining

Pipelining divides nested loops into sequential steps, which can be performed in parallel [31]. The possibility of pipelining depends on the architecture of the hardware accelerator.

6.4 Loop tiling

Loop tiling partitions the cache into tiles, which load as a whole chunk from memory, to avoid that data in the cache to supersede before the usage [31]. In terms of ASIC- and FPGA based hardware accelerators, loop tiling deals with the efficient usage of available on-chip memory to increase the locality of data [30].

7 GRAPHICS PROCESSING UNIT

GPUs enabled the great breakthrough of DNNs in the field of computer vision, since GPUs provide a massive parallelization of the computations required by DNNs such as AlexNet [26]. For this reason GPUs are the most widely used hardware accelerator used by DNNs architects [1, 34, 45]. Compared to ASICs and FPGAs, GPUs provide a high throughput and are easy to use. However, recent innovations in the structure of neural networks, such as sparsity, BNNs and SNNs result in a weak algorithmic efficiency using GPUs. This structures can be handled more efficiently using FPGAs and ASICs [34, 57].

7.1 Frameworks

GPUs are very popular since many frameworks already exist which offer a high-level API for creating, training, and testing neural networks. Some examples are: Caffe, Torch, Theano, Caffe2, PyTorch, TensorFlow, MXNet, CoreML, CNTK, and TensorRT [52]. This simplification and the support of training algorithms due to the floating-point computation is one of the reasons why most DNNs, which have won the ImageNet competition, use GPUs [1, 21, 26, 44, 46, 47, 57].

7.2 Architecture

GPUs consists of many floating-point arithmetic units for vector processing in combination with a high bandwidth memory. This is essential for DNNs, since most DNNs are based on dense floating-point general matrix multiplications (GEMMs) of 32-bit floating-point data [34, 45]. Figure 11 visualizes a typical temporal SIMD/SIMT structure used by GPUs. Temporal SIMD architectures utilize a centralized control for a large number of ALUs which share registers and memory. In temporal architectures, each ALU can fetch data from memory and load data back to memory. However, a single ALUs in a GPU cannot communicate with each other directly. For this reason, communication between individual ALUs is done using shared memory. This bottleneck increases the memory traffic, which is power intensive and decreases the possible throughput [45]. Compared to that spatial structures implemented in FPGAs and ASICs avoid this bottleneck [45].

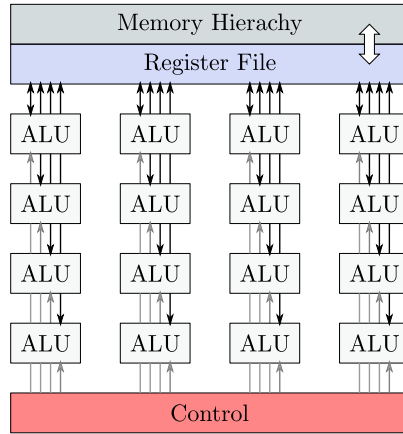


Fig. 11. Temporal SIMD/SIMT architecture. Adapted from [45, Fig. 17].

7.3 Embedded GPUs

Embedded GPU based hardware accelerator provides an easy to use, cost-efficient platform for fast development since many popular frameworks for neural networks are also suitable for embedded GPUs. Platforms like NVIDIA Jetson combine an ARM CPU with a small power-optimized GPU, for parallel image processing [23].

7.4 Overview

Table 1 presents an overview of existing Nvidia GPUs. It shows that the introduction of tensor processing units, specially designed for AI, leads to huge performance gains. Besides, smaller data types are supported, which offers an additional performance gain.

8 APPLICATION SPECIFIC INTEGRATED CIRCUITS

ASICs provide the greatest degree of freedom for designing a neural network hardware accelerator. ASIC designs require a huge effort and need an in-depth knowledge of chip design and neural networks [31, 53]. However, due to the great design freedom of ASICs, ASIC based hardware accelerators can outperform GPUs and FPGAs in both speed, and power efficiency [53, 57]. The big disadvantage of ASICs is a long time to market, and the high initial costs, which hinder them from keeping up with the rapid progress in the field of deep neural networks [53].

Table 2 gives an overview of existing ASIC based hardware accelerators. It shows that DaDianNao achieves the highest performance density and has moderate power consumption. However, the metrics are given cryptically as a multiple of the GPU performance.

8.1 DianNao family

The architecture of DianNao uses global input data, weight, and output data buffers in combination with a neural function unit (NFU) for computations. The NFU consists of 16-bit fixed-point multipliers in the first stage, an adder tree in the second stage, followed by a non-linearity unit. The second stage additionally contains shifters and max operators for pooling layers. The goal of DianNao is to support neural networks of any scale [5, 6]. In DaDianNao, adding additional buffers next to the neurons reduces the data movement. Additionally, the single computational block is split into four identical NFUs. DaDianNao targets high performance [6, 7]. ShiDianNao improves the DianNao structure by using the fact that in CNNs a whole vector shares the same weight. This

Table 1. Overview of GPU-based hardware accelerator implementations.

Name	Area [mm ²]	feature size [nm]	Quanti- zation	Bit width	Tensor unit	Throughput [TOPS] ^(a)	Freq. [MHz]	Power [W]	$[\frac{\text{GOPS}}{\text{mm}^2}]^{(b)}$
V100 ¹ [37]	815	12	float	64		7.8	1530	300	9.57
			float	32		15.7	1530	300	19.26
			mixed	32-8	X	125	1530	300	153.37
T4 ¹ [38]	545	12	float	32		8.1	1590	70	14.81
			float	16	X	65	1590	70	119.26
			fixed	8	X	130	1590	70	238.53
RTX 2080 ² [38]	545	12	fixed	4	X	260	1590	70	477.06
			float	32		10.6	1710	225	19.45
			float	16	X	84.8	1710	225	155.6
RTX 3080 ² [39]	628.4	8	fixed	8	X	169.6	1710	225	311.19
			fixed	4	X	322.2	1710	225	591.2
			float	32		29.8	1710	320	47.42
TX 1 ³ [36]	N/A	20	float	16	X	119	1710	320	189.3
			fixed	8	X	238	1710	320	378.7
			fixed	4	X	476	1710	320	757.47
AGX Xavier ⁴ [35]	N/A	12	float	32		0.512	998	15	N/A
			float	16		1.024	998	15	N/A
AGX Xavier ⁴ [35]	N/A	12	float	16	X	11	1370	30	N/A
			int	8	X	22	1370	30	N/A

(a) The fictitious peak performance is used due to a lack of comparable data

(b) Performance density

NVIDIA ¹Tesla, ²GeForce, ³Jetson (Tegra), ⁴Jetson (Volta)

condition is used to reduce the required DRAM accesses and to reduce the power consumption, therefore [6].

8.2 EIE: Efficient Inference Engine

EIE is designed for operating on sparse networks. That brings the benefit that faster and more power-efficient SRAM replaces power-hungry DRAM memory. The architecture of EIE utilizes a central control unit, which controls a PE array. Each PE computes one slice of the compressed network. The central control unit broadcasts non-zero activation values to the PEs, which use an activation queue, for buffering input data [18].

Table 2. Overview of existing ASIC-based hardware accelerator implementations.

Name	Area [mm ²]	feature size [nm]	Quanti- zation	Bit width	Throughput [GOPS] ^(a)	Frequency [MHz]	Power [mW]	[$\frac{\text{GOPS}}{\text{mm}^2}$] ^(b)
DaDianNao [7]	4335*	28	fixed	16	1586288*	606	48380*	366*
EIE [18]	40.8	45	fixed	16	102	800	590	2.5
Cambricon-X [61]	6.38	65	fixed	16	544	1000	954	85.26
Eyeriss [8]	16	65	fixed	16	67.2**	200	278	4.2
TPU [25]	331***	28	fixed	8	92000	700	40000	69.48

(a) The fictitious peak performance is used due to a lack of comparable data

(b) Performance density

*For 64 nodes. In paper only stated as a multiple of NVIDIA K20M

**Stated in the paper in GMAC. Therefore the double value is used

***4 dies

8.3 Eyeriss

The core architecture of Eyeriss is a spatial array with 168 PEs in a 12×14 shape, including a 108-kB GLB, an RLC CODEC, and a ReLU module. Each PE can communicate with the neighboring PEs or the GLB. A PE contains a MAC unit a local buffer called spads and a control unit [8].

8.4 Tensor processing unit

A matrix multiply unit, which includes 256×256 MACs, is the core element of the tensor processing unit. The multiplication unit has an input bit width of 8, the accumulation unit a bit width of 16. The matrix multiplication unit is followed by 4096, 256-element, 32-bit accumulators. Additionally, a dedicated activation module for applying the non-linearity and a dedicated normalization and pooling module is used. A 24 MiB on-chip buffer is used for storing the results. The weights are loaded directly from DRAM into a specialized weight FIFO [25].

9 FIELD PROGRAMMABLE GATE ARRAYS

FPGAs offer a cost-efficient method to develop custom hardware solutions since FPGAs provide the possibility to implement a reprogrammable logic circuit. Realizing a task in a logic circuit offers the advantage of excessive parallelization compared to a software solution in a CPU, which leads to a considerable speed advantage and higher power efficiency [16, 54]. Additionally, the implementation of an accelerator in hardware reduces the necessary memory accesses since the intermediate results can be passed on directly to the next PE. In addition to configurable logic blocks (CLBs), modern FPGAs implement hundreds to thousands of digital signal processors (DSPs) and various options for storing data locally. Embedded block-RAMs (BRAMs) and FIFOs can be used for larger amounts of data. In contrast, LUTs are well suited to be used as distributed RAM for smaller amounts of data. That enables FPGA based neural network implementations to minimize the delay due to memory access.

For this reason FPGAs offer the possibility to realize a low latency inference in combination with a high energy efficiency ($\sim 10 - 50\text{GOP/s/W}$) [30]. All in all, FPGAs represent a tunable balance between performance and energy consumption [52]. However, greater freedom of design also inevitably leads to a more complex solution and more development effort.

9.1 Algorithmic optimization in FPGA based hardware accelerators

FPGA based hardware accelerators widely use the algorithmic optimization methods described in Section 5. However, the efficiency of these methods depends on the structure of the hardware accelerator. GEMM is well suited for architectures that support batching, such as single computation engines. Nevertheless, for streaming architectures, GEMM is not applicable since multiple batches would require an enormous tile size of the buffers [30].

Basically FPGAs would be well suited for the use of FFT. However, FFT is more efficient for large kernel sizes and the trend in DNN structures is towards small kernel sizes [21, 41, 45, 46].

9.2 Mapping Convolutional Neural Networks on FPGAs

The task of mapping a DNN having millions of parameter and ten to hundred layers is a complex multidimensional optimization problem [30]. The goal is to minimize off-chip memory access since these lead to high energy consumption and latency [30, 45]. However, high hardware utilization is required to achieve high throughput. That is achieved by using loop optimization methods such as loop unrolling, loop tiling, loop recording, and loop pipelining (See Section 6).

On-chip memory is usually not sufficient for the large amounts of data required for input data, temporary results, and weights of a DNNs. For this reason, the following three typically memory hierarchies are used: External memory, on-chip buffers (cache), and registers for the individual PEs. Since full unrolling of large DNNs is often not possible due to the limited resources, intermediate results are buffered locally. The less intermediate results are buffered, and the earlier they can be processed, the less data movement is required. That leads to fewer off-chip memory accesses, and the efficiency increases, therefore. Additionally, suitable tile size is necessary to store the intermediate results [30].

A difficulty for the use of FPGAs is the complex design flow to get an efficient implementation of a neural network structure in an FPGA. High level synthesis (HLS) based on widely used programming languages such as C, C++, and OpenCL promises to simplify and accelerate the design flow. These offer the possibility of automating the design flow based on the parameters of each layer. That enables neural network designers without any hardware experience to create a customized hardware implementation [52]. Venieris et al. provide an overview of existing CNN-to-FPGA tool flows, and features such as performance, arithmetic precision, portability supported NN-models, and optimization is compared [52]. However, abstraction comes with a loss in efficiency [33].

9.3 Hardware architectures

Typical automated toolflows use one of following 2 hardware structures: Single computation engines and streaming architectures [52].

9.3.1 Single computation engines. Single computation engines are usually implemented in a systolic array [52]. Figure 12 visualizes the typical structure of systolic arrays. It shows that in a systolic array, each PE can exchange data with the adjacent PEs. Which enables systolic arrays to achieve a high data rate between the PEs. The PEs of systolic arrays can either consist of MAC-units, convolutional PEs which is shown in Figure 12b or of SIMD units [52, 55]. However, weights, inputs, and outputs are transferred from or to the external memory, which leads to a bottleneck. Adding

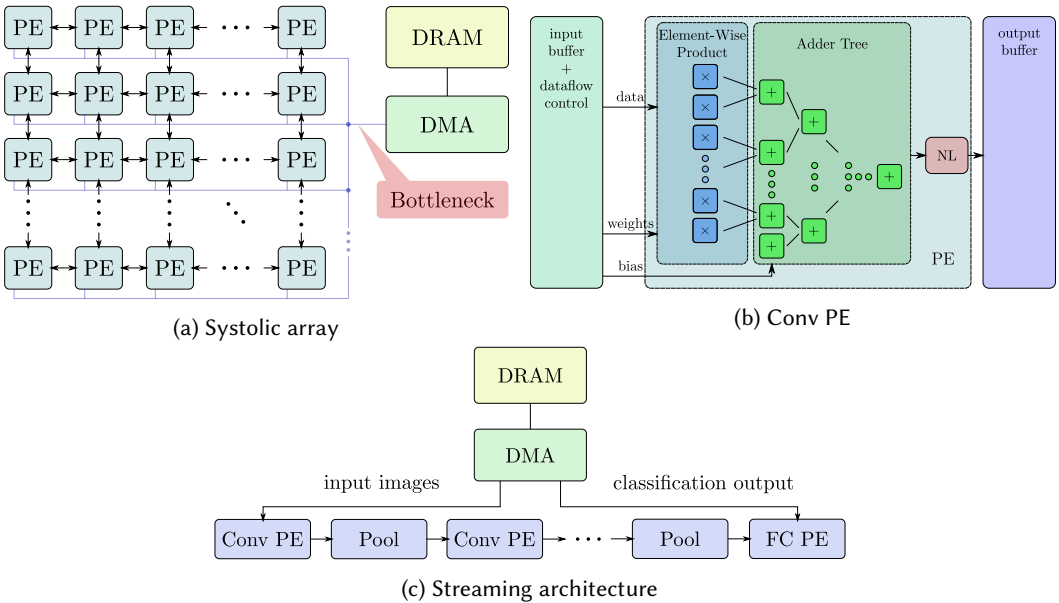


Fig. 12. Typical hardware architectures used in an FPGA-based hardware accelerator. a Typical structure of systolic array. PEs can communicate with adjacent PEs. Weights, input, and output data are transferred using DMA. Adapted from [1, Fig. 4]. c Typical structure of a streaming architecture. Each PE computes an independent layer. The data is pipelined through the accelerator. Adapted from [52, Fig. 1]. b Typical structure of convolutional processing element. It consists of N multipliers depending on the kernel size, followed by an adder tree, to sum, the intermediate results and a non-linearity element. Adapted from [40, Fig. 4].

additional buffers for weights and input data can reduce problem[55]. The benefit of using systolic arrays is that they can adapt easily for different DNN structures. However, the flexibility comes with the price of reduced efficiency, which leads to a very different throughput using different DNNs [15, 51].

9.3.2 Streaming architectures. Figure 12c visualizes the structure of streaming architectures, which implement an independent specialized PE for each layer. The data is pipelined through the hardware accelerator. However, it has to be ensured that the PEs are designed so that the calculations of individual layers take the same amount of time to prevent idling. The weights are stationary required for a specific layer, which reduces global data transfer. That enables streaming architectures to implement a highly efficient accelerator. However, compared to single computation engines, streaming architectures are less flexible since they are optimized for a single DNN structure. Additionally, it is difficult to increase or decrease the size of the accelerator, if the FPGA type changes [52].

Processing elements. MACs are heavily used in neural networks since it is the basic operation for fully connected and convolutional layer. As the name indicates MACs consists of a multiplier and an accumulator [45]. However, when accumulating the results of different input channels of a CNN, an accumulation without a multiplication is required. In this case, costly multipliers are unused. For this reason, especially in streaming architectures, a specialized convolutional PE is used, which consists of N multipliers depending on the kernel size, followed by an adder tree, to sum the intermediate results of each input channel. Usually, a non-linearity unit follows the adder

tree [30, 40]. Additionally, the latency decreases by using Conv-PEs since less pipelined stages are required. Figure 12b visualizes the structure of a convolutional processing element [40].

9.4 Overview of existing implementations

In 2015 Zhang et al. introduced a 32-bit floating-point single computation engine on a high-end Xilinx FPGA using HLS. The accelerator uses 64 two-level unrolled PEs operating with 100MHz. Each PE has 7 inputs for activations and weights and 1 bias input. Additionally, buffers are applied to input and output. The size of on-chip buffers is determined using loop tiling. A MicroBlaze CPU controls the data flow of the hardware accelerator. The system has a power consumption of 18.61W. Due to the high precision PEs, no accuracy loss occurs. However, the high precision results in a low throughput of 61.62GOPS [58]. For this reason, this design has compared to others a poor implementations efficiency.

In 2016 Zhang et al. created a hardware-software co-design library called Caffeine, which reduces the design effort for creating a hardware accelerator based on Caffe. It has similar performance metrics compared to the previous version. However, in this version 16 bit fixed point quantization is available, which increases the performance by a factor of 3.8.

Also, in 2016 Qiu et al. presented a 16-bit fixed-point single computation engine specialized for convolution operations. This structure brings the benefit that convolutions computations are more efficient. However, due to the non-optimal data flow for FC layers, they can only be calculated with low-performance [40]. The overall throughput is similar to the 16-bit fixed-point implementation introduced by Zhang et al. [40, 59]. Additionally, the throughput of a 4/8 bit Quantization is estimated, which would achieve a $2.6 \times$ higher throughput with a reasonable low accuracy loss of 0.4% [40].

Wei et al. created an automated systolic array architecture synthesis in 2017. The unique thing in this implementation is that each PE contains a SIMD vector accumulator. The systolic array structure aims to maximize the on-chip data reuse. For this reason, additional buffers for input, output, and weights are placed at the systolic array edges [55]. This implementation achieves compared to Zhang et al. a very high throughput of 461GOPS using 32-bit floating-point, and 16-bit fixed-point implementation [55, 58, 59]. However, a comparison of the implementations is difficult since different FPGA vendors are used which utilize different DSP structures.

Aydonat et al. use 16-bit floating-point operations for accumulation and 16 bit fixed point operation for multiplication to reduce the overhead from using FP32 DSPs. The hardware accelerator structure uses serial PEs, optimized for computing Winograd transformed feature maps [3]. This design outperforms the 16-bit implementation of Wei et al. by a factor of 1.18 in terms of throughput [3, 55].

Liang et al. introduced in 2018 a BNN for FPGA. Bit-level XNOR and shift operations replace costly multiplications, due to the binarization of both weights and activation. That enables a massive boost in throughput. However, an extensive drop in accuracy occurs for complex DNNs such as AlexNet. The use of full precision computations in the last layer counters this effect. Nevertheless, an accuracy loss of 14.7% occurs using AlexNet [29].

Chung et al. and Fowers et al. presented the Brainwave neural processing unit (NPU) in 2018. Which is a single thread SIMD instruction set architecture paired with 96000 MACs for vector-matrix and vector-vector operations. The target of the hardware accelerator is a maximal performance for cloud-scale real-time applications. The quantization is a 9-bit pseudo floating point quantization,

including a sign bit, a 5-bit exponent, and a 2-5 bit mantissa. While 128 numbers share a common exponent. This quantization results in an accuracy loss of 1-2%. When performing ResNet-50 on an Arria 10 1150, the brainwave NPU achieves a lower real-time latency (batch 1) than an Nvidia P40 GPU using int8 quantization [11, 14].

Nguyen et al. present a real-time YOLO object detection streaming architecture on a Xilinx XC7V485. A specialized PE is implemented for each layer. The images are pipelined through the processing elements. A binary quantization of the weights and a 4 to 6-bit quantization of the activations are used to minimize the footprint of the accelerator. That leads to an accuracy loss of 2.5%. This implementation achieves a $1.24\times$ higher frame rate with an $11.5\times$ higher energy efficiency compared to an NVIDIA GTX Titan X GPU using Sim-YOLO-v2 [33].

Table 3 provides an overview of the introduced related FPGA based hardware accelerator implementations. It compares the design methodology and the quantization with their influence on the throughput and utilization efficiency to estimate implementation efficiency. The throughput divided by the available logic elements (LEs) is used as a measure for the utilization efficiency because it shows how efficiently the available resources are used to generate maximum throughput. In some related works, the throughput divided by the number of DSPs is used for implementation efficiency [52]. However, the number of DSPs and LEs correlate anyway, which leads to a similar comparison. Historically, a single LE is defined as a four-input LUT, a programmable register, and a carry chain. However, in modern FPGAs LUTs usually have more than 4 inputs. For this reason, an adjustment factor is used for the conversion into a fictitious number of 4 input LUTs. Nevertheless, a comparison of implementations using different FPGA-vendors is difficult anyway.

From Table 3 it can be observed that the influence of using a HLS or register transfer level (RTL) based design methodology has a minor impact on the implementation efficiency. Furthermore, quantization is crucial for implementation efficiency. However, the efficiency gain comes with an accuracy loss. The throughput gain is much higher than the accuracy loss if a suitable combination of weight and activation bit width is used, though. Furthermore, latest implementations show that FPGAs are able to challenge GPUs in real-time computation tasks [14, 33].

10 COMPARISON

The decision of which neural network hardware accelerator is the most suitable depends on the application. For this reason, a separated comparison is done in terms of throughput, latency, accuracy, power, and usability. Since accuracy, throughput, and energy cost is a trade off [16]. Based on that, the aim of the comparison is not to determine which platform is generally the best, but rather which platform should be selected for which application.

A distinction between throughput and latency is necessary since these parameters can differ greatly depending on the hardware accelerator and batch size. For real-time applications latency is the most important measure. Throughput is of minor interest. The opposite is true for all other applications.

10.1 Throughput

Due to the high clock frequency, the large die size, and the SIMD architecture, which provides high parallelism, GPUs can offer up to some Tera FLOP/s throughput [25, 29]. Due to the use of specialized PEs and a lower bit width, ASIC based implementations achieve higher throughput compared to general-purpose GPUs. However, the introduction of specialized tensor processing units in GPUs, Table 2 and Table 1 show that GPUs achieve a similar performance density compared to TPU using 8 bit fixed point quantization. Nevertheless, TPU is more power-efficient, and the die

Table 3. Overview of existing FPGA-based hardware accelerator implementations.

Article	Quantization	Bit width		Accuracy loss ^(a)	Throughput [GOPS]	Freq. [MHz]	Power [W]	FPGA		$\frac{\text{GOPS}}{\text{kLE}}$
		W	X					Type	# LE	
Zhang et al. [58]	float	32	32	0%	61.62	100	18.61	XC7VX485	485.7k	0.13
Zhang et al. [59]	float	32	32	0%	96	200	25	XC7U060	726k	0.13
	fixed	16	16	N/A	365	25				
Qiu et al. [40]	fixed	16	16	0%	187.8	150	9.63	XC7Z045	350k	0.54
		4	8	0.4%	495**		N/A			1.41
Hegde et al. [23]	fixed	16	16	N/A	11.5	180	19	XC7Z045	350k	0.03
Wei et al. [55]	float	32	32	0%	461	221.65	N/A	Arria 10	1150k	0.4
	fixed	8	16	2%	1171	231.85	N/A			1.02
Guo et al. [15]	fixed	16	16	0.06%	187.80	150	9.63	XC7Z045	350k	0.54
		8	8	0.62%	19.2	100	2	XC7Z020	85k	0.23
Aydonat et al. [3]	float	16	16	0%	1382	45	303	Arria 10	1150k	1.2
Liang et al. [29]	binary	1*	1*	0.46-14.7%	9396	150	26.2	Stratix V	695k	13.52
Chung et al. [11]	float	9	9	1-2%	39500	300	125	Stratix 10	2800k	14.1
Ma et al. [30]	fixed	16	16	2%	348	150	N/A	Stratix V	622k	0.56
					715	200		Arria 10	1150k	0.62
Venieris et al. [51]	fixed	16	16	N/A	48.53	125	<5	XC7Z020	85k	0.57
					155.81		<5	XC7Z045	350k	0.45
Nguyen et al. [33]	binary	1*	4-6*	2.5%	1877	200	18.29	XC7VX485	485.7k	3.36

(a) compared to 32bit floating-point

*Not applied to first and/or last layers.

** Estimated value.

is manufactured with a larger feature size.

Compared to GPUs and ASICs, FPGA based implementation achieve a lower throughput. That comes from the fact that the variable routing of an FPGA limits the bandwidth and is area intensive.

10.2 Latency

As mentioned before, latency is critical for real-time applications, which are implemented mainly in edge devices. In real-time applications, the ability to batch the input data is limited to meet the real-time requirements [11, 62].

How efficient GPUs process neural networks depends on the size of the matrices and thus on the batch size. For this reason, GPUs lose performance in real-time applications, since the GEMM based optimization is less efficient using a small batch size [11]. Therefore, ASIC or FPGA based

implementations, which do not rely on GEMM can achieve a better throughput to latency ratio. Chung et al. showed that FPGA based hardware accelerators can outperform an Nvidia Tesla Xp GPUs in real-time applications if fixed-point quantization is used [11]. Besides note that the higher design freedom of ASICs and FPGAs enables more efficient processing of compact models. However, for a more complex data flow additional area is necessary.

10.3 Accuracy

The accuracy of a DNN mainly depends on the neural network structure, the dataset, and the training methods. These parameters are independent of the hardware accelerator. Many neural network hardware accelerators use quantization to increase performance density. However, that influences the accuracy of the required computations and can lead to an accuracy loss. Generally, for both floating-point and fixed-point quantization, using a bit width of 16 bits or higher does not lead to an accuracy loss. Otherwise, retraining can be used to recover the accuracy to minimize the accuracy loss due to quantization. Retraining can even lead to a performance gain for some not optimized neural network structures such as AlexNet [64].

Table 3 gives an overview of the used bit width and resulting accuracy loss. It shows that, if a minor accuracy loss is permitted, a more aggressive quantization method, such as binary weights in combination with 6-bit activations, can be used. That leads to an enormous gain in terms of throughput and latency [33].

10.4 Power

Power consumption is a key-element for battery-powered embedded systems such as mobile devices, robots, etc. [29]. However, for data centers, power consumption is also of interest since energy and cooling costs make up a significant part of the running costs of a data center [32]. For this reason, FPGAs and ASICs become more attractive for data centers since they offer advantages in the utilization of computation capacity and high energy efficiency [11, 14]. FPGAs usually achieve more than 10 times higher energy efficiency compared to GPUs [29]. However, comparing Table 1 and Table 3 shows that this is only true if the GPU uses floating-point-32. Table 2 visualizes that except for TPU, all other ASIC implementations have a very high power efficiency. However, TPU has a better power efficiency compared to GPUs and high throughput. Therefore, TPU is a favorable solution for data centers.

For power-critical embedded devices, ASICs or FPGAs is a feasible choice.

10.5 Usability

For GPUs many frameworks already exist which offer a high-level API for creating, training, and testing neural networks (See Section 7). For using the tensor cores of a GPU a C++ library called TensorRT is required, which is integrated into TensorFlow. It also includes a parser to import existing models from most other neural network training APIs. For this reason GPUs provide the best usability compared to ASICs and FPGAs.

However, TPU supports the high-level TensorFlow framework, which greatly increases usability. Since the TPU introduced by Jouppi et al. uses int8 quantization, training is not recommended [25]. However, next-generation TPUs support floating-point operations and training, therefore. Since other ASIC implementations do not support high-level APIs they have poor usability. For FPGAs some frameworks based on HLS and OpenCL exist (See Section 9). Nevertheless, more in-depth knowledge of neural networks and hardware design is necessary for implementing an efficient FPGA based hardware accelerator. Just like ASICs, FPGAs usually do not support training. Compared to ASICs, FPGAs have the advantage that they have a short design cycle, which is crucial in the field of neural networks because progress is fast [54].

11 CONCLUSIONS

In the last years, the accuracy achieved by deep neural networks is increasing rapidly. For this reason, DNNs are widely used for many AI applications. However, the increasing accuracy bases on an exponential expansion of the parameter number, which leads to an exponential increase of computational load. This computational effort can hardly be handled by CPUs. For this reason, optimized neural network hardware accelerators are required. Optimization can be achieved through algorithmic optimization, quantization, and DNN compression. However, one has to be aware that the efficiency of optimization methods depends on the hardware accelerator's architecture, supported quantization, and dataflow. For this reason, a general assessment of the efficiency of individual techniques is not possible. Besides, quantization optimization and DNN compression methods can lead to an accuracy loss. However, post-optimization retraining can reduce accuracy loss.

Since general-purpose GPUs provide a huge number of floating-point ALUs, GPUs are well-suited for the training of neural networks. However, the SIMD/SIMT architecture of GPUs lead to many DRAM accesses. In combination with the use of power-hungry floating-point computations, the high number of DRAM accesses lead to the high power consumption of GPUs. For this reason, in modern GPUs, additional tensor processing units are added, which provide an improved data flow and a quantization down to int4.

In contrast to general-purpose GPUs, ASIC-based neural network hardware accelerators are extremely optimized for tensor computations. For this reason, ASICs provide both high throughput and a high power-efficiency. However, creating a fully custom neural network hardware accelerator requires a long design time and an in-depth knowledge of chip design and neural networks. Since the progress in the field of deep learning is fast, the time to market is essential. Additionally, a high-level API, which simplifies the implementation of a designed neural network, is required to enable all neural network architects to use the accelerator.

FPGAs offer a programmable logic circuit that enables a highly optimized design with comparable low design effort and initial costs. Additionally, some frameworks exist based on HLS, which automatically generates an optimized hardware implementation based on the neural network design. Recent advances in sparsity and quantization enabled FPGAs to achieve a throughput comparable to general-purpose GPUs while having a higher power efficiency. Besides, FPGAs offer the possibility to optimize the dataflow, to reduce the required memory bandwidth for sparse neural network structures. For this reason, FPGAs are a promising candidate for the use in real-time edge systems. They also offer the advantage over ASICs that they are reprogrammable. Which is a considerable benefit due to the rapid progress in the field of deep learning.

It can be concluded that FPGAs are a feasible choice for a real-time application with a limited power budget. ASICs are well suited for the use in data centers since the high initial costs and the long time-to-market is an acceptable risk to achieve a higher performance and energy efficiency. For all other applications, GPUs is the most feasible choice.

REFERENCES

- [1] Kamel Abdelouahab, Maxime Pelcat, Jocelyn Serot, and François Berry. 2018. Accelerating CNN inference on FPGAs: A Survey. *arXiv preprint arXiv:1806.01683* (2018).
- [2] Andrew Lavin and Scott Gray. 2016. Fast Algorithms for Convolutional Neural Networks. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2016), 4013–4021.
- [3] Utku Aydonat, Shane O'Connell, Davor Capalija, Andrew C. Ling, and Gordon R. Chiu. 2017. An OpenCL™ Deep Learning Accelerator on Arria 10. *Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable*

- Gate Arrays* (2017), 55–64. <https://doi.org/10.1145/3020078.3021738>
- [4] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. 2020. YOLOv4: Optimal Speed and Accuracy of Object Detection. *arXiv preprint arXiv:2004.10934* (2020).
 - [5] Tianshi Chen, Zidong Du, Ninghui Sun, Jia Wang, Chengyong Wu, Yunji Chen, and Olivier Temam. 2014. DianNao: A Small-Footprint High-Throughput Accelerator for Ubiquitous Machine-Learning. *ACM SIGARCH Computer Architecture News*. 42, 1 (2014), 269–284.
 - [6] Yunji Chen, Tianshi Chen, Zhiwei Xu, Ninghui Sun, and Olivier Temam. 2016. DianNao family: Energy-Efficient Hardware Accelerators for Machine Learning. *Commun. ACM* 59, 11 (2016), 105–112. <https://doi.org/10.1145/2996864>
 - [7] Yunji Chen, Tao Luo, Shaoli Liu, Shijin Zhang, Liqiang He, Jia Wang, Ling Li, Tianshi Chen, Zhiwei Xu, Ninghui Sun, and Olivier Temam. 2014. DaDianNao: A Machine-Learning Supercomputer. *47th Annual IEEE/ACM International Symposium on Microarchitecture* (2014), 609–622. <https://doi.org/10.1109/MICRO.2014.58>
 - [8] Yu-Hsin Chen, Tushar Krishna, Joel S. Emer, and Vivienne Sze. 2017. Eyeriss: An Energy-Efficient Reconfigurable Accelerator for Deep Convolutional Neural Networks. *IEEE Journal of Solid-State Circuits* 52, 1 (2017), 127–138. <https://doi.org/10.1109/JSSC.2016.2616357>
 - [9] Jian Cheng, Pei-song Wang, Gang Li, Qing-hao Hu, and Han-qing Lu. 2018. Recent advances in efficient computation of deep convolutional neural networks. *Frontiers of Information Technology & Electronic Engineering* 19, 1 (2018), 64–77. <https://doi.org/10.1631/FITEE.1700789>
 - [10] Francois Chollet. 2017. Xception: Deep Learning With Depthwise Separable Convolutions. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2017), 1251–1258.
 - [11] Eric Chung, Michael Haselman, Maleen Abeydeera, Logan Adams, Hari Angepat, Christian Boehn, Derek Chiou, Oren Firestein, Alessandro Forin, Kang Su Gatlin, Mahdi Ghandi, Jeremy Fowers, Stephen Heil, Kyle Holohan, Ahmad El Hussein, Tamas Juhasz, Kara Kagi, Ratna K. Kovvuri, Sitaram Lanka, Friedel van Megen, Dima Mukhortov, Prerak Patel, Kalin Ovtcharov, Brandon Perez, Amanda Rapsang, Steven Reinhardt, Bitu Rouhani, Adam Sapek, Raja Seera, Sangeetha Shekar, Balaji Sridharan, Gabriel Weisz, Lisa Woods, Michael Papamichael, Phillip Yi Xiao, Dan Zhang, Ritchie Zhao, Doug Burger, Adrian Caulfield, Todd Massengill, Ming Liu, Daniel Lo, and Shlomi Alkalay. 2018. Serving DNNs in Real Time at Datacenter Scale with Project Brainwave. *IEEE Micro* 38, 2 (2018), 8–20. <https://doi.org/10.1109/MM.2018.022071131>
 - [12] Matthieu Courbariaux, Itay Hubara, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. 2016. Binarized Neural Networks: Training Deep Neural Networks with Weights and Activations Constrained to +1 or -1. *arXiv preprint arXiv:1602.02830* (2016).
 - [13] By Lei Deng, Guoqi Li, Song Han, Luping Shi, and Yuan Xie. 2020. Model Compression and Hardware Acceleration for Neural Networks: A Comprehensive Survey. *Proc. IEEE* 108, 4 (2020), 485–532. <https://doi.org/10.1109/JPROC.2020.2976475>
 - [14] Jeremy Fowers, Kalin Ovtcharov, Michael Papamichael, Todd Massengill, Ming Liu, Daniel Lo, Shlomi Alkalay, Michael Haselman, Logan Adams, Mahdi Ghandi, et al. 2018. A Configurable Cloud-Scale DNN Processor for Real-Time AI. *2018 ACM/IEEE 45th Annual International Symposium on Computer Architecture (ISCA)* (2018), 1–14.
 - [15] Kaiyuan Guo, Lingzhi Sui, Jiantao Qiu, Song Yao, Song Han, Yu Wang, and Huazhong Yang. 2017. Angel-Eye: A Complete Design Flow for Mapping CNN onto Customized Hardware. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 37, 1 (2017), 35–47. <https://doi.org/10.1109/ISVLSI.2016.129>
 - [16] Kaiyuan Guo, Shulin Zeng, Jincheng Yu, Yu Wang, and Huazhong Yang. 2019. [DL] A Survey of FPGA-Based Neural Network Inference Accelerators. *ACM Transactions on Reconfigurable Technology and Systems* 12, 1 (2019).
 - [17] Yiwen Guo, Anbang Yao, and Yurong Chen. 2016. Dynamic Network Surgery for Efficient DNNs. *Advances in neural information processing systems* (2016), 1379–1387.
 - [18] Song Han, Xingyu Liu, Huizi Mao, Jing Pu, Ardavan Pedram, Mark A. Horowitz, and William J. Dally. 2016. EIE: Efficient Inference Engine on Compressed Deep Neural Network. *ACM SIGARCH Computer Architecture News* 44 (2016), 243–254. <https://doi.org/10.1109/ISCA.2016.30>
 - [19] Song Han, Huizi Mao, and William J. Dally. 2016. Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantization and Huffman Coding. *4th International Conference on Learning Representations, ICLR* (2016).
 - [20] Song Han, Jeff Pool, John Tran, and William Dally. 2015. Learning both Weights and Connections for Efficient Neural Network. *Advances in neural information processing systems* (2015), 1135–1143.
 - [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2016), 770–778.
 - [22] Zhezhi He and Deliang Fan. 2019. Simultaneously Optimizing Weight and Quantizer of Ternary Neural Network Using Truncated Gaussian Approximation. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2019), 11438–11446.

- [23] Gopalakrishna Hegde, Siddhartha, and Nachiket Kapre. 2017. CaffePresso: Accelerating Convolutional Networks on Embedded SoCs. *ACM Transactions on Embedded Computing Systems* 17, 1 (2017), 1–26. <https://doi.org/10.1145/3105925>
- [24] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, Quoc V. Le, and Hartwig Adam. 2019. Searching for MobileNetV3. *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)* (2019).
- [25] Norman P. Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa, Sarah Bates, Suresh Bhatia, Nan Boden, Al Borchers, Rick Boyle, Pierre-luc Cantin, Clifford Chao, Chris Clark, Jeremy Coriell, Mike Daley, Matt Dau, Jeffrey Dean, Ben Gelb, Tara Vazir Ghaemmaghami, Rajendra Gottipati, William Gulland, Robert Hagmann, C. Richard Ho, Doug Hogberg, John Hu, Robert Hundt, Dan Hurt, Julian Ibarz, Aaron Jaffey, Alek Jaworski, Alexander Kaplan, Harshit Khaitan, Daniel Killebrew, Andy Koch, Naveen Kumar, Steve Lacy, James Laudon, James Law, Diemthu Le, Chris Leary, Zhuyuan Liu, Kyle Lucke, Alan Lundin, Gordon MacKean, Adriana Maggiore, Maire Mahony, Kieran Miller, Rahul Nagarajan, Ravi Narayanaswami, Ray Ni, Kathy Nix, Thomas Norrie, Mark Omernick, Narayana Penukonda, Andy Phelps, Jonathan Ross, Matt Ross, Amir Salek, Emad Samadiani, Chris Severn, Gregory Sizikov, Matthew Snelham, Jed Souter, Dan Steinberg, Andy Swing, Mercedes Tan, Gregory Thorson, Bo Tian, Horia Toma, Erick Tuttle, Vijay Vasudevan, Richard Walter, Walter Wang, Eric Wilcox, and Doe Hyun Yoon. 2017. In-Datacenter Performance Analysis of a Tensor Processing Unit. *Proceedings of the 44th Annual International Symposium on Computer Architecture* (2017), 1–12. <https://doi.org/10.1145/3079856.3080246>
- [26] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. 2012. ImageNet Classification with Deep Convolutional Neural Networks. *Advances in neural information processing systems* (2012), 1097–1105.
- [27] Vadim Lebedev and Victor Lempitsky. 2016. Fast ConvNets Using Group-Wise Brain Damage. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2016), 2554–2564.
- [28] Edward H. Lee, Daisuke Miyashita, Elaina Chai, Boris Murmann, and S. Simon Wong. 2017. LogNet: Energy-efficient neural networks using logarithmic computation. *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (2017), 5900–5904.
- [29] Shuang Liang, Shouyi Yin, Leibo Liu, Wayne Luk, and Shaojun Wei. 2018. FP-BNN: Binarized neural network on FPGA. *Neurocomputing* 275 (2018), 1072–1086. <https://doi.org/10.1016/j.neucom.2017.09.046>
- [30] Yufei Ma, Yu Cao, Sarma Vrudhula, and Jae-sun Seo. 2018. Optimizing the Convolution Operation to Accelerate Deep Neural Networks on FPGA. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems* 26, 7 (2018), 1354–1367. <https://doi.org/10.1109/TVLSI.2018.2815603>
- [31] Sparsh Mittal. 2018. A survey of FPGA-based accelerators for convolutional neural networks. *Neural Computing and Applications* 32, 4 (2018), 1109–1139. <https://doi.org/10.1007/s00521-018-3761-1>
- [32] Sanaa Hamid Mohamed, Taisir E. H. El-Gorashi, and Jaafar M. H. Elmirghani. 2019. A Survey of Big Data Machine Learning Applications Optimization in Cloud Data Centers and Networks. *arXiv preprint arXiv:1910.00731* (2019).
- [33] Duy Thanh Nguyen, Tuan Nghia Nguyen, Hyun Kim, and Hyuk-Jae Lee. 2019. A High-Throughput and Power-Efficient FPGA Implementation of YOLO-CNN for Object Detection. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems* 27, 8 (2019), 1861–1873. <https://doi.org/10.1109/TVLSI.2019.2905242>
- [34] Eriko Nurvitadhi, Suchit Subhaschandra, Guy Boudoukh, Ganesh Venkatesh, Jaewoong Sim, Debbie Marr, Randy Huang, Jason Ong Gee Hock, Yeong Tat Liew, Krishnan Srivatsan, and Duncan Moss. 2017. Can FPGAs Beat GPUs in Accelerating Next-Generation Deep Neural Networks? *Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays* (2017), 5–14. <https://doi.org/10.1145/3020078.3021740>
- [35] Nvidia. [n.d.]. Jetson AGX Xavier and the new era of autonomous machines. https://info.nvidia.com/rs/156-OFN-742/images/Jetson_AGX_Xavier_New_Era_Autonomous_Machines.pdf
- [36] Nvidia. 2015. NVIDIA Tegra 4 Family GPU Architecture. https://www.nvidia.com/docs/IO/116757/Tegra_4_GPU_Whitepaper_FINALv2.pdf
- [37] Nvidia. 2017. Nvidia Tesla V100 GPU Architecture. <https://images.nvidia.com/content/volta-architecture/pdf/volta-architecture-whitepaper.pdf>
- [38] Nvidia. 2018. NVIDIA-Turing-Architecture-Whitepaper. <https://www.nvidia.com/content/dam/en-zz/Solutions/design-visualization/technologies/turing-architecture/NVIDIA-Turing-Architecture-Whitepaper.pdf>
- [39] NVIDIA. 2020. NVIDIA-ampere-GA102-GPU-Architecture-Whitepaper-V1. <https://www.nvidia.com/content/dam/en-zz/Solutions/geforce/ampere/pdf/NVIDIA-ampere-GA102-GPU-Architecture-Whitepaper-V1.pdf>
- [40] Jiantao Qiu, Sen Song, Yu Wang, Huazhong Yang, Jie Wang, Song Yao, Kaiyuan Guo, Boxun Li, Erjin Zhou, Jincheng Yu, Tianqi Tang, and Ningyi Xu. 2016. Going Deeper with Embedded FPGA Platform for Convolutional Neural Network. *Proceedings of the 2016 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays* (2016), 26–35. <https://doi.org/10.1145/2847263.2847265>
- [41] Joseph Redmon and Ali Farhadi. 2018. YOLOv3: An Incremental Improvement. *arXiv preprint arXiv:1804.02767* (2018).
- [42] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. 2018. MobileNetV2: Inverted Residuals and Linear Bottlenecks. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*

- (2018), 4510–4520.
- [43] Ahmad Shawahna, Sadiq M. Sait, and Aiman El-Maleh. 2018. FPGA-Based Accelerators of Deep Learning Networks for Learning and Classification: A Review. *IEEE Access* 7 (2018), 7823–7859. <https://doi.org/10.1109/ACCESS.2018.2890150>
- [44] Karen Simonyan and Andrew Zisserman. 2015. Very Deep Convolutional Networks for Large-Scale Image Recognition. *International Conference on Learning Representations* (2015).
- [45] Vivienne Sze, Yu-Hsin Chen, Tien-Ju Yang, and Joel S. Emer. 2017. Efficient Processing of Deep Neural Networks: A Tutorial and Survey. *Proc. IEEE* 105, 12 (2017), 2295–2329. <https://doi.org/10.1109/JPROC.2017.2761740>
- [46] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander A. Alemi. 2017. Inception-v4, Inception-ResNet and the Impact of Residual Connections on Learning. *Thirty-first AAAI conference on artificial intelligence* (2017).
- [47] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. 2015. Going Deeper With Convolutions. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2015), 1–9.
- [48] Mingxing Tan and Quoc V. Le. 2019. EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks. *Proceedings of the 36th International Conference on Machine Learning* 97 (2019), 6105–6114.
- [49] R. Tapiador, A. Rios-Navarro, A. Linares-Barranco, Minkyu Kim, Deepak Kadetotad, and Jae-sun Seo. 2017. Comprehensive Evaluation of OpenCL-based Convolutional Neural Network Accelerators in Xilinx and Altera FPGAs. *Advances in Computational Intelligence. IWANN 2017. Lecture Notes in Computer Science* 10306 (2017).
- [50] Amirhossein Tavanaei, Masoud Ghodrati, Saeed Reza Kheradpisheh, Timothee Masquelier, and Anthony S. Maida. 2018. Deep Learning in Spiking Neural Networks. *Neural Networks* 111 (2018), 47–63. <https://doi.org/10.1016/j.neunet.2018.12.002>
- [51] Stylianos I. Venieris and Christos-Savvas Bouganis. 2019. fpgaConvNet: Mapping Regular and Irregular Convolutional Neural Networks on FPGAs. *IEEE transactions on neural networks and learning systems* 30, 2 (2019), 326–342. <https://doi.org/10.1109/TNNLS.2018.2844093>
- [52] Stylianos I. Venieris, Alexandros Kouris, and Christos-Savvas Bouganis. 2018. Toolflows for Mapping Convolutional Neural Networks on FPGAs: A Survey and Future Directions. *ACM Computing Surveys (CSUR)* 51 (2018), 1–39.
- [53] Erwei Wang, James J. Davis, Ruizhe Zhao, Ho-Cheung Ng, Xinyu Niu, Wayne Luk, Peter Y. K. Cheung, and George A. Constantinides. 2019. Deep Neural Network Approximation for Custom Hardware. *Comput. Surveys* 52, 2 (2019), 1–39. <https://doi.org/10.1145/3309551>
- [54] Teng Wang, Chao Wang, Xuehai Zhou, and Huaping Chen. 2018. A Survey of FPGA Based Deep Learning Accelerators: Challenges and Opportunities. *arXiv preprint arXiv:1901.04988* (2018), 1–10.
- [55] Xuechao Wei, Yu, Cody, Hao, 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. *Proceedings of the 54th Annual Design Automation Conference 2017* (2017), 1–6.
- [56] Shmuel Winograd. 1990. *Arithmetic complexity of computations* (3. printing ed.). CBMS-NSF regional conference series in applied mathematics, Vol. 33. Society for Industrial and Applied Mathematics, Philadelphia, Pa.
- [57] Xiaowei Xu, Yukun Ding, Sharon Xiaobo Hu, Michael Niemier, Jason Cong, Yu Hu, and Yiyu Shi. 2018. Scaling for edge inference of deep neural networks. *Nature Electronics* 1, 4 (2018), 216–222. <https://doi.org/10.1038/s41928-018-0059-3>
- [58] Chen Zhang, Peng Li, Guangyu Sun, Yijin Guan, Bingjun Xiao, and Jason Cong. 2015. Optimizing FPGA-based Accelerator Design for Deep Convolutional Neural Networks. *Proceedings of the 2015 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays* (2015), 161–170. <https://doi.org/10.1145/2684746.2689060>
- [59] Chen Zhang, Guangyu Sun, Zhenman Fang, Peipei Zhou, Peichen Pan, and Jason Cong. 2016. Caffeine: Toward Uniformed Representation and Acceleration for Deep Convolutional Neural Networks. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 38, 11 (2016), 2072–2085. <https://doi.org/10.1109/TCAD.2017.2785257>
- [60] Qingchen Zhang, Laurence T. Yang, Zhikui Chen, and Peng Li. 2018. A survey on deep learning for big data. *Information Fusion* 42 (2018), 146–157. <https://doi.org/10.1016/j.inffus.2017.10.006>
- [61] Shijin Zhang, Zidong Du, Lei Zhang, Huiying Lan, Shaoli Liu, Ling Li, Qi Guo, Tianshi Chen, and Yunji Chen. 2016. Cambricon-X: An accelerator for sparse neural networks. *49th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)* (2016), 1–12. <https://doi.org/10.1109/MICRO.2016.7783723>
- [62] Xiaofan Zhang, Junsong Wang, Chao Zhu, Yonghua Lin, Jinjun Xiong, Wen-mei Hwu, and Deming Chen. 2018. DNNBuilder: an automated tool for building high-performance DNN hardware accelerators for FPGAs. *2018 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)* (2018), 1–8. <https://doi.org/10.1145/3240765.3240801>
- [63] Xiangyu Zhang, Jianhua Zou, Kaiming He, and Jian Sun. 2016. Accelerating Very Deep Convolutional Networks for Classification and Detection. *IEEE transactions on pattern analysis and machine intelligence* 38, 10 (2016), 1943–1955. <https://doi.org/10.1109/TPAMI.2015.2502579>
- [64] Chenzhuo Zhu, Song Han, Huizi Mao, and William J. Dally. 2017. Trained Ternary Quantization. *5th International Conference on Learning Representations ICLR* (2017).