Optimization of FPGA-based CNN Accelerators Using Metaheuristics

Sadiq M. Sait^{1, 2}, Aiman El-Maleh^{1, 2}, Mohammad Altakrouri¹, and Ahmad Shawahna¹

Department of Computer Engineering, King Fahd University of Petroleum and Minerals, Dhahran-31261, KSA.

²Interdisciplinary Research Center for Intelligent Secure Systems, King Fahd University of Petroleum and Minerals, Dhahran-31261, KSA.

{sadiq, aimane, g201707010, g201206920}@kfupm.edu.sa

Abstract-In recent years, convolutional neural networks (CNNs) have demonstrated their ability to solve problems in many fields and with accuracy that was not possible before. However, this comes with extensive computational requirements, which made general central processing units (CPUs) unable to deliver the desired real-time performance. At the same time, field-programmable gate arrays (FPGAs) have seen a surge in interest for accelerating CNN inference. This is due to their ability to create custom designs with different levels of parallelism. Furthermore, FPGAs provide better performance per watt compared to other computing technologies such as graphics processing units (GPUs). The current trend in FPGAbased CNN accelerators is to implement multiple convolutional layer processors (CLPs), each of which is tailored for a subset of layers. However, the growing complexity of CNN architectures makes optimizing the resources available on the target FPGA device to deliver the optimal performance more challenging. This is because of the exponential increase in the design variables that must be considered when implementing a Multi-CLP accelerator as CNN's complexity increases. In this paper, we present a CNN accelerator and an accompanying automated design methodology that employs metaheuristics for partitioning available FPGA resources to design a Multi-CLP accelerator. Specifically, the proposed design tool adopts simulated annealing (SA) and tabu search (TS) algorithms to find the number of CLPs required and their respective configurations to achieve optimal performance on a given target FPGA device. Here, the focus is on the key specifications and hardware resources, including digital signal processors (DSPs), block random-access memories (BRAMs), and off-chip memory bandwidth. Experimental results and comparisons using four well-known benchmark CNNs are presented demonstrating that the proposed acceleration framework is both encouraging and promising. The SA-/TS-based Multi-CLP achieves $1.31\times$ - 2.37× higher throughput than the state-of-the-art Single-/Multi-CLP approaches in accelerating AlexNet, SqueezeNet 1.1, VGGNet, and GoogLeNet architectures on the Xilinx VC707 and VC709 FPGA boards.

Index Terms—Convolutional Neural Network, FPGA, Metaheuristics, Simulated Annealing, Tabu Search, Combinatorial Optimization, NP-Hard Problems.

I. Introduction

ONVOLUTIONAL neural network (CNN) is a powerful method used for processing data with predefined grid-like topology, such as 1-D time series in speech recognition [1], 2-D image retrieval in face detection and recognition [2], [3], and in intelligent transportation systems [4]–[6], to name a few. The CNN proved its effectiveness when it was used to win the ImageNet challenge, an annual competition for visual object recognition, in 2012, by dropping the classification error record from 26% to 15%, which was a significant improvement at the time [7]. Since then, more studies and applications have emerged to enhance CNN application in different fields [8]–[10].

In general, CNNs consist of an input layer, an output layer, and multiple intermediate hidden layers. The key operations involved in the construction of CNN layers include convolution (CONV), pooling (POOL), and inner product. The CONV layer extracts the unique features from the input image. Specifically, the first CONV layer in the architecture captures low-level features, such as edges, colors, gradients, orientations, etc., while subsequent CONV layers adapt to extract higher-level features, resulting in a wholesome understanding of the processed image.

On the other hand, the POOL layer reduces the dimensionality of the processed data, extracts rotational and positional invariant features, and suppresses noise by applying the max, avg, or min functions. Last, but not least, the inner product layer, also known as the dense or the fully connected (FC) layer, is a layer whose neurons are connected to all the neurons of its preceding layer, hence the name. It is worth noting that the last FC layer in the architecture is also referred to as the classification layer because it is responsible for decision-making, such as the class score. Here, we must emphasize that there are other layers that can also be used in constructing

This preprint has not undergone peer review (when applicable) or any post-submission improvements or corrections. The Version of Record of this article is published in The Journal of Supercomputing, and is available online at https://doi.org/10.1007/s11227-022-04787-8.

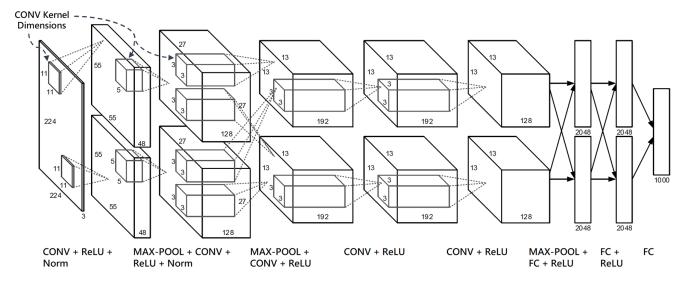


Fig. 1: An illustration of AlexNet architecture. The input/output feature maps to/from each layer(s) are shown as matrices with their respective dimensions and the operations performed are indicated below.

CNNs such as normalization (Norm) and rectified linear unit (ReLU). The Norm layer is used to smooth the data, whereas the ReLU activation layer helps in learning and modeling complex data.

CNN layers are represented by a set of matrices that reflect the main elements and the relationship between them according to the network structure. Specifically, each layer receives the output matrix of the previous layer as input, which is a set of 2-D arrays named feature maps (FMs). Then, it performs its operation to produce the FMs for the subsequent layer. Figure 1 shows a well-known CNN architecture called AlexNet [7]. It consists of 5 CONV layers each of which is followed by a ReLU activation function, interspersed by 2 Norm layers, 3 MAX-POOL layers, and concluded by 3 FC layers.

In CNNs, as the name implies, CONV layers are the most critical ones, and their operations constitute over 90% of the total computation time [11]. The CONV layer receives a set of N FMs from the previous layer. Then, it convolves the input FMs (IF) with a set of M small kernels, also known as filters, using a shifting window that slides over the IF with a stride of size S. The kernel is a matrix of numbers called weight parameters whose values are learned during the training phase by the backpropagation algorithm. Each kernel is used to compute one specific output FM, that is, the m-th kernel produces the m-th output FM as demonstrated in Figure 2. Thus, the number of output FMs (OF) equals the number of CONV layer kernels. Additionally, each kernel has one bias term (B). Each bias is added to every element in its corresponding output FM to produce the final OF, which form the IF for the next layer.

The number of CONV kernels, their size, and the number of channels are defined by the CNN designers depending on the type of convolution. Specifically, standard convolution and point-wise convolution employ M kernels of size $K \times K$ and 1×1 , respectively, each of which contains N channels [12]. Conversely, depth-wise convolution uses N single-channel kernels of size $K \times K$, where each kernel is applied to one input FM to produce the corresponding output FM, i.e., M = N. To determine the value of the neuron at position (r, c) of m-th output FM, the following computation is performed

$$\mathbf{OF}[m][r][c] = \mathbf{B}[m] + \sum_{n=0}^{N-1} \sum_{i=0}^{K-1} \sum_{j=0}^{K-1} \mathbf{W}[m][n][i][j] \times \mathbf{IF}[n][S \times r + i][S \times c + j]$$
(1)

Here, $\mathbf{B}[m]$ represents the bias parameter of the m-th kernel, $m \in [M]$, $\mathbf{W}[m][n][i][j]$ denotes the weight parameter at position (i,j) in the n-th channel of the m-th kernel, $n \in [N]$, i and $j \in [K]$, $\mathbf{IF}[n][x][y]$ indicates the neuron at position (x,y) of the n-th input FM, and S is the stride size. By applying Equation (1) to compute all neurons of OF, we can see that the main structure of CONV layer consists of 6 nested loops as shown in Algorithm 1. One can also note that these loops contain a massive number of the multiply-accumulate (MAC) operations. Precisely, each layer performs $M \times R \times C \times N \times K^2$ MAC operation. Overall, AlexNet requires about 1.46 Giga operations (GOPs) to process a single RGB image of size 224×224 pixels.

The intensive computational requirements of CNNs limit their real-time applications on general central processing units (CPUs). On the other hand, the computational efficiency of graphics processing units (GPUs) and field-programmable gate arrays (FPGAs) make them excellent platforms for accelerating CNNs. However, the critical need for lower power

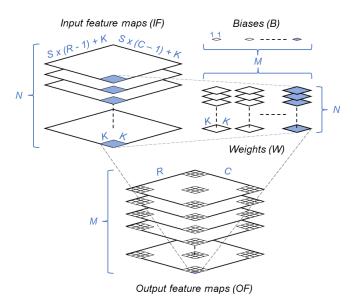


Fig. 2: An illustration of the convolution layer. The blue feature in the last output feature map is computed by taking the dot-product of the blue weights with the blue highlighted portion of the input feature maps and then adding the blue shaded bias value to the result.

consumption in today's applications, due to CNN deployment on battery-powered devices such as drones and Internet of things (IoT) [13], [14], makes FPGAs more suitable for CNN acceleration [8]. This is because of the high power efficiency, also known as performance per watt, that FPGAs offer.

With regards to model parameters, each standard CONV layer learns $M \times N \times K^2$ weight parameters. For example, AlexNet's first CONV layer receives 3 input FMs, uses 11×11 kernels, and produces 96 output FMs, resulting in 34,848 parameters. Overall, AlexNet has over 60 million parameters which need about 240 MB of memory space. The memory requirements for FMs and model parameters exceed what commercially available FPGAs can provide in on-chip memory. Thus, they must be stored in off-chip memory and transferred to on-chip memory during computation on need basis. The significant amount of storage and external memory bandwidth required become a performance bottleneck. Therefore, many attempts have been made to maximize the throughput of CNN applications by designing a hardware accelerator known as a convolutional layer processor (CLP) [15]–[19].

The CLP optimizes the implementation of CONV layers by applying loop transformations to Algorithm 1. Specifically, the loop unrolling technique is used to maximize the parallelism of CONV computations by replicating the hardware resources of MAC unit. Furthermore, loop tilling and local memory promotion techniques are adopted to reduce off-chip memory accesses and maximize data sharing and reuse. Here, we must emphasize that the CLP design is parametrized by the number

```
Algorithm 1: Convolution Algorithm.
```

```
Input: The input feature maps (IF), weight (W) and bias
            (\mathbf{B}) parameters, number of input feature maps (N)
            and output feature maps (M), size of each output
            feature map, rows (R) and columns (C), kernel size
            (K), and window stride size (S).
    Output: The output feature maps (OF).
   Procedure CONV (IF, W, B, N, M, R, C, K, S):
        for m \leftarrow 0 to M-1 do
 1
             for r \leftarrow 0 to R-1 do
2
                  for c \leftarrow 0 to C-1 do
 3
                       \mathbf{OF}[m][r][c] \leftarrow \mathbf{B}[m]
 4
                       for n \leftarrow 0 to N-1 do
 5
 6
                            for i \leftarrow 0 to K-1 do
                                 for j \leftarrow 0 to K-1 do
 7
                                      \mathcal{P} \leftarrow \mathbf{W}[m][n][i][j] \times
                                             \mathbf{IF}[n][S \times r + i][S \times c + j]
                                      \mathbf{OF}[m][r][c] \leftarrow \mathbf{OF}[m][r][c] + \mathcal{P}
10
                                 end
                            end
11
12
                       end
                  end
13
             end
14
15
        end
        return OF
16
```

of input/output FMs to/from CONV layers as well as their dimensions. However, different CONV layers, even those in the same architecture, vary considerably in their configurations (N, M, R, C, K, and S).

To overcome this problem, three main design schemes have been proposed in the literature. In the first approach, such as in [15], a CLP is modeled for each CONV layer. In this way, a CNN of L CONV layers is accelerated using L CLPs that process L independent images in a pipelined fashion. Even though this approach optimizes the computations of each CONV layer, it suffers from non-negligible latency and bandwidth overheads. This is due to the need for orchestrating the off-chip memory accesses for a large number of CLPs. Additionally, dividing the limited on-chip memory among many CLPs reduces the overall data locality. Last, but not least, implementing a dedicated controller for each CLP leaves them with insufficient resources for computation.

The second approach, on the other hand, designs a single, unified CLP based on the optimal parameters that achieve the lowest overall latency [16], [17]. The globally optimized CLP is then used to iteratively process CONV layers, one layer at a time. However, a CLP that provides the best performance across all layers is not necessarily optimal in utilizing its hardware resources. This is due to the radically varying CONV configurations. Considering the three CONV layers shown in Figure 3 as an example, one can note that the Single-CLP accelerator in Figure 3a is underutilized while processing L1 and the portions of L3.

To alleviate this issue, the third approach implements multiple CLPs and optimally distributes the hardware resources

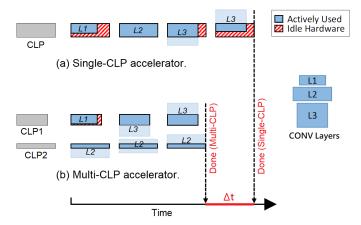


Fig. 3: The operation of (a) Single-CLP accelerator and (b) Multi-CLP accelerator on CNN with three CONV layers. The dimensions of the hardware (CLP, CLP1,and CLP2) and the CONV layers (L1, L2,and L3) are represented by the size and shape of the boxes.

among them [18], [19]. In Figure 3b, the same hardware resources used to implement the Single-CLP accelerator in Figure 3a are partitioned into two CLPs (CLP1 and CLP2). Thereafter, by mapping L1 and L3 to CLP1 and L2 to CLP2, the Multi-CLP design achieves Δt reduction in the overall execution time compared to the Single-CLP design.

The problem now is, how to determine the appropriate number of CLPs to use, on what basis to assign a CONV layer to a CLP, how many hardware resources each CLP can utilize, and how to determine the appropriate tiling factors to assign to each CONV layer? This intractable design space makes finding near-optimal configurations by exhaustive search algorithms almost impossible. Thus, intelligent techniques are needed to efficiently explore the design space for the optimal configurations of Multi-CLP design that improve the throughput of CNN applications.

Metaheuristics algorithms have proved their capabilities for finding optimal solutions for several NP-hard problems, with very efficient performance [20]. However, very few studies exist in the current literature that have proposed and exploited the use of metaheuristics algorithms in optimizing CNNs. Specifically, the standard genetic algorithm was employed in [17] to explore the design space for an optimal Single-CLP accelerator. On the other hand, the authors in [21] adopted the simulated annealing algorithm to improve the training of the LeNet-5 architecture.

In this paper, we present a systematic methodology for the optimization of the throughput of an FPGA-based accelerator with Multi-CLP implementation. Specifically, we present an analytical and empirical design scheme for accurate modeling of cost, in terms of hardware resources used, and performance for a Multi-CLP design given its parameters. Then, we employ simulated annealing and tabu search metaheuristic algorithms

to search for the optimal Multi-CLP design, constrained by the given CNN architecture and target FPGA specifications. As a case study, we implemented CNN accelerators for four well-known benchmark architectures, namely, AlexNet [7], SqueezeNet 1.1 [22], VGGNet [23], and GoogLeNet [24] on VC707 and VC709 FPGA boards and compared them with previous approaches. Our implementation achieves a performance of 113.92 Giga floating-point operations (GFLOPs) under 100 MHz working frequency. The key contributions of this work are summarized as follows:

- Designing a CNN accelerator with loop tiling, loop unrolling, and loop pipelining techniques.
- Modeling CNN implementation and identifying the critical optimization variables for multiple CLPs design on FPGA platforms.
- Clarifying the complexity of finding the optimal configurations for multiple CLPs accelerator, i.e., defining the problem statement and the resultant intractable design space.
- Introducing a metric to quantify the dimensional mismatches between CLP dimensions and CONV layer dimensions, and use it to improve CNN throughput.
- Presenting computational performance and resource usage models for cost estimation of a candidate multiple CLPs design.
- Proposing a metaheuristic-based optimization framework that employs the estimation models to efficiently explore the design space using simulated annealing and tabu search algorithms, and find the optimal multiple CLPs design for CNNs on FPGAs based on performance and resources constraints.
- Validating the solutions obtained and comparing them with related works.

The rest of the paper is organized as follows. Section II reviews previous work related to accelerating CNNs using FPGAs. In Section III, we discuss the optimization techniques applied to achieve efficient implementation of CNNs on FPGA platforms. Furthermore, we present an analytical and empirical model to accurately predict the performance and hardware resources required for a given design variant. Section IV presents the methodology used for exploring the design space to achieve an efficient implementation of Multi-CLP hardware accelerator for CONV layers. The results of our experiments are presented and discussed in Section V, and the paper is concluded in Section VI.

II. RELATED WORK

In this section, we review the existing approaches targeting the optimization of CONV implementation on FPGA platforms. Although CNN structures perform well for their intended applications, they have the potential to be further optimized with minimal impact on accuracy. For example, the values of FMs and weight parameters are originally represented as 32-bit floating-point numbers. However, it has been demonstrated that fewer bits can be used to represent these values without a noticeable accuracy drop [25]. This reduces the hardware requirements for CNN implementation as well as reduces the inference latency.

The authors in [26] proposed an accelerator for LeNet-5 architecture to perform handwritten digits classification. The proposed strategy is based on three major aspects; loop parallelization to utilize resources, fixed-point data optimization to find the minimum number of bits that maintains accuracy level, and finally implementing MAC approximate units through logic blocks such as look-up tables (LUTs) and flip-flops (FFs) rather than using high-precision digital signal processors (DSPs). With these optimizations, the authors achieve less memory usage and reduced network latency. However, the solution is problem specific and deals with a relatively small CNN. Large CNNs would require advanced techniques to achieve the full potential of FPGA resources.

In [16], the authors adopted the polyhedral-based data dependence analysis [27] to optimize the computations and memory access operations in CONV layers. Specifically, they proposed an analytical design scheme to estimate the computational performance of a given CONV design. The loop transformations are employed to derive all possible CONV designs. The goal behind employing these techniques is to fully utilize the hardware resources provided by the target FPGA platform for effective acceleration. After enumerating all candidate solutions, the roofline performance model [28] is used to identify the optimal design for each layer.

Here, we must emphasize that different CONV layers have different structures, and therefore varying optimal loop unrolling and loop tiling factors. Considering AlexNet architecture discussed in Figure 1, the optimal unrolling factors $\langle T_n, T_m \rangle$ for the second and third CONV layers are $\langle 24, 20 \rangle$ and $\langle 5, 96 \rangle$, respectively, where T_n and T_m are the unrolling factors for the input FMs and output FMs, respectively. Hence, designing a CLP to accelerate these differently structured layers requires complex hardware implementation to reconfigure the computational units and their interconnects.

To overcome this issue, the authors in [16] designed the CLP based on the uniform cross-layer unrolling factors. More precisely, they found the unified unrolling factors $\langle T_n', T_m' \rangle$ that maximize the overall performance. Accordingly, the computational engine of the CLP is implemented as T_m' duplicated tree-shaped poly structures. These structures receive identical T_n' values from the input FMs and each structure receives T_n' weights from one of the T_m' kernels. Each structure

multiplies its data using $T_n^{'}$ digital multipliers. Then, an adder tree is used to accumulate the multipliers outputs as well as the previous partial result. Furthermore, they used the double-buffering technique to overlap the off-chip memory and CLP buffers data transfers with computation.

Suda et al. [17] proposed an OpenCL-based framework for accelerating CNN inference. Focusing on CONV layers, the authors reformulated the CONV operation into a matrix multiplication operation. Specifically, they organized CONV kernels as a 2-D matrix of size $M \times (N \times K \times K)$. Similarly, the input FMs were flattened and rearranged as a 2-D matrix of $(N \times K \times K)$ rows and $(R \times C)$ columns. In this way, the output of CONV layer is calculated by multiplying these two matrices.

To speed up CNN operations, the authors followed the same strategy in [16] and modeled each type of CNN layer using unified loop unrolling factors. Precisely, they unrolled the output matrix of CONV operation in both dimensions by the factor T_{out} . Thus, they designed a CLP consisting of $T_{out} \times T_{out}$ structures, each of which computes an output feature. Accordingly, kernel weights and input FMs matrices are tiled into blocks of size $T_{out} \times T_{out}$. On each iteration, a tile from the kernel weight matrix and a tile from the input FMs matrix are fetched into the on-chip memory. Then, each CLP structure performs T_{out} MAC operations on a row of kernel weights tile and a column of input FMs tile.

To further improve the throughput of CONV layer, the computations on the inputs to CLP structures were also unrolled using the factor T_{in} so as to perform T_{in} MAC operations, out of the T_{out} , in parallel. To find the optimal unrolling configuration, the authors model the execution time of each layer as a function of the unrolling factors. Then, the standard genetic algorithm was used to explore the design space for the minimum overall execution time considering FPGA resource constraints. Hence, this approach is referred to as GA-based Single-CLP. It is noteworthy that the high-level synthesis tool employed to compile their OpenCL codes to hardware restricted T_{in} to be from the set $\{1, 2, 4, 8, 16\}$ and T_{out} to be integer multiplicative of T_{in} .

The FPGA-based CNN accelerators discussed earlier employed a single globally-optimized CLP design to maximize the overall throughput. However, using a CNN accelerator with uniform unrolling factors leads to sub-optimal performance for some CONV layers due to their significantly varying dimensions, which affects the overall performance. For example, the cross-layer optimization for AlexNet layers in [16] increases the total execution cycles of the layer-based optimization by 44, 442 cycles. Moreover, following the methodology in [16] to derive the optimal Single-CLP design for SqueezeNet 1.1 implementation on Virtex-7 690T FPGA, the results show

that the dynamic utilization of CLP's MAC units is less than 77% [18].

To improve the throughput of the Single-CLP design, Shen et al. [18] proposed to partition the available hardware resources between multiple CLPs. In doing so, they introduced a two-step iterative algorithm that searches for candidate partitions in the first step and then optimizes the tiling factors of each CONV layer in the second step. Specifically, the algorithm starts with a predefined target performance aiming to find a design with such a performance.

During the first phase, it generates several candidate partitions of computational resources. For each candidate design, the number of CLPs is set to the partition size. Then, the methodology in [16] is adopted to compute $\langle T_n, T_m \rangle$ for each CLP. The last process in this phase distributes CNN layers on the adopted CLPs. To facilitate this process, the authors proposed to order memory-bounded layers based on their need for computation and communication, whereas computational-bounded layers are ordered based on the difference between the number of input FMs and output FMs.

With such ordering, the authors assume that similar behavioral layers will be neighboring. Thus, they constrain the algorithm to only assign neighboring layers to the same CLP. At the end of the first phase, the performance of candidate designs is evaluated. If no one meets the target performance, the algorithm slightly reduces the target performance and repeats the first phase. When a valid design is found, the algorithm moves on to the second phase and computes $\langle T_r, T_c \rangle$ based on the methodology in [16]. The parameters for the design with the minimum bandwidth requirement is then used to configure a generic CLP template in high-level synthesis.

In this work, we improve the previously discussed works by efficiently partitioning hardware resources between multiple CLPs. In doing so, we employ intelligent metaheuristics to help explore the intractable design space for multiple CLPs design that is more optimized than that achieved with conventional iterative algorithms. This results in better utilization of FPGA compute resources, improving CNN throughput.

Furthermore, unlike the work in [17] which constraints the unrolling factors to be from a predefined set, this work allows the unrolling factors to be whatever value that yields the highest performance. Last, but not least, the authors in [18] have specified the number of CLPs in multiple CLPs design to be 6 at most. This is because they want to keep the optimization time of their algorithm within an acceptable amount of time. This paper proposes a metaheuristic-based optimization framework that is efficient in exploring design space and takes less than a minute to provide an optimal Multi-CLP design on a single CPU.

III. CNN ACCELERATOR DESIGN ON FPGA PLATFORMS

In this section, we discuss the optimization techniques applied to achieve efficient implementation of CNNs on FPGA platforms. A hardware accelerator referred to as CLP is typically designed to improve CONV layer throughput. Here, the focus is on CONV layers because they are the most computationally intensive layers. Furthermore, we present an analytical and empirical model to accurately predict the performance and hardware resources required for a given design variant.

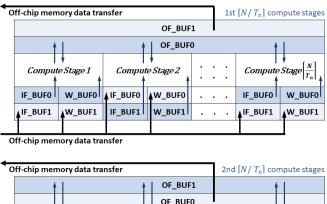
A. Optimizing CLP Computation and Memory Access

The computation of the CONV layer on FPGAs requires a memory space to store input FMs, output FMs, and all kernel weights. For instance, the second CONV layer in the VGGNet architecture demands 12.25 MB, 0.14 MB, and 12.25 MB for input FMs, kernel weights, and output FMs, respectively. However, FPGA platforms usually contain limited on-chip memory. For example, Xilinx Virtex-7 485T FPGA can store a maximum of 4.52 MB of data on-chip. Thus, we employ the loop tiling technique to overcome this issue. Initially, all the data required for calculation is stored in the external memory. Then, the CONV operation is iteratively performed on a small portion of data, called a tile or block. Each tile is loaded and cached in on-chip buffers before being fed into the computational engine of the CLP.

However, loop tiling opens several design challenges that affect CLP performance. Improper tiling may degrade the efficiency of data reuse and the parallelism of data processing. Therefore, deciding whether to tile a loop and with what factor plays an important role in the performance that can be achieved. Typically, CONV kernels are small in size, $K \leq 11$. Hence, the loops over the kernel dimensions, the loop iterators i and j in Algorithm 1, are not tiled.

On the other hand, the loops over the N input FMs, the M output FMs, and the dimensions, rows (R) and columns (C), of each output FM are tilted with the factors T_n , T_m , T_r , and T_c , respectively. In other words, each of these loops is transformed into two loops; an outer loop that iterates over the tiles, and an inner loop that iterates over the elements in each tile. Note that we denote the loop iterators over the tile elements of N, M, R, and C as nt, mt, rt, and ct, respectively.

Additionally, the CONV algorithm is a good candidate for parallelism. Thus, we adopt the loop unrolling technique to speed up CONV operations. Specifically, loop unrolling utilizes the available FPGA computation resources to maximize the parallelism of the MAC units. Even though loop unrolling considerably improves CLP throughput, it imposes complex connection topologies and affects CLP operating frequency. To



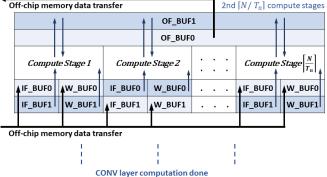


Fig. 4: CLP computation and ping-pong buffer structure.

mitigate these shortcomings, we only unroll input FMs and output FMs loops based on their tiling factors.

More precisely, the loops indicated by the iterators nt and mt are unrolled with the factors T_n and T_m , respectively. In this way, all their operations are performed in parallel. To avoid loop-carried dependence, we transformed unrolled loops to the innermost level. Moreover, the operations of the unrolled loops are fully-pipelined to improve system throughput. The CONV algorithm after optimization for loop tiling, loop unrolling, and loop pipelining is demonstrated in Algorithm 2. Note that the nested loops are reordered to maximize data reuse, thus reducing off-chip memory accesses. Determining the optimal value of tiling factors for each CNN architecture will be discussed later in Section IV.

The optimized CONV algorithm can be divided into two stages; the memory data transfer stage, and the data computation stage. In the data transfer phase, highlighted in red and blue, the input FMs on-chip buffer (IF_BUF) and weight parameters on-chip buffer (W_BUF) are filled with a block of input FMs and a block of kernel weights, respectively. Later, the output FM data block contained in the output FMs on-chip buffer (OF_BUF) is copied back to external memory as shown in line 23. Note that loop tiling factors control the size of these buffers, and thus the amount of data that is moved for each buffer refill or write-out.

On the other hand, the operations of the computation phase, shaded in yellow, are unrolled based on T_n and T_m tiling factors. Thus, loop tiling factors also control how CLP computational engine is constructed. Specifically, the computational engine

is modeled as T_m duplicated MAC tree tiles. Each of which receives T_n different weights but they all share the same T_n input features. Each MAC tree tile multiplies its data using T_n digital multipliers. Finally, an adder tree is used to accumulate the output of the multipliers with the previous partial result.

To alleviate, or even prevent, blocking of CLP computations due to external memory data transfer, CLP adopts the double-buffering scheme for all on-chip buffers. In other words, each on-chip buffer can be logically considered as two independent sets operating in a ping-pong manner as illustrated in Figure 4. In the first compute stage, the computational engine processes the features and weights from the input buffer set 0 (IF_BUF0 and W_BUF0). During the same time, the features and weights required for the second compute stage are copied from external memory to input buffer set 1 (IF_BUF1 and W_BUF1).

The next compute stage does the same but uses the opposite input buffers. That is why it is called the ping-pong buffer structure. With regards to the resulting output FMs, the output features from the first $\lceil N/T_n \rceil$ compute stages are stored in the output buffer set 0 (OF_BUF0). During these stages, the content of the OF_BUF1, which are the results of the previous $\lceil N/T_n \rceil$ compute stages, is transferred to off-chip memory. Note that the CLP repeats this entire process several times to cover all the computations of CONV layer. As evident from Figure 4, double-buffering causes data transfer time to overlap with computation.

B. Multi-CLP Accelerator Design

In the previous section, we discussed an optimized CLP design that processes CONV layers of a CNN architecture iteratively one after another. Given that CNN layers have radically varying dimensions, designing a single CLP to process all layers may be ideal for a particular CONV layer, or some CONV layers, but causes under-utilization of computational resources in others, affecting overall performance. This is because of the mismatch between the CLP dimensions (T_n) and T_m and the CONV layers dimensions (N) and (N).

Considering AlexNet layers illustrated in Figure 1 as a motivating example, we found that, on average, 78% of the computational resources were idle, or doing useless work, while processing the first CONV layer. Note that we followed the methodology in [16] to find the cross-layer optimized T_n and T_m . The utilization of g-th CLP computational resources while it is processing CONV layer ℓ can be quantified by

$$\lambda_{\ell}^{(g)} = \frac{\left(N_{\ell} / T_n^{(g)}\right) \times \left(M_{\ell} / T_m^{(g)}\right)}{\left\lceil N_{\ell} / T_n^{(g)} \right\rceil \times \left\lceil M_{\ell} / T_m^{(g)} \right\rceil} \tag{2}$$

To illustrate the computation resource utilization problem, Figure 5 shows a simple example of a CLP designed with factors $\langle T_n, T_m, T_r, T_c \rangle$ equal to $\langle 2, 3, 2, 2 \rangle$. The CLP

Algorithm 2: Optimized Convolution Algorithm.

```
Input: The input feature maps (IF), weight (W) and bias (B) parameters, number of input feature maps (N) and their tiling
            factor (T_n), number of output feature maps (M) and their tiling factor (T_n), size of each output feature map, rows (R)
            and columns (C), and their corresponding tiling factors, (T_r) and (T_c), kernel size (K), and window stride size (S).
   Output: The output feature maps (OF).
   Procedure CONV (IF, W, B, N, T_n, M, T_m, R, T_r, C, T_c, K, S):
        for r \leftarrow 0 to R-1 by T_r do
            for c \leftarrow 0 to C-1 by T_c do
                 for m \leftarrow 0 to M-1 by T_m do
 3
                      \mathbf{OF\_BUF} \leftarrow \operatorname{Broadcast}(\mathbf{B}[m:m+T_m],T_r,T_c)
                      for n \leftarrow 0 to N-1 by T_n do
                           \mathbf{IF\_BUF} \leftarrow \mathbf{IF}[n:n+T_n][S \times r:S \times (r+T_r-1)+K][S \times c:S \times (c+T_c-1)+K]
                           W_BUF \leftarrow W[m: m + T_m][n: n + T_n][0: K][0: K]
                           for i \leftarrow 0 to K-1 do
 8
                               for j \leftarrow 0 to K-1 do
                                    for rt \leftarrow 0 to \min(T_r - 1, R - r - 1) do
 10
                                         for ct \leftarrow 0 to \min(T_c - 1, C - c - 1) do
 11
                                              for all mt \leftarrow 0 to T_m - 1 in parallel do
                                                                                                                                                 ▶ Unroll
 12
                                                  for all nt \leftarrow 0 to T_n - 1 in parallel do
                                                                                                                                                 ▶ Unroll
 13
                                                       \mathcal{P} \leftarrow \mathbf{IF\_BUF}[nt][S \times rt + i][S \times ct + j] \times \mathbf{W\_BUF}[mt][nt][i][j]
                                                       \mathbf{OF\_BUF}[mt][rt][ct] \leftarrow \mathbf{OF\_BUF}[mt][rt][ct] + \mathcal{P}
 15
                                                  end
 16
                                        end end
 17
 18
                                    end
19
                               end
20
21
                           end
22
                      end
                      \mathbf{OF}[m:m+T_m][r:r+T_r][c:c+T_c] \leftarrow \mathbf{OF\_BUF}
23
24
                 end
25
            end
26
        end
        return OF
27
```

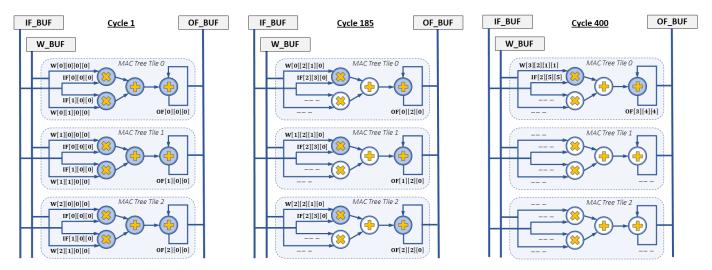


Fig. 5: An example of a Single-CLP design with $\langle T_n, T_m, T_r, T_c \rangle$ factors set to $\langle 2, 3, 2, 2 \rangle$. The CLP processes a CONV layer with $\langle N, M, R, C, K, S \rangle$ configurations equal to $\langle 3, 4, 5, 5, 2, 1 \rangle$.

is used to process a CONV layer whose configurations $\langle N, M, R, C, K, S \rangle$ are $\langle 3, 4, 5, 5, 2, 1 \rangle$. One can notice that N and M are not a perfect multiple of T_n and T_m , respectively. In particular, T_n needs two iterations to cover N. The first iteration covers the first two input FMs, leaving one input FM for the second iteration. The same goes for T_m where one output FM is left for the second iteration.

During the first cycle, the CLP works on complete tiles, therefore all of the T_m MAC tree tiles are in full use. On the other hand, the CLP uses a partially filled tile of input FMs in cycle 185, during which only T_m (out of $T_n \times T_m$) multipliers and adders do useful calculations. To make matters worse, during the last cycle, cycle 400, which computes the output feature OF[3][4][4], neither IF_BUF nor W_BUF is

completely full. Thus, only half of the resources of the first MAC tree tile are used. The combined effect of T_n and T_m dimensional mismatches leads to utilizing only 50% of the computational resources based on Equation (2).

To overcome the issue of under-utilization, and thus improve CNN throughput, the CLP can be designed with reconfigurable unrolling factors to make it work well with different layers. However, such a design requires the construction of complex hardware structures, which leaves the CLP with insufficient resources for computation. Instead, we partition the available hardware resources among multiple small, specialized CLPs. Specifically, the Multi-CLP design adopts G CLPs to process CNN layers, where $G \in [1, L]$ and L is the number of CONV layers. Each CONV layer is bound to a single CLP. In other words, a fixed set of layers, referred to as \mathfrak{L}_g , is assigned to g-th CLP, $\mathfrak{L}_g \subseteq [L]$. Later, each CLP sequentially processes its assigned CONV layers.

In this context, we define an episode as a single pass through CLP layers. To avoid intra-episode data dependencies, the output FMs produced from CONV layer ℓ during the i-th episode are not used in episode i. Instead, they are used as input for CONV layer $(\ell+1)$ in episode (i+1). Therefore, CLPs need to synchronize before starting a new episode. The advantage of Multi-CLP accelerator comes from various dimensions supported that can accommodate CONV layers of varying dimensions. Furthermore, by applying data pipelines to the CLP level as well, it becomes possible to work concurrently on L independent images. Since a single CONV layer is processed for a given input image in an episode, processing an image requires L episodes. Note that this nature of backto-back processing allows data transfer for one layer to be overlapped with computation for another.

The problem now is, how to determine the appropriate number of CLPs to use, on what basis to assign a CONV layer to a CLP, and how many hardware resources each CLP can utilize? Here, we must emphasize that this research aims to optimize CNN throughput. To achieve that, each CONV layer must be assigned to the CLP that most closely matches its dimensions. Therefore, we need to look for a CONV-CLP assignment that maximizes resource utilization. Additionally, CNN throughput is constrained by the CLP that takes the longest time to finish the *episode*. Thus, the *episode* time for that CLP should be as short as possible. Furthermore, the size of on-chip buffers is inversely related to the required off-chip bandwidth. Therefore, we need to ensure that the Multi-CLP design uses sufficient on-chip buffers to minimize duplicate data transfers.

To fulfill the aforementioned requirements, in Section III-C, we present an analytical and empirical model to help estimate the performance and cost, in terms of hardware resources, for a

candidate Multi-CLP design. Then, in Section IV, we introduce an optimization framework that employs the proposed model as well as metaheuristics to explore the intractable design space for the optimal Multi-CLP design, constrained by the given CNN architecture and target FPGA specifications.

C. Design Cost and Performance

The Multi-CLP design for a CNN is characterized by (i) the number of CLPs, (ii) the assignment of CONV layers to CLPs, (iii) the unrolling factors $\langle T_n, T_m \rangle$ of each CLP, and (iv) the tiling parameters $\langle T_r, T_c \rangle$ of each CONV layer. Considering that synthesis, placement, and routing of a candidate Multi-CLP design may take several minutes or even hours, one can note that it is infeasible to perform these processes at each design point for selecting the candidate with the highest performance. Hence, this section presents an analytical and empirical design scheme for accurate modeling of cost, in terms of hardware resources used, and performance for a Multi-CLP design given its parameters. Here, the focus is on the key specifications and hardware resources of the target FPGA platform, including DSPs, block random access memories (BRAMs), and off-chip memory bandwidth (BW).

1) Computational Performance: To evaluate the efficiency of a Multi-CLP design, the speed of each CLP, in terms of computation cycles, must be considered. Based on Algorithm 2, the number of cycles needed to process CONV layer ℓ assigned to q-th CLP is modeled as follows. Along the rows and columns dimensions of the output FMs, the algorithm iterates over the tiles and the elements in each tile. Thus, it runs for $R_\ell \times C_\ell$ cycles to process loop iterators r, c, rt, and ct. On the other hand, input FMs and output FMs are unrolled using the factors $T_n^{(g)}$ and $T_m^{(g)}$, respectively. Therefore, processing loop iterators n, m, nt, and mt need $\lceil N_{\ell} / T_n^{(g)} \rceil \times \lceil M_{\ell} / T_m^{(g)} \rceil$ cycles. Note that the two innermost loops, highlighted in yellow, only take 1 cycle because they are unrolled. Finally, loop iterators i and j are neither tiled nor unrolled, and thus, require $K_{\ell} \times K_{\ell}$ cycles. Altogether, the number of cycles required to process CONV layer ℓ on g-th CLP is calculated as

$$Comp_Cycle_{\ell}^{(g)} = \left\lceil \frac{N_{\ell}}{T_n^{(g)}} \right\rceil \times \left\lceil \frac{M_{\ell}}{T_m^{(g)}} \right\rceil \times R_{\ell} \times C_{\ell} \times K_{\ell}^2$$
 (3)

Accordingly, the computational performance of the g-th CLP design in processing CONV layer ℓ is defined as the total number of operations, i.e., multiply and accumulate operations, required by the CLP to process the layer over the total number of cycles needed to do so as follows

$$Comp_Perf_{\ell}^{(g)} = \frac{2 \times N_{\ell} \times M_{\ell} \times R_{\ell} \times C_{\ell} \times K_{\ell}^{2}}{Comp_Cycle_{\ell}^{(g)}}$$

$$= \frac{2 \times N_{\ell} \times M_{\ell}}{\left\lceil N_{\ell} / T_{n}^{(g)} \right\rceil \times \left\lceil M_{\ell} / T_{m}^{(g)} \right\rceil}$$
(4)

2) DSP Slice Usage: The use of DSP slices in each CLP is dominated by the T_m MAC tree tiles that work in parallel to improve computational throughput. Each MAC tree tile consists of T_n parallel multipliers and an adder tree. Specifically, each multiplier performs multiplication on an input activation feature and an input kernel weight. Then, the resultant products are accumulated with the old partial result using a binary adder tree consisting of T_n adders. Thus, CLP computational engine contains $T_n \times T_m$ multipliers and adders.

It is noteworthy that the number of CLP slices required depends on the operation type and data representation. To estimate how many DSP slices are needed for a given CLP design, we empirically measure the number of DSPs used to implement the multiplier and adder for each representation in the design space and use a lookup table model to calculate the total usage. Based upon this, the number of DSP slices required to implement the computational engine of *g*-th CLP is determined by

$$DSP_Usage^{(g)} = DSP(Q_{IF}, Q_W) \times T_n^{(g)} \times T_m^{(g)}$$
 (5)

where \mathcal{Q}_{IF} and \mathcal{Q}_{W} indicate whether input features and weights, respectively, are in single-precision floating-point representation (FP32) or quantized to low-precision representation such as 16-bit fixed-point format (FxP16), and the function $\mathrm{DSP}(\mathcal{Q}_{IF},\mathcal{Q}_{W})$ returns the number of DSP slices required to implement a single multiplier and a single adder using the specified format. For instance, the FP32 adder and multiplier comprise 2 and 3 DSP slices, respectively [29]. That is, $\mathrm{DSP}(\mathrm{FP32},\mathrm{FP32})$ equals 5. On the other hand, the FxP16 adder and multiplier can be implemented on the same DSP slice, and therefore $\mathrm{DSP}(\mathrm{FxP16},\mathrm{FxP16})$ equals 1. For a Multi-CLP design consisting of G CLPs, the number of DSP slices required is the sum of the DSPs used by all CLPs

$$DSP_Usage = \sum_{g=0}^{G-1} DSP_Usage^{(g)}$$
 (6)

3) BRAM Usage: Each CLP in Multi-CLP design requires three on-chip buffers, namely, IF_BUF, W_BUF, and OF_BUF as discussed in Section III-A. These buffers are typically constructed using BRAMs. To accurately model BRAM usage in each CLP, the following must be considered; (i) BRAM capacity and configuration capability, (ii) data precision, (iii) buffer scheme (single/double buffering), (iv) number of banks needed per buffer, (v) number of read/write ports required for each bank, (vi) CLP unrolling factors (T_n, T_m) , and (vii) tiling parameters (T_r, T_c) and configurations (K, S) for each CONV layer assigned to the CLP.

In this work, we employ the double-buffering scheme for all buffers to overlap off-chip data transfer with computation. Furthermore, we adopt the memory banking method to allow each buffer to provide the required number of input/output channels. Accordingly, the bank size must be large enough to support the two most demanding successive tiles. If considering the g-th CLP that processes CONV layers in subset \mathfrak{L}_g , the above requirement on minimum bank depth is translated as

$$\operatorname{Min_Depth}_{d}^{(g)} = \max \Big(\operatorname{MFP}(d, \ell) + \operatorname{MFP}(d, \ell + 1), \\ 2 \times \operatorname{MFP}(d, \ell) : \forall \ell \in \mathfrak{L}_g \Big)$$
 (7)

where the function MFP (d,ℓ) returns the memory footprint, in terms of data size, for layer ℓ tiles on each bank in the buffer type $d, d \in \{IF, W, OF\}$, and \mathfrak{L}_g is a set of all CONV layers assigned to g-th CLP, $\mathfrak{L}_g \subseteq [L]$. The number of BRAMs per bank (BPB) is calculated by

$$BPB_d^{(g)} = \left\lceil \frac{Min_Depth_d^{(g)}}{ADDR(Q_d)} \right\rceil$$
 (8)

Here, the ADDR(\mathcal{Q}_d) is a function that returns the appropriate BRAM depth, i.e., the number of addresses, to store data represented with the width specified by \mathcal{Q}_d . The BRAM model considered in this paper is based on the RAMB18E1 primitive provided by Xilinx 7 series FPGAs [30]. This memory can be configured ($depth \times width$) as a 512 × 36, 1024×18 , 2048×9 , 4096×4 , 8192×2 , or 16384×1 . Thus, ADDR(FP32) is equal to 512, while ADDR(FxP16) is equal to 1024. Thus, the total number of BRAMs required to construct the buffer type d of g-th CLP is calculated by

$$BRAM_Usage_d^{(g)} = BPB_d^{(g)} \times \mathcal{B}_d^{(g)}$$
 (9)

where $\mathcal{B}_d^{(g)}$ represents the number of memory banks needed to design the buffer type d of g-th CLP. Details of the specification of BRAMs and the number of banks needed for each buffer type are presented next.

IF_BUF: The CLP computational engine requires T_n channels from IF_BUF. Therefore, this buffer is designed using T_n banks, i.e., $\mathcal{B}_{IF}^{(g)} = T_n^{(g)}$. When processing CONV layer ℓ , its footprint on each of these banks is

$$MFP(IF,\ell) = (K_{\ell} + S_{\ell} \times (T_{r\ell} - 1)) \times (K_{\ell} + S_{\ell} \times (T_{c\ell} - 1))$$
(10)

Note that each of IF_BUF BRAMs must provide two independent ports; a read-only port for supplying CLP computational engine and a write-only port for loading off-chip memory data. Therefore, the type of these BRAMs is set to simple dual-port mode [30].

W_BUF: The weights buffer is similar to input FMs buffer in terms of operating mode. But, unlike IF_BUF which is shared among all MAC tree tiles, the digital multipliers of MAC tree tiles use different kernel weights, and thus require independent memory banks. Accordingly, the W_BUF is organized into

 $T_n \times T_m$ banks, i.e., $\mathcal{B}_W^{(g)} = T_n^{(g)} \times T_m^{(g)}$. The memory footprint of CONV layer ℓ on each of these banks is

$$MFP(W, \ell) = K_{\ell} \times K_{\ell} \tag{11}$$

OF_BUF: In each computation cycle, the old output feature is loaded and accumulated with the result of MAC tree tile. After that, the new output feature is stored back in the output FMs buffer. As we note, the accumulator needs a read port and a write port. However, these ports use the same address for both operations. Additionally, the resulting output features from the previous stage must be transferred to off-chip memory. Therefore, the OF_BUF uses true dual-port BRAMs in write-first mode [30]. One can also note that OF_BUF must be organized into T_m banks. That is, $\mathcal{B}_{OF}^{(g)}$ equals $T_m^{(g)}$. Thus, the memory footprint of CONV layer ℓ on each of these banks is calculated as

$$MFP(OF, \ell) = T_{r\ell} \times T_{c\ell}$$
 (12)

Together, the total number of BRAM resources required to construct the on-chip buffers of the Multi-CLP design is calculated as

BRAM_Usage =
$$\sum_{g=0}^{G-1} \sum_{d \in \{IF, W, OF\}} BRAM_Usage_d^{(g)}$$
 (13)

4) Bandwidth Usage: The computational performance discussed in Section III-C1 represents the maximum number of operations that can be performed with the available computational resources per cycle. Thus, it is also known as the computational roof. However, one critical problem is that CLP computation may be blocked by data transfer for memory-bound layers. Consequently, such layers cannot achieve the best performance supported by computational resources. To find the attainable computational performance in such a case, we adopt the roofline model [28].

The roofline model is used to provide performance estimates for a candidate implementation on a particular system. It relates the attainable performance to the peak performance the system can achieve and the off-chip DRAM memory traffic as illustrated in Figure 6. The input/output bandwidth roof determines the maximum computational performance supported by the memory system for a given computation to communication (CTC) ratio, which is calculated as $CTC \times BW$. Thus, the attainable performance for CONV layer ℓ when it is processed by g-th CLP is given by

$$\operatorname{Perf}_{\ell}^{(g)} = \min \left(\operatorname{Comp_Perf}_{\ell}^{(g)}, \ CTC_{\ell}^{(g)} \times BW \right)$$
 (14)

where $\operatorname{Comp_Perf}_{\ell}^{(g)}$ is the computational roof indicated in Equation (4), the BW refers to the off-chip memory bandwidth, and the $CTC_{\ell}^{(g)}$ describes the off-chip memory traffic that the CLP needs during CONV layer processing. The CTC can be

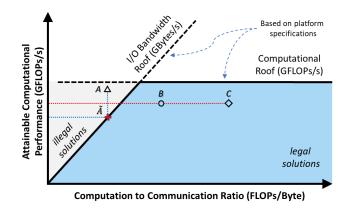


Fig. 6: An illustration of roofline model basis. As we can see, implementation A is memory-bounded, i.e., system throughput is constrained by the communication with off-chip memory, while implementations B and C are computation-bounded, and thus system throughput can only be improved by using additional computational resources. One can also note that design C has better data reuse compared to design B, and thus it is more preferable than design B.

calculated by dividing the total number of operations by the total amount of off-chip data accessed as follows

$$CTC_{\ell}^{(g)} = \frac{2 \times N_{\ell} \times M_{\ell} \times R_{\ell} \times C_{\ell} \times K_{\ell}^{2}}{\sum_{d \in \{IF, W, OF\}} \alpha_{d, \ell}^{(g)} \times \mathcal{B}_{d}^{(g)} \times MFP(d, \ell)}$$
(15)

Here, $\mathcal{B}_d^{(g)}$ denotes the number of memory banks in the buffer type d of g-th CLP, MFP (d,ℓ) refers to the memory footprint for CONV layer ℓ tiles on these memory banks as mentioned in Equations (10), (11), and (12), and $\alpha_{d,\ell}^{(g)}$ indicates the number of off-chip memory accesses needed for buffer type d during the computation of CONV layer ℓ by g-th CLP, such that

$$\alpha_{IF,\ell}^{(g)} = \alpha_{W,\ell}^{(g)} = \frac{N_{\ell}}{T_n^{(g)}} \times \frac{M_{\ell}}{T_m^{(g)}} \times \frac{R_{\ell}}{T_{r\ell}} \times \frac{C_{\ell}}{T_{c\ell}}$$

$$\alpha_{OF,\ell}^{(g)} = \frac{M_{\ell}}{T_m^{(g)}} \times \frac{R_{\ell}}{T_{r\ell}} \times \frac{C_{\ell}}{T_{c\ell}}$$
(16)

To estimate the bandwidth required to support the maximum computation speed, we focus on the peak bandwidth that a CLP needs. Specifically, the minimum bandwidth required to achieve the peak computational performance for CONV layer ℓ which is processed by q-th CLP is determined by

$$\operatorname{Min_BW}_{\ell}^{(g)} = \frac{\operatorname{Comp_Perf}_{\ell}^{(g)}}{CTC_{\ell}^{(g)}}$$
 (17)

For a memory bound layer, the speed of g-th CLP design in processing CONV layer ℓ is defined by the data transfer cycles instead of the computation cycles mentioned in Equation (3). Thus, the arithmetic utilization of Multi-CLP design processes

a CNN with L CONV layers using G CLPs is given by

$$\text{Utilization} = \frac{\sum_{g \in [G]} \sum_{\ell \in \mathfrak{L}_g} \lambda_{\ell}^{(g)} \times \text{Cycle}_{\ell}^{(g)}}{G \times \max\left(\text{Cycle}^{(g)} : \forall g \in [G]\right)}$$
(18)

where $\lambda_{\ell}^{(g)}$ is the utilization ratio of g-th CLP computational resources while it is processing CONV layer ℓ that is computed as mentioned in Equation (2).

IV. DESIGN SPACE EXPLORATION

In this section, we discuss the difficulty of determining the value of design variables for multiple CLPs used to accelerate CNNs on FPGAs. Then, we present our proposed approach to finding the optimal Multi-CLP design for a given CNN architecture and resource budget.

A. Problem Statement and Formulation

The implementation of a CNN architecture on an FPGA platform using multiple CLPs while optimizing their performance is a non-trivial task. This is because it requires optimizing the utilization of both processing units and memory units to minimize the execution cycles. However, the number of cycles required depends on the number of CLPs used, the assignment of CONV layers to CLPs, the loop unrolling factors $\langle T_n, T_m \rangle$ for each CLP, and the tilling parameters $\langle T_r, T_c \rangle$ for each CONV layer, which imposes an intractable design space. Furthermore, we must ensure that the peak memory bandwidth and resources utilized for a given candidate design comply with the given hardware constraints.

Thus, determining the best values of the design variables through exhaustive search is infeasible, especially when the number of variables is large and/or FPGA resources are limited. Additionally, the number of possible assignments for CONV layers to CLPs is exponential in the number of layers. Considering the unrolling factors $\langle T_n, T_m \rangle$ as a motivating example, the search space of a given CNN structure with L CONV layers would consist of $L^L \times \max(N_\ell : \forall \ell \in [L])^L \times \max(M_\ell : \forall \ell \in [L])^L$ possible solutions, where N_ℓ and M_ℓ are the number of input FMs and output FMs for CONV layer ℓ , respectively.

For such a computationally hard problem, there is no existing deterministic algorithm that can find the optimal solution. Adhoc heuristic solutions have been employed by Shen et al. [18] to find good solutions. However, iterative non-deterministic heuristics can be employed to efficiently traverse the search space as they have been proven to find excellent solutions to such hard problems [31]. In this work, we employ simulated annealing (SA) and tabu search (TS) algorithms to find the optimal Multi-CLP design for a given CNN structure. The

TABLE I: An example of assigning 5 CONV layers to 5 CLPs with initial random values to T_n , T_m , T_r , and T_c .

Config	T_n	T_m	L_1	L_2	L_3	L_4	L_5	# DSPs
CLP_1	3	5	0	0	0	1	1	75
CLP_2	4	1	0	1	0	0	0	20
CLP_3	1	1	1	0	1	0	0	5
CLP_4	7	2	0	0	0	0	0	-
CLP_5	3	6	0	0	0	0	0	-
T_r			3	5	2	4	1	
	T_c		3	3	3	4	2	

SA and TS are adopted due to their efficiency in finding high-quality solutions to such kinds of problems, simplicity of implementation, ability to provide a globally optimal solution in a reasonable time, the ease of tuning their parameters, and providing a balance between the exploration of solution space and the exploration of information obtained [32].

The space of solutions comprises finding optimal values for the following four elements; (i) the number of CLPs to be used, (ii) the assignment of CONV layers to CLPs, (iii) the values of T_n and T_m for each CLP, (iv) the values of T_r and T_c for each CONV layer. The solution representation is composed of L+2 rows as shown in Table I. Each one of the first L rows represents a single CLP, while the last two rows represent the T_r and T_c values of each CONV layer. On the other hand, CLP rows contain L+2 columns. The first two columns store the current values of T_n and T_m , while the remaining L columns indicate the assignment of layers to CLPs by storing 0 (unassigned) or 1 (assigned). Here, we must emphasize that a column representing a layer assignment must have a single 1, while a row can have multiple 1s.

Initially, the search begins with a random assignment of CONV layers to CLPs, and random T_n and T_m values for each CLP provided they do not violate the hardware resource constraints. The candidate solution represented in Table I shows that the total number of DSPs is 100, which can be calculated based on Equation (6) assuming features and weights are in 32-bit floating-point representation. One can notice that CLP_4 and CLP_5 do not require any DSP since no layer is assigned to them. Note that in the case where more than one layer is assigned to a CLP, the CLP processes these layers sequentially one after another. Thus, the number of cycles required by a CLP is the sum of the number of cycles needed to process all its layers.

The proposed optimization algorithm explores the design space by generating neighboring solutions from the current one, aiming to reach the optimal candidate. In our formulation, a neighboring solution is generated by performing a single move from the current solution. A move or step is defined as a single change in the current solution either by changing the layer assignment to a different CLP, or by changing one of the

 T_n , T_m , T_r , or T_c parameters for a CLP. An optimal solution is defined as a Multi-CLP design that requires the least number of cycles to process all layers. The number of cycles required for the candidate design S is calculated as

$$\operatorname{Cycle}(\mathcal{S}) = \max \left(\sum_{\ell \in \mathfrak{L}_q} \operatorname{Cycle}_{\ell}^{(g)} : \forall g \in [G] \right)$$
 (19)

where $\operatorname{Cycle}_{\ell}^{(g)}$ is the number of cycles required to process CONV layer ℓ in g-th CLP, G is the number of CLPs in the candidate solution \mathcal{S} , and \mathfrak{L}_g is a set of all CONV layers assigned to g-th CLP, $\mathfrak{L}_g \subseteq [L]$. Note that the number of cycles for memory-bounded layers is defined by the data transfer cycles, while the number of cycles for computational-bounded layers is defined by the computational roof as discussed in Section III-C. Based on the discussion above, the optimization problem can be summarized as

minimize
$$\operatorname{Cycle}(\mathcal{S})$$

subject to $\operatorname{DSP_Usage}_{\mathcal{S}} \leq \operatorname{DSP}_{max}$ (20)
 $\operatorname{BRAM_Usage}_{\mathcal{S}} \leq \operatorname{BRAM}_{max}$

Here, the DSP_Usage_S and BRAM_Usage_S are the total number of DSP and BRAM resources required for the solution S, which are computed as discussed in Equations (6) and (13), respectively, the DSP_{max} and BRAM_{max} are the number of DSPs and BRAMs available in the target FPGA platform. For each candidate solution, the memory optimization to find the values of T_r and T_c for each CONV layer is based on minimizing the peak memory bandwidth required as discussed in Section III-C4. The parameters T_r and T_c in turn would set the buffer sizes for each CLP, and thus they must comply with the available on-chip memory budget as demonstrated in Section III-C3.

B. Optimization of Multi-CLP Designs

In the section, we describe the employment of simulated annealing and tabu search metaheuristic algorithms in optimizing the implementation of FPGA-based CNN accelerators based on multiple CLPs.

1) Simulated Annealing (SA): The SA is an effective optimization tool based on inspiration derived from annealing of metals [33], [34]. Its power lies in its simplicity of implementation, and the ease of tuning its parameters. During search, the moves that cause a decrease in the cost function, also known as good moves, are always accepted. On the other hand, the moves that lead to higher cost, also known as bad moves or uphill moves, are sometimes accepted with a probability that depends on a parameter called temperature (T). Initially, that is when T is high, the search is almost random. As T decreases, the probability of accepting bad moves decreases

Algorithm 3: Simulated Annealing Algorithm.

```
Input: The initial solution (S_0), initial temperature (T_0),
             cooling rate (\alpha), annealing temperature constant
             (\beta), total allowed time (MaxTime), and time until
             next update (M).
    Output: The best admissible solution (S^*).
    Procedure SA (S_0, T_0, \alpha, \beta, MaxTime, M):
         T \leftarrow T_0
         Time \leftarrow 0
         repeat
              \mathcal{S} \leftarrow \texttt{Metropolis}(\mathcal{S}, T, M)
              Time \leftarrow Time + M
              T \leftarrow \alpha \times T
              M \leftarrow \beta \times M
         until Time \ge MaxTime
         return S^*
   Procedure Metropolis (\mathcal{S}, T, M):
12
         repeat
              S_{new} \leftarrow \text{Neighbor}(S)
13
                                                             Do a single move
              \Delta h \leftarrow \text{Cost}(\mathcal{S}_{new}) - \text{Cost}(\mathcal{S})

    □ Using Eq. (19)

14
              if (\Delta h < 0) or (random < e^{-\Delta h/T}) then
15
                                                           ▷ Accept the solution
16
17
              end
              M \leftarrow M - 1
18
         until M=0
19
         return S
20
```

and the search begins to become greedy. At zero temperature, the search becomes totally greedy, and only good moves are accepted.

The basic structure of SA algorithm is shown in Algorithm 3, the core of which is the Metropolis procedure. The Metropolis procedure simulates the annealing process at a given temperature value T [35]. It receives the current temperature T and the current solution S as input and improves it. It is also provided with the value M, which is a multiplication factor that increases the amount of time annealing is applied at a given temperature T [36]. The temperature is initialized to a value T_0 (explained below), and is slowly decreased using the parameter α . As temperature is lowered, the amount of time spent in annealing at a given temperature is gradually increased, using the parameter β , $\beta > 1$. The time spent in each call to Metropolis procedure is stored in variable Time. SA algorithm stops when the time limit is reached, i.e., $Time \geq MaxTime$.

The Metropolis procedure generates a local new solution \mathcal{S}_{new} from a given solution \mathcal{S} using the procedure Neighbor. Two perturbation schemes are used to generate a new neighbor out of the current solution; either (i) by randomly changing the assignment of a single CONV layer to CLPs, or (ii) by mutating T_n or T_m of a randomly selected CLP. These two perturbations are performed probabilistically, by making changes to T_n or T_m with 80% probability while making changes of CONV layer assignment to CLPs with 20% probability. The rationale for

this is that because of the observation that optimizing the values of T_n and T_m for CLPs has greater effect on the solution cost. Additionally, during later iterations, changing CONV layer assignment to CLPs results in more invalid solutions. This is because adding new CLPs would require additional DSPs which may exceed the DSP constraint.

The function $\operatorname{Cost}(\mathcal{S})$ returns the cost of a given solution \mathcal{S} , which represents the number of cycles required to complete the execution for the given solution representation as mentioned in Equation (19). The new solution is accepted if its cost, $\operatorname{Cost}(\mathcal{S}_{new})$, is better than the cost of the current solution \mathcal{S} . The new solution is also accepted by $\operatorname{Metropolis}$ probabilistically when its cost is higher than the cost of the current solution. This is done by generating a random number in the range 0 to 1. If the generated number is smaller than $(e^{-\Delta h/T})$, then the inferior solution is accepted, resulting in hill-climbing, i.e., transition from a low-cost solution to a high-cost solution, where e is the Euler's number, and Δh is the change in cost.

2) Tabu Search (TS): The TS is an another non-deterministic algorithm that also starts from an initial feasible solution and generates a list of candidate solutions, that forms what is known as the candidate list (V*), by making neighborhood moves, and then selecting from this list the solution that is best among all candidates in the current iteration [37]–[40]. To avoid move reversals, that is recycling back to previously visited solutions, a device called tabu list (TL) is used that stores some attribute(s) of these moves. The TL has a given size and can be viewed as a window or queue on accepted moves. The moves stored in TL are not allowed as they may undo previous moves returning back to the same solution.

The tabu status is overridden when certain criteria, known as aspiration criteria (AC), are satisfied. The AC temporarily override the tabu status if the move causes the cost to be less than the aspiration level (AL), a value defined based on the AC adopted. For instance, the "best cost" criterion defines the AL as the cost of the best admissible solution. Accordingly, if the tabu move leads to a solution whose cost is less than AL, which is an indication that there is no cycling back to a previously visited solution, then the AC override the tabu status and accept the move. Additionally, the AL is updated to the cost of the new solution.

An algorithmic description of TS metaheuristic is shown in Algorithm 4. The procedure starts from a feasible solution, a neighborhood $\aleph(\mathcal{S})$ is defined for each solution \mathcal{S} . Among all the possible neighborhood solutions, N neighbor solutions are randomly selected and added to the candidate list using the procedure $\mathtt{Sample}(\aleph(\mathcal{S}),N)$. Note that the number of possible neighborhood solutions is typically much larger than the candidate list size, i.e., $|\mathbf{V}^*| = N \ll |\aleph(\mathcal{S})|$. From these N

Algorithm 4: Tabu Search Algorithm.

```
Input: A set of feasible solutions (\Omega), candidate list size
             (N), tabu list size (M), aspiration criteria (AC),
             total allowed time (MaxTime).
    Output: The best admissible solution (S^*).
    Procedure TS (\Omega, N, M, AC, MaxTime):
         Start with an initial feasible solution S, S \in \Omega
         Initialize best admissible solution and current time
 2
           S^* \leftarrow S
            Time \leftarrow 0
         Initialize tabu list and aspiration level
            \mathbf{TL} \leftarrow \{\}
            \mathbf{AL} \leftarrow \mathsf{Cost}(\mathcal{S}^*)
                                                   ▷ Assuming AC is "best cost"
         repeat
              Generate neighborhood solutions of S, \aleph(S)
              Select N neighbor solutions and add them to
                candidate list
                    \mathbf{V}^* \leftarrow \texttt{Sample}(\aleph(\mathcal{S}), N)
              Choose best solution in candidate list
                    \mathcal{S}^* \leftarrow \text{Best}(\mathbf{V}^*)
              if Move(S, S^*) \notin TL or Cost(S^*) < AL then
                   Update tabu list
                      \mathbf{TL} \leftarrow \mathbf{TL}[1:M] + \{ \texttt{Move}(\mathcal{S}^*, \mathcal{S}) \}
10
                   S \leftarrow S^*
                                                           if Cost(S) < AL then
11
                      \mathbf{AL} \leftarrow \mathtt{Cost}(\mathcal{S})
12
                                                        ▶ Update aspiration level
13
                   Time \leftarrow Time + 1
14
              end
15
16
         until Time = MaxTime
         return S^*
17
```

neighborhood solutions in the candidate list, the best solution is chosen for consideration as the next solution \mathcal{S}^* . The selected solution could have a cost higher than the cost of the current solution resulting in hill-climbing.

Two perturbation schemes are used to generate neighborhood solutions; either (i) by randomly changing the assignment of a single CONV layer to CLPs, or (ii) by mutating T_n or T_m of a randomly selected CLP. For each of the perturbation schemes, a separate TL is maintained. The attributes stored in TLs are related to the perturbation scheme. In the case where a perturbation consists of changing the assignment of a CONV layer to CLPs, the CONV layer ID and the CLP ID to which it is assigned are stored. For the case of changing parameter T_n or T_m , the attributes saved are the CLP ID and the new value of the parameter. These two perturbations are done with the same probabilities as in the case of SA, i.e., 80% for changing T_n or T_m , and 20% for changing CONV layer assignment to CLPs.

V. EXPERIMENTS AND RESULTS

In this section, we describe our experimental settings. Then, we employ the proposed optimization framework to design multiple CLPs, Multi-CLP in short, to accelerate four widely used typical CNNs on FPGAs. Next, we discuss the results and findings. Additionally, we compare the Multi-CLP accelerator

designed using the presented metaheuristic-based framework with the state-of-the-art Single-CLP and Multi-CLP design frameworks.

A. Experimental Settings

To demonstrate the versatility of metaheuristics in designing Multi-CLP accelerators for CNNs, we conduct extensive experiments on the ImageNet large-scale visual recognition challenge (ILSVRC) dataset [41], which is known as one of the most popular image classification benchmarks. The ILSVRC dataset contains 50 thousand images for validation, all of which are natural high-resolution images. Each image is annotated as one of 1,000 classes. To meet the configurations required for input data to CNNs, images are resized and then center cropped before being used as input to the network.

For those experiments, we use our metaheuristic-based optimization technique discussed in Section IV to design Multi-CLP accelerators for four benchmark architectures, namely, AlexNet [7], SqueezeNet 1.1 [22], VGGNet [23], and GoogLeNet [24], as they are well-known in the field of image classification. The number of CONV layers in AlexNet, VGGNet, SqueezeNet 1.1, and GoogLeNet architectures is 10, 13, 26, and 57, respectively. To investigate the applicability of the proposed framework in designing Multi-CLP accelerators for CNNs with different bit-precision levels, we consider two formats; 32-bit floating-point and 16-bit fixed-point.

The Multi-CLP design is implemented as a parameterized Verilog hardware description language (Verilog-HDL). The design is synthesized and implemented using Xilinx Vivado Design Suite (v2020). Our application targets the VC707 and VC709 FPGA boards, both boards operate at a frequency of 100 MHz but each has different hardware resources. Table II summarizes the hardware specifications of these two boards in terms of FPGA platform, DSPs, 18 Kb BRAM (BRAM18K), LUTs, FFs, and external memory.

The optimization problem described in Equation (20) employs the computational performance and resource utilization models discussed in Section III-C to find the highest-throughput Multi-CLP design. The maximum number of DSPs (DSP $_{max}$) and BRAMs (BRAM $_{max}$) available for the optimization process is set to 80% of the board's resources as in [18]. This is due to the need for some hardware resources to implement, for example, the soft-core processor and controllers. Accordingly, the DSP $_{max}$ and BRAM $_{max}$ are set to 2,240 and 1,648 for the VC707 board, respectively, while maximum of 2,880 DSPs and 2,352 BRAMs are allowed for the optimization procedure when working on the VC709 board.

With regards to the optimization procedure, the SA metaheuristic algorithm needs to start from a high temperature T, which is set in our experiments to 25,000. The parameter α for

TABLE II: FPGA boards hardware resources.

Specification	VC707 [42]	VC709 [43]		
Platform	Virtex-7 VX485T	Virtex-7 VX690T		
DSP Slices	2,800	3,600		
BRAM18K	2,060	2,940		
LUTs	303,600	433, 200		
FFs	607, 200	866, 400		
External Memory	1GB DDR3 SODIMM	2 × 4GB DDR3 SODIMM		

updating the temperature is set to the value 0.99. On the other hand, the constant β which controls the time spent at each annealing temperature is set to 1.005. The stopping criterion is set to 1,000 iterations. For TS metaheuristic algorithm, the size of the candidate list used is 20, as it has been found based on experimental analysis that increasing this number does not make significant improvements while decreasing it has a negative effect on the quality of the obtained solutions. Based on the empirical analysis, the size of the tabu list chosen was 7, as this gave the best results, while the stopping criterion is set to 1,000 iterations. For each performed experiment, we have reported the best results obtained from 10 runs.

B. Results and Discussion

In the first experiment, we study the effectiveness of the proposed optimization framework in finding the optimal configurations that maximize the throughput of the Multi-CLP design. Hence, this experiment uses the computational performance and resource usage models as well as the maximum resources available on the FPGA to find the solution with the minimum execution cycles.

The behavior of a single run of SA algorithm for AlexNet CNN acceleration using multiple CLPs on 485T FPGA is shown in Figure 7. The figure shows the variation in the execution time, which is defined as the total number of cycles divided by the operating frequency. The number of cycles is the cost function, illustrated in Equation (19), that SA algorithm tries to optimize for 1,000 iterations when run with the parameters mentioned in the previous section. We refer to this approach as SA-based Multi-CLP. As can be seen, since the algorithm starts with a high temperature, the fluctuation in costs is large, which gets reduced over time as the temperature is reduced and better solutions are discovered with minimal cost variations toward the end. Final solution configurations can be found in Table III.

Similarly, the behavior of a single run of TS algorithm for the same CNN and the same FPGA board is also illustrated. This approach to designing a Multi-CLP based on TS algorithm

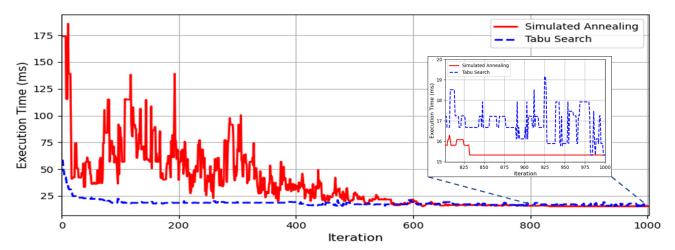


Fig. 7: Cost function optimization progress for a run of 1,000 iterations for AlexNet on 485T FPGA using simulated annealing and tabu search metaheuristic algorithms.

is indicated by the TS-based Multi-CLP. As can be seen from Figure 7, the cost reduces very quickly in the first few iterations as there are several neighboring candidates to choose from in each iteration, and then, it reduces at a slower rate in subsequent iterations with hill-climbing behavior, until reaching the final solution shown in Table III.

From the figure, we can deduce the effectiveness of these algorithms in finding the optimal solution to the Multi-CLP design problem. Both algorithms were able to efficiently explore the intractable design space and find Multi-CLP design solutions with approximately the same low level of execution time, but each has a different final solution. It is noteworthy that the optimization framework is very fast, it explores the search space in less than a minute on a general-purpose processor.

Next, we use the SA-based Multi-CLP and TS-based Multi-CLP approaches to accelerate AlexNet architecture on 485T FPGA. The obtained designs are then compared with the state-of-the-art Single-CLP [16] design and Multi-CLP [18] design, as shown in Table III. The table shows the number of CLPs adopted by each technique and the parallelism factors for each CLP. Additionally, it shows the assignment of CONV layers to CLPs. The number of cycles each CLP spends in processing its layers is also provided. Note that when layer assignment column contains more than one layer in a row, the cycle count provided is the total cycles needed to process all of those layers.

The table also shows the execution time of an image's CONV layer, which represents the time interval in which the system can receive a new image for processing as discussed in Section III-B. The execution time is calculated by dividing design cycles count, illustrated in Equation (19), over design operating frequency. Specifically, the cycle count for each CLP is the number of cycles required to execute all of its assigned

layers. On the other hand, CLPs in Multi-CLP design work concurrently. Therefore, the overall cycle count for such a design is the maximum cycle count for its CLPs. The last two columns in Table III present the number of DSP slices and BRAMs utilized for each technique.

As evidenced by the results, the number of cycles required in the case of a single CLP implementation is almost $1.3 \times$ more than that of the multiple CLPs implementations. The results also show that SA-based Multi-CLP and TS-based Multi-CLP were able to find implementation configurations with a fewer number of cycles (less execution time) compared to Multi-CLP even though all designs use the same number of CLPs and DSP slices. This demonstrates the effectiveness of the proposed SA-based Multi-CLP and TS-based Multi-CLP in balancing the computational resources and workloads assigned to each CLPs such that each CLP can be kept busy most of the time. Although the implementation configurations obtained using SA and TS metaheuristic algorithms show small improvements, the results indicate the importance of using metaheuristics in finding different configurations for a given CNN using a systematic methodology.

Table IV shows similar experiments performed for AlexNet on 690T FPGA. The number of cycles required in the case of a single CLP implementation is about $1.5\times$ more than that of multiple CLPs implementations. For the implementation configuration obtained using SA-based Multi-CLP, we got the same number of cycles as obtained by Multi-CLP. Although both techniques result in designs with same number of cycles and CLPs, these designs have completely different configurations for their adopted CLPs. On the other hand, implementation configuration obtained by TS-based Multi-CLP has used only 4 CLPs and results in a higher execution time.

In Table V, we summarize the results presented in Ta-

TABLE III: Single-CLP and Multi-CLP accelerators for AlexNet on 485T FPGA using 32-bit floating-point data.

Technique	CLP No.	Unrolling	g Factors	Layer	Cycles	Execution	DSPs	BRAMs
rechnique	CLF NO.	T_n	T_m	Mapping	Cycles	Time (ms)	DSFS	DKANIS
Single-CLP [16]	CLP_1	7	64	$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c} -732\times10^{3}\\ -510\times10^{3}\\ -510\times10^{3}$	20.06	2,240	618
	GL D	9	0.4	5a , 5b	170×10^{3}			
	$-CLP_1$	3	24	$\begin{bmatrix} 1a & 1b \\ -a & -a \end{bmatrix}$	$1,464 \times 10^3$	-		
Multi-CLP	$-CLP_2$	8	19	$\begin{bmatrix} 2a & 2b \\ -2 & -2 \end{bmatrix}$	$1,531 \times 10^3$	15.58	2 240	731
[18]	$-CLP_3$	1	96	$\frac{3a}{1} - \frac{3b}{1} - \frac{3b}{1}$	$1,558 \times 10^{3}$	15.56	2,240	191
	CLP_4	2	64	$\begin{vmatrix} 4a & 4b \\ 5a & 5b \end{vmatrix}$	$\begin{bmatrix} -7.5 & -7.5 $			
	CLP_1	3	24	$-\frac{1a}{4a}$	$-\frac{732 \times 10^3}{779 \times 10^3} -$		2,240	
SA-based Multi-CLP	CLP_2	CLP_2 3	24	$\begin{bmatrix} -1b \\ -4b \end{bmatrix}$	$\begin{bmatrix} -732 \times 10^{3} \\ -779 \times 10^{3} \end{bmatrix} - \begin{bmatrix} -732 \times 10^{3} \\ -779 \times 10^{3} \end{bmatrix} - \begin{bmatrix} -732 \times 10^{3} \\ -779 \times 10^{3} \end{bmatrix}$	15.31		644
(ours)	CLP_3	16	11	$\begin{bmatrix} 2a, 2b \\ -5a \end{bmatrix}$	$\begin{array}{c} 1,312\times 10^{3} \\ -219\times 10^{3} \end{array}$	15.51	2,240	044
	CLP_4	16	8	$\begin{bmatrix} 3a, 3b \\ -5b \end{bmatrix}$	$\begin{bmatrix} 1,168 \times 10^{3} \\ -292 \times 10^{3} \end{bmatrix}$			
	CLP_1	5	24	$\begin{bmatrix} - & 3a \\ - & - & - \\ - & 4a \\ - & 5a \end{bmatrix}$	$\begin{bmatrix} -633 \times 10^3 \\ -75 \times 10^3 \\ -75 \times 10^3 \\ -75 \times 10^3 \end{bmatrix} -$			
TS-based	CLP_2	3	12	$\begin{vmatrix} \\ 1a \end{vmatrix}$	$1,464 \times 10^{3}$	15.00	0.040	C49
Multi-CLP	CLP_3	3	12	$1b^{-}$	$1,464 \times 10^{3}$	15.32	2,240	648
(ours)	CLP_4	8	32	$\begin{bmatrix} 2a & , & 2b \\ & & 3b \\ & & 4b \\ & & 5b \end{bmatrix}$	$\begin{array}{c c} -875 \times 10^{\overline{3}} \\ -875 \times 10^{\overline{3}} \\ -292 \times 10^{\overline{3}} \\ -219 \times 10^{\overline{3}} \\ -146 \times 10^{\overline{3}} \end{array}$			

bles III and IV to get an insight into the required bandwidth, arithmetic unit utilization, CONV layers throughput, and performance. We can see that SA-/TS-based Multi-CLP provides $1.34\times$ and $1.54\times$ more throughput improvement than Single-CLP design on 485T and 690T FPGAs, respectively. This improvement comes from its ability to make better use of the available arithmetic units. Specifically, Single-CLP can provide a useful work to the multipliers and adders only 74% and 64% of the time on 485T and 690T FPGAs, respectively, whereas SA-/TS-based Multi-CLP improves the arithmetic units utilization to about 96% and 98% on 485T and 690T FPGAs, respectively. Note that the arithmetic utilization is computed as discussed in Equation (18). Furthermore, SA-/TS-based Multi-CLP improves CONV layer throughput of Multi-CLP design by $1.02\times$ on 485T.

On the other hand, there is a trade-off between off-chip memory bandwidth and on-chip buffer size. Using large buffers saves off-chip memory bandwidth, whereas adopting small buffers results in a high bandwidth requirement. Here, we must emphasize that the proposed optimization framework gives a higher cost to bandwidth requirement than that given to buffer size. In other words, when more than one feasible design with the same low number of execution cycles is encountered, the design with minimum bandwidth requirement is chosen as optimal. We have adopted this strategy to minimize power expenses since on-chip memory modules consume less power than off-chip memory modules [44].

The results for single CLP and multiple CLPs accelerators for SqueezeNet 1.1 on 485T and 690T FPGAs with 16-bit fixed-point data are shown in Table VI. For SqueezeNet 1.1 accelerators on 485T FPGA, the single CLP implementation requires about 1.9× more cycles than the multiple CLPs implementation. Even more importantly, though, the results show that as the number of CONV layers increases, the search space also increases, making it more challenging to obtain an optimal multiple CLPs design using conventional iterative algorithms. Specifically, SA-based Multi-CLP processes images in a 1.04 ms and a 1.43 ms shorter time than Multi-CLP on 485T and 690T FPGAs, respectively. Additionally, although TS-based Multi-CLP produced a solution that requires more

TABLE IV: Single-CLP and Multi-CLP accelerators for AlexNet on 690T FPGA using 32-bit floating-point data.

Technique	CLP No.	Unrollin	g Factors	Layer	Cycles	Execution	DSPs	BRAMs
recinique	CLI NO.	T_n	T_m	Mapping	Cycles	Time (ms)	DSIS	DIANIS
Single-CLP [16]	CLP_1	9	64	$\begin{bmatrix} 1a & , & 1b & \\ -2a & , & 2b & \\ -3a & , & 3b & \\ -4a & , & 4b & \\ -5a & , & 5b & \end{bmatrix}$	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	17.69	2,880	758
	CLP_1	1	64	5a , $5b$	$1,168 \times 10^{3}$			
	CLP_2	1	96	4a , $4b$	$1,168 \times 10^{3}$			
Multi-CLP	CLP_3	2	64	3a , $3b$	$1,168 \times 10^{3}$	11.68	2,880	1,238
[18]	CLP_4	1	48	1a	$1,098 \times 10^3$	11.00	2,000	1,230
	CLP_5	1	48	1b	$1,098 \times 10^3$			
	CLP_6	3	64	2a , $2b$	$1,166 \times 10^{3}$			
	CLP_1	3	16	1a	$1,098 \times 10^3$			
	CLP_2	3	16	1b	$1,098 \times 10^{3}$			
	CLP_3	12	8	2a	$1,166 \times 10^{3}$		2,880	1,344
SA-based Multi-CLP	CLP_4	6	16	2b	$1,166\times10^3$	11.68		
(ours)	CLP_5	16	16	$\begin{bmatrix} 3a & , & 3b \\ 4a & , & 4b \\ 5a & & & \\ & & & \\ \end{bmatrix}$	$ \begin{array}{cccccccccccccccccccccccccccccccccccc$	-	2,000	2,000
	CLP_6	8	4	5 <i>b</i>	$1,168 \times 10^3$			
	CLP_1	6	32	$\begin{bmatrix} - & -\frac{2a}{3b} & - & - \\ - & -\frac{5a}{5a} & - & - \end{bmatrix}$	$ \begin{bmatrix} -583 \times 10^{3} \\ -795 \times 10^{3} \\ -795 \times 10^{3} \end{bmatrix} $			
TS-based Multi-CLP (ours)	CLP_2	3	48	$\begin{bmatrix} - & - & - & - & - & - & - & - & - & - $	$\begin{array}{c} -366\times10^{\overline{3}} \\ -523\times10^{\overline{3}} \\ -523\times10^{\overline{3}} \\ -292\times10^{\overline{3}} \end{array}$	11.81	2,880	1,348
	CLP_3	3	16	1b	$1,098 \times 10^3$]		
	CLP_4	6	32	$\begin{bmatrix} -\frac{2b}{4a}, & 4b \end{bmatrix}$	$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$			

TABLE V: AlexNet accelerators performance on 485T and 690T FPGAs using 32-bit floating-point data.

		Virtex-7	VX485T		Virtex-7 VX690T				
Technique	BW (GB/s)	Arith. Util.	Thr. (img/s)	Perf. (GFLOPs/s)	BW (GB/s)	Arith. Util.	Thr. (img/s)	Perf. (GFLOPs/s)	
Single-CLP [16]	1.40	74.1%	48.85	65.05	1.78	64.0%	55.40	73.77	
Multi-CLP [18]	1.38	95.6%	63.98	85.20	1.49	98.1%	85.55	113.92	
SA-based Multi-CLP (ours)	1.42	97.4%	65.32	86.98	1.35	98.1%	85.55	113.92	
TS-based Multi-CLP (ours)	1.33	95.7%	65.27	86.91	1.30	98.1%	84.67	112.75	

cycles than that of SA-based Multi-CLP, it processes images in a shorter time and with fewer DSP slices than Multi-CLP.

On the other hand, when comparing the performance of SA-based Multi-CLP in accelerating SqueezeNet 1.1 on VC707 and VC709 boards, we can notice that it used the available

resources to the maximum extent to improve the network throughput. Specifically, with a $1.29\times$ increase in computational resources in VC709 compared to VC707, SA-based Multi-CLP increased SqueezeNet 1.1 throughput by $1.30\times$, which demonstrates the scalability of the proposed accelerator.

TABLE VI: Single-/Multi-CLP accelerators for SqueezeNet on 485T and 690T FPGAs using 16-bit fixed-point data.

			Virtex	-7 VX485T				Virte	x-7 VX690T	
Technique	CLP	Unre	olling	Layer	Cl	CLP	Unre	olling	Layer	Cooler
	No.	T_n	T_m	Mapping	Cycles	No.	T_n	T_m	Mapping	Cycles
Single-CLP [16]	CLP_1	32	68	1 - 26	349×10^3	CLP_1	32	87	1 - 26	331×10^3
	CLP_1	6	16	2, 3, 6, 5	179×10^3	CLP_1	8	16	2, 6, 3, 5	125×10^{3}
	CLP_2	3	64	1, 8, 9, 12	183×10^{3}	CLP_2	3	64	1	115×10^{3}
Multi-CLP	CLP_3	4	64	all others	165×10^{3}	CLP_3	11	32	all others	133×10^{3}
[18]	CLP_4	8	64	7, 4, 16, 19	176×10^{3}	CLP_4	8	64	7, 4, 16	145×10^3
	CLP_5	8	128	26, 22, 25, 13	185×10^3	CLP_5	5	256	19, 26, 22, 25	144×10^{3}
	CLP_6	16	10	10	183×10^{3}	CLP_6	16	26	13, 10	141×10^{3}
	CLP_1	2	7	14	176×10^{3}	CLP_1	64	16	5, 23, 25, 26	139.5×10^3
	CLP_2	16	32	7, 13, 20, 25	179×10^{3}	CLP_2	3	64	1,12	132×10^{3}
	CLP_3	8	16	2, 16, 17	180×10^{3}	CLP_3	32	8	11,13	138×10^{3}
	CLP_4	13	19	5,10	148×10^{3}	CLP_4	8	22	2, 16, 24	139×10^{3}
SA-based Multi-CLP	CLP_5	43	13	26	$181 imes 10^3$	CLP_5	8	32	3, 4, 9	138×10^{3}
(ours)	CLP_6	11	4	3,8	176×10^{3}	CLP_6	8	32	10, 18, 20	139×10^{3}
	CLP_7	3	64	1, 6, 23, 15	177×10^{3}	CLP_7	8	43	19, 22	138×10^3
	CLP_8	16	13	4, 18, 24, 12	181×10^{3}	CLP_8	8	32	7, 6, 14	138×10^{3}
	CLP_9	- 8	26	9, 21, 22	172×10^{3}	CLP_9	7	8	8, 15	$\boxed{93 \times 10^{\overline{3}}}$
	CLP_{10}	16	8	11, 19	177×10^{3}	CLP_{10}	14	4	17, 21	129×10^{3}
	CLP_1	8	_ 26_	3, 22, 6	179×10^{3}	CLP_1	3	66	21,1	132×10^{3}
	CLP_2	8	22	7	170×10^{3}	CLP_2	12	4	5	$[138 \times 10^{3}]$
	CLP_3	11	53	26,9	182×10^{3}	CLP_3	16	32	6, 11, 7, 20, 25	141×10^{3}
	CLP_4	16	10	10	176×10^{3}	CLP_4	43	1	17	$85 \times 10^{\overline{3}}$
TS-based	CLP_5	3	_ 66	1,14,21,23	183×10^{3}	CLP_5	16	16	2, 4, 14	135×10^{3}
Multi-CLP	CLP_6	16	_ 22	4,18,25	175×10^{3}	CLP_6	16	2	33	100×10^{3}
(ours)	CLP_7	10	8	8,11,20	$oxed{183 imes 10^3}$	CLP_7	16	32	all others	139×10^{3}
	CLP_8	_ 11	_ 19	2,13,24	183×10^{3}	CLP_8	88	_ 22	19, 12, 24	133×10^3
	CLP_9	6	_ 32	_ 15, 16, 19	179×10^{3}	CLP_9	33	41	16	$\begin{array}{c c} 141 \times 10^3 \end{array}$
	CLP_{10}	10	8	5, 12, 17	178×10^{3}	CLP_{10}	8	_ 26	10	141×10^{3}
	22110			o, 12, 11	2.0 % 10	CLP_{11}	47	16	8, 26	140.5×10^3

In Tables VII and VIII, we show the configurations of single CLP and multiple CLPs accelerators designed for VGGNet and GoogLeNet on 690T FPGA with 16-bit fixed-point representation. Note that no results have been reported in Single-CLP [16] and Multi-CLP [18] accelerators for these two large CNNs. Therefore, with the sake of brevity, we only present the results for SA-based Single-/Multi-CLP as it turns out from previous experiments that SA algorithm provides more optimized solutions than TS algorithm in this specific problem.

For VGGNet, SA-based Multi-CLP reduces the total cycles compared to SA-based Single-CLP by about 676×10^3 cycles. Here, we see a significant difference compared to the results from AlexNet, which is due to the variance in the dimensions between CONV layers in VGGNet, where, for example, the first CONV layer has $\langle N, M \rangle$ as $\langle 3, 64 \rangle$, and those for last CONV layer is $\langle 512, 512 \rangle$. Hence, using a single CLP results in a large DSP resource under-utilization during the processing of some layers. For GoogLeNet, it requires around $2\times$ more cycles for the single CLP implementation than that of multiple CLPs. For

large CNNs, such as GoogLeNet which has 57 CONV layers, it becomes more challenging to find an optimized multiple CLPs configuration without using metaheuristics algorithms.

Additionally, we modeled the behavior of multiple CLPs design in Verilog and synthesized it using Xilinx Vivado Design Suite targeting the Xilinx Virtex-7 485T FPGA. We kept all synthesis properties to their default. Table IX shows a summary of the implementation and performance results for the proposed SA-based Multi-CLP approach and compares them with those reported in Single-CLP [16], GA-based Single-CLP [17], and Multi-CLP [18] approaches when they are used to accelerate AlexNet architecture. For each technique, the table provides the FPGA platform used for implementation, operating frequency, precision adopted to represent FMs and weights, design entry employed to describe the accelerator, hardware resource usage after placement and routing in terms of DSPs, BRAMs, FFs, and LUTs, and power efficiency (a.k.a., performance per Watt).

When comparing resource usage estimates with resource usage from Vivado's implementation report for SA-based Multi-CLP, one can note that the model underestimated the

TABLE VII: SA-based Single-/Multi-CLP accelerators for VGGNet on 690T FPGA using 16-bit fixed-point data.

Technique	CLP	Unro	olling	Layer	Cycles
rechnique	No.	T_n	T_m	Mapping	$(\times 10^3)$
SA-based Single-CLP (ours)	CLP_1	44	65	1 – 13	6,631
a.,	CLP_1	8	64	1, 8, 4	5 , 955
SA-based Multi-CLP	CLP_2	32	19	3, 7, 11	5,503
(ours)	CLP_3	48	2	13	4,976
	CLP_4	64	26	all others	5,841

DSP slices count by 53 DSP slices per CLP, on average. The reason is that Vivado's utilization report for SA-based Multi-CLP design accounts not only for DSPs used in CLP's computational engine, but also for those used in the control logic, address calculations, and loop indexing as well. When considering only those DSPs used to implement the computational engines for SA-based Multi-CLP, we found that the estimates match those reported by Vivado utilization report.

On the other hand, the implemented and predicted BRAM counts are perfectly matched. This is because in the implemented parametrized Verilog modules, we instantiate the exact required memory type, mode, and count using Xilinx HDL language templates. However, this does not apply to Single-CLP, GA-based Single-CLP, and Multi-CLP because they use high-level synthesis tools to compile C, OpenCL, and C++ codes into HDL codes, respectively. As can be seen from the results, designing with multiple CLPs makes better use of available resources than a single CLP design. Due to the dimensional mismatch issue in Single-CLP design, this design scheme could not take advantage of more than 82% of DSP slices. One can also note that the proposed SA-based Multi-CLP increases the usage of FF and LUT resources in Single-CLP by about 5% for each additional CLP. All this increase in resource usage contributes to improving the power efficiency by $3.65\times$.

Compared to GA-based Single-CLP, the proposed SA-based Multi-CLP has effectively made use of the available resources resulting in a $3.19\times$ increase in power efficiency even though GA-based Single-CLP uses 8-bit fixed-point weights and 16-bit fixed-point FMs while our accelerator uses 32-bit floating-point data and operates at a lower frequency. This is due to the high-level synthesis tool employed in [17] which constrains T_{in} to be from the set $\{1,2,4,8,16\}$ and T_{out} to be integer multiplicative of T_{in} , where T_{out} and T_{in} are the unified unrolling factors for the output matrix of CONV operation and the input vectors to CLP structures, respectively, which are determined by GA to minimize execution time. In contrast, this work allows unrolling factors to be whatever value yields the highest performance.

TABLE VIII: SA-based Single-/Multi-CLP accelerators for GoogLeNet on 690T FPGA using 16-bit fixed-point.

Technique	CLP	Unro	olling	Layer	Cycles	
тесппіque	No.	T_n	T_m	Mapping	$(\times 10^3)$	
SA-based Single-CLP (ours)	CLP_1	45	64	1 - 57	1,330	
	CLP_1	51	2	11, 23, 31, 43, 28	620	
	CLP_2	1	21	32	470	
	CLP_3	10	64	3, 19, 53	636	
	CLP_4	78	1	17,21,40	571	
	CLP_5	26	$-\frac{1}{2}$	8,15	564	
	CLP_6	8	20	36, 16	594	
	CLP_7	4	13	54	635	
	$\bar{C}\bar{L}P_8$	2	40	10,45	608	
	CLP_9	18	3	37, 38, 51, 56, 57	634	
	$\bar{C}\bar{L}\bar{P_{10}}$	7 - 7	13	18,34,50	587	
SA-based	$C\bar{L}P_{11}$	8	13	30,47	631	
Multi-CLP	CLP_{12}	9	32	12	635	
(ours)	CLP_{13}	14	1	27	464	
	CLP_{14}	3	66	1	637	
	CLP_{15}	1	11	49	122	
	CLP_{16}	6	11	13, 39, 44	555	
	CLP_{17}	2	38	35, 26, 48	636	
	CLP_{18}	6	27	25,42,55	602	
	CLP_{19}	12	20	9, 14, 24, 46	575	
	$\bar{C}\bar{L}\bar{P_{20}}$	4		5		
	CLP_{21}	9	8	7,52	521	
	$\bar{C}\bar{L}\bar{P_{22}}$		- - -	$-\frac{1}{2},\frac{1}{20},\frac{1}{41}$	631	
	CLP_{23}	99		4,6,22	 596	

On the other hand, the usage of DSP, FF, and LUT resources for Multi-CLP and the proposed SA-based Multi-CLP is almost consistent. However, the additional 168 BRAMs used in Multi-CLP led to a 0.4 Watt increase in its power consumption, which also caused its power efficiency to be about 0.87 GOPs/Watt lower than that of the proposed SA-based Multi-CLP.

VI. CONCLUSION

CNNs have shown great performance in a wide range of machine learning applications. However, they require significant computational power that cannot be met by general purpose processors. Nowadays, FPGAs have been found to deliver excellent performance in accelerating CNNs as they provide the best performance per Watt. Most of the existing FPGA-based accelerators use a single CLP for processing all CONV layers. However, the current trend in CNN acceleration is using multiple CLPs to make efficient use of hardware resources.

Tachnique	Technique Platform		Precision	Design		Resources		Power Efficiency	
Technique	Flatioriii	(MHz)	Frecision	Entry	DSP	BRAM	FF	LUT	(GOPs/Watt)
Single-CLP [16]	Virtex-7 VX485T	100	32-bit floating-point	С	2,309 (82%)	689 (34%)	219, 815 (36%)	146, 325 (48%)	3.31
GA-based Single-CLP [17]	Stratix-V GSD8	120	(8/16)-bit fixed-point	OpenCL	727 (37%)	1, 480 (58%)	N/A	119, 622 (46%)	3.79
Multi-CLP [18]	Virtex-7 VX485T	100	32-bit floating-point	C++	2,443 (87%)	812 (39%)	270, 991 (45%)	176, 876 (58%)	11.21
SA-based Multi-CLP (ours)	Virtex-7 VX485T	100	32-bit floating-point	Verilog	2, 452 (88%)	644 (31%)	269, 053 (44%)	176, 449 (58%)	12.08

TABLE IX: Implementation and performance summary of Single-/Multi-CLP accelerators optimized for AlexNet.

The multiple CLPs design involves optimizing the number of CLPs used, assignment of CONV layers to CLPs, the parallelization factors $\langle T_n, T_m \rangle$ for each CLP, and the tilling parameters $\langle T_r, T_c \rangle$ for each CONV layer. This imposes an intractable design space where it is almost impossible to find near-optimal configurations through exhaustive search algorithms. Therefore, in this work, we employed SA and TS algorithms to design CNN accelerators using multiple CLPs on FPGAs targeting the optimization of the number of execution cycles.

Experimental results on accelerating AlexNet, VGGNet, SqueezeNet 1.1, and GoogLeNet architectures demonstrated the effectiveness of the proposed optimization framework in providing efficient accelerators. The results also show the importance of using metaheuristics in finding different configurations for a given CNN using a systematic methodology. Compared to a single CLP design, SA-based Multi-CLP and TS-based Multi-CLP accelerate CNN computations by about $1.31\times-2.37\times$ when targeting Xilinx Virtex-7 FPGAs. Our implementation achieves a performance of 113.92 GFLOPs under 100 MHz working frequency.

In future work, we consider improving the quality of resource usage model in predicting DSP slices required by taking into account those DSPs used in control logic, loop indexing, address calculation, etc. Additionally, we target the adoption of data quantization technique for further acceleration. Reducing the bit-precision level for features and weights in turn allows multipliers and accumulators to be implemented using other logical elements not only using DSPs. This increase in computing units enhances CNN throughput further.

VII. ACKNOWLEDGMENT

The authors would like to thank King Fahd University of Petroleum and Minerals (KFUPM) for supporting this research and providing the computing facilities.

REFERENCES

- [1] Xinhui Hu, Xugang Lu, and Chiori Hori. "Mandarin speech recognition using convolution neural network with augmented tone features". In: *The 9th International Symposium on Chinese Spoken Language Processing*. 2014, pp. 15–18. DOI: 10.1109/ISCSLP.2014.6936674.
- [2] Mohamed Khalil-Hani and Liew Shan Sung. "A convolutional neural network approach for face verification". In: 2014 International Conference on High Performance Computing & Simulation (HPCS). 2014, pp. 707–714. DOI: 10.1109/HPCSim.2014.6903759.
- [3] Sachin Sudhakar Farfade, Mohammad J Saberian, and Li-Jia Li. "Multi-view face detection using deep convolutional neural networks". In: *Proceedings of the 5th ACM on International Conference on Multimedia Retrieval*. 2015, pp. 643–650. DOI: 10.1145/2671188.2749408.
- [4] Jilong Zheng, Yaowei Wang, and Wei Zeng. "CNN Based Vehicle Counting with Virtual Coil in Traffic Surveillance Video". In: 2015 IEEE International Conference on Multimedia Big Data. 2015, pp. 280–281. DOI: 10.1109/BigMM.2015.56.
- [5] Ruochen Wang and Zhe Xu. "A pedestrian and vehicle rapid identification model based on convolutional neural network". In: *Proceedings of the 7th International Conference on Internet Multimedia Computing and Service*. 2015, pp. 1–4. DOI: 10.1145/2808492.2808524.
- [6] Mian Mian Lau, King Hann Lim, and Alpha Agape Gopalai. "Malaysia traffic sign recognition with convolutional neural network". In: 2015 IEEE International Conference on Digital Signal Processing (DSP). 2015, pp. 1006–1010. DOI: 10.1109/ICDSP.2015.7252029.

- [7] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. "Imagenet classification with deep convolutional neural networks". In: *Advances in neural information processing systems* 25 (2012).
- [8] Ahmad Shawahna, Sadiq M. Sait, and Aiman El-Maleh. "FPGA-Based Accelerators of Deep Learning Networks for Learning and Classification: A Review". In: *IEEE Access* 7 (2019), pp. 7823–7859. DOI: 10.1109/ACCESS. 2018.2890150.
- [9] Xin Feng, Youni Jiang, Xuejiao Yang, et al. "Computer vision algorithms and hardware implementations: A survey". In: *Integration* 69 (2019), pp. 309–320. DOI: 10.1016/j.vlsi.2019.07.005.
- [10] Deepak Ghimire, Dayoung Kil, and Seong-heum Kim. "A Survey on Efficient Convolutional Neural Networks and Hardware Acceleration". In: *Electronics* 11.6 (2022), p. 945. DOI: https://doi.org/10.3390/electronics11060945.
- [11] Jason Cong and Bingjun Xiao. "Minimizing computation in convolutional neural networks". In: *International* conference on artificial neural networks. Vol. 8681. Springer. 2014, pp. 281–290. DOI: 10.1007/978-3-319-11179-7_36.
- [12] Andrew G Howard, Menglong Zhu, Bo Chen, et al. "Mobilenets: Efficient convolutional neural networks for mobile vision applications". In: *arXiv preprint arXiv:1704.04861* (2017). DOI: 10.48550/arXiv.1704. 04861.
- [13] Gwo-Jiun Horng, Min-Xiang Liu, and Chao-Chun Chen. "The Smart Image Recognition Mechanism for Crop Harvesting System in Intelligent Agriculture". In: *IEEE Sensors Journal* 20.5 (2020), pp. 2766–2781. DOI: 10. 1109/JSEN.2019.2954287.
- [14] He Jiang, Xiaoru Li, and Fatemeh Safara. "IoT-based agriculture: Deep learning in detecting apple fruit diseases". In: *Microprocessors and Microsystems* (2021), p. 104321. DOI: 10.1016/j.micpro.2021.104321.
- [15] Huimin Li, Xitian Fan, Li Jiao, et al. "A high performance FPGA-based accelerator for large-scale convolutional neural networks". In: 2016 26th International Conference on Field Programmable Logic and Applications (FPL). 2016, pp. 1–9. DOI: 10.1109/FPL.2016.7577308.
- [16] Chen Zhang, Peng Li, Guangyu Sun, et al. "Optimizing FPGA-based accelerator design for deep convolutional neural networks". In: *Proceedings of the 2015 ACM/SIGDA international symposium on field-programmable gate arrays.* 2015, pp. 161–170. DOI: 10.1145/2684746.2689060.
- [17] Naveen Suda, Vikas Chandra, Ganesh Dasika, et al. "Throughput-optimized OpenCL-based FPGA accelera-

- tor for large-scale convolutional neural networks". In: *Proceedings of the 2016 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays.* 2016, pp. 16–25. DOI: 10.1145/2847263.2847276.
- [18] Yongming Shen, Michael Ferdman, and Peter Milder. "Maximizing CNN accelerator efficiency through resource partitioning". In: 2017 ACM/IEEE 44th Annual International Symposium on Computer Architecture (ISCA). 2017, pp. 535–547. DOI: 10.1145/3079856. 3080221.
- [19] Wenyan Lu, Guihai Yan, Jiajun Li, et al. "FlexFlow: A Flexible Dataflow Accelerator Architecture for Convolutional Neural Networks". In: 2017 IEEE International Symposium on High Performance Computer Architecture (HPCA). 2017, pp. 553–564. DOI: 10.1109/HPCA.2017. 29
- [20] Ibrahim H Osman and James P Kelly. "Meta-heuristics: an overview". In: *Meta-heuristics* (1996), pp. 1–21. DOI: 10.1007/978-1-4613-1361-8 1.
- [21] LM Rasdi Rere, Mohamad Ivan Fanany, and Aniati Murni Arymurthy. "Simulated annealing algorithm for deep learning". In: *Procedia Computer Science* 72 (2015), pp. 137–144. DOI: 10.1016/j.procs.2015.12.114.
- [22] Forrest N Iandola, Song Han, Matthew W Moskewicz, et al. "SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and< 0.5 MB model size". In: *arXiv* preprint arXiv:1602.07360 (2016). DOI: 10.48550/arXiv. 1602.07360.
- [23] Karen Simonyan and Andrew Zisserman. "Very deep convolutional networks for large-scale image recognition". In: arXiv preprint arXiv:1409.1556 (2014). DOI: 10.48550/arXiv.1409.1556.
- [24] Christian Szegedy, Wei Liu, Yangqing Jia, et al. "Going deeper with convolutions". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2015, pp. 1–9.
- [25] Ahmad Shawahna, Sadiq M. Sait, Aiman El-Maleh, et al. "FxP-QNet: A Post-Training Quantizer for the Design of Mixed Low-Precision DNNs With Dynamic Fixed-Point Representation". In: *IEEE Access* 10 (2022), pp. 30202– 30231. DOI: 10.1109/ACCESS.2022.3157893.
- [26] Mannhee Cho and Youngmin Kim. "FPGA-Based Convolutional Neural Network Accelerator with Resource-Optimized Approximate Multiply Accumulate Unit". In: *Electronics* 10.22 (2021), p. 2859. DOI: 10.3390/electronics10222859.
- [27] Louis-Noel Pouchet, Peng Zhang, Ponnuswamy Sadayappan, et al. "Polyhedral-based data reuse optimization for configurable computing". In: *Proceedings of the ACM/SIGDA international symposium on Field pro-*

- grammable gate arrays. 2013, pp. 29–38. DOI: 10.1145/2435264.2435273.
- [28] Samuel Williams, Andrew Waterman, and David Patterson. "Roofline: an insightful visual performance model for multicore architectures". In: *Communications of the ACM* 52.4 (2009), pp. 65–76. DOI: 10.1145/1498765. 1498785.
- [29] Xilinx. *Vivado Design Suite Product Guide: Floating-Point Operator v7.1 [Online]*. Available: https://docs.xilinx.com/v/u/en-US/pg060-floating-point (2020).
- [30] Xilinx. *User Guide: 7 Series FPGAs Memory Resources* [Online]. Available: https://docs.xilinx.com/v/u/en-US/ug473_7Series_Memory_Resources (2019).
- [31] M Sait Sadiq and Youssef Habib. "Iterative Computer Algorithms with Applications in Engineering: Solving Combinatorial Optimization Problems". In: *IEEE, Los Alamitos, CA* (1999), p. 387.
- [32] Sadiq M Sait and Habib Youssef. *VLSI physical design automation: theory and practice*. Vol. 6. World Scientific, 1999.
- [33] Scott Kirkpatrick, C Daniel Gelatt Jr, and Mario P Vecchi. "Optimization by simulated annealing". In: *science* 220.4598 (1983), pp. 671–680. DOI: 10.1126/science.220.4598.671.
- [34] Vladimir Cerny. "Thermodynamical approach to the traveling salesman problem: An efficient simulation algorithm". In: *Journal of optimization theory and applications* 45.1 (1985), pp. 41–51. DOI: 10.1007/BF00940812.
- [35] Nicholas Metropolis, Arianna W Rosenbluth, Marshall N Rosenbluth, et al. "Equation of state calculations by fast computing machines". In: *The journal of chemical physics* 21.6 (1953), pp. 1087–1092. DOI: 10.1063/1. 1699114.
- [36] Habib Youssef, Sadiq M Sait, and Hakim Adiche. "Evolutionary algorithms, simulated annealing and tabu search: a comparative study". In: *Engineering Applications of Artificial Intelligence* 14.2 (2001), pp. 167–181. DOI: 10.1016/S0952-1976(00)00065-8.
- [37] Fred Glover. "Tabu search—part I". In: *ORSA Journal on computing* 1.3 (1989), pp. 190–206. DOI: 10.1287/ijoc.1.3.190.
- [38] Fred Glover. "Tabu search—part II". In: *ORSA Journal on computing* 2.1 (1990), pp. 4–32. DOI: 10.1287/ijoc.2.
- [39] Fred Glover and Eric Taillard. "A user's guide to tabu search". In: *Annals of operations research* 41.1 (1993), pp. 1–28. DOI: 10.1007/BF02078647.

- [40] Fred Glover and Manuel Laguna. "Tabu search". In: *Handbook of combinatorial optimization*. Springer, 1998, pp. 2093–2229. DOI: 10.1007/978-1-4613-0303-9_33.
- [41] Olga Russakovsky, Jia Deng, Hao Su, et al. "Imagenet large scale visual recognition challenge". In: *International journal of computer vision* 115.3 (2015), pp. 211–252. DOI: 10.1007/s11263-015-0816-y.
- [42] Xilinx. *User Guide: VC707 Evaluation Board for the Virtex-7 FPGA [Online]*. Available: https://docs.xilinx.com/v/u/en-US/ug885_VC707_Eval_Bd (2019).
- [43] Xilinx. *User Guide: VC709 Evaluation Board for the Virtex-7 FPGA [Online]*. Available: https://docs.xilinx.com/v/u/en-US/ug887-vc709-eval-board-v7-fpga (2019).
- [44] Paulo Garcia, Deepayan Bhowmik, Robert Stewart, et al. "Optimized memory allocation and power minimization for FPGA-based image processing". In: *Journal of Imaging* 5.1 (2019), p. 7. DOI: 10.3390/jimaging5010007.