





SqueezeJet: High-Level Synthesis Accelerator Design for Deep Convolutional Neural Networks

Panagiotis G. Mousouliotis^(✉)  and Loukas P. Petrou 

Division of Electronics and Computer Engineering,
Department of Electrical and Computer Engineering, Faculty of Engineering,
Aristotle University of Thessaloniki, 54124 Thessaloniki, Greece
pmousoul@ece.auth.gr, loukas@eng.auth.gr

Abstract. Deep convolutional neural networks have dominated the pattern recognition scene by providing much more accurate solutions in computer vision problems such as object recognition and object detection. Most of these solutions come at a huge computational cost, requiring billions of multiply-accumulate operations and, thus, making their use quite challenging in real-time applications that run on embedded mobile (resource-power constrained) hardware. This work presents the architecture, the high-level synthesis design, and the implementation of SqueezeJet, an FPGA accelerator for the inference phase of the SqueezeNet DCNN architecture, which is designed specifically for use in embedded systems. Results show that SqueezeJet can achieve 15.16 times speed-up compared to the software implementation of SqueezeNet running on an embedded mobile processor with less than 1% drop in top-5 accuracy.

Keywords: DCNN accelerator · FPGA · High-level synthesis

1 Introduction

Since the impressive results of AlexNet deep convolutional neural network (DCNN) in the Image-Net Large-Scale Vision Recognition Challenge (ILSVRC) in [1], DCNN research activity has seen exponential growth with the trend being deeper architectures accompanied by higher accuracies [2, 3]. Following this trend, research in DCNN FPGA accelerators provides solutions that use high-end costly FPGA devices and aim at the datacenter rather than the mobile applications [4–6]. An exception to the most-accurate-network trend in the DCNN architecture research, is SqueezeNet¹ (SqN) [7, 12], an AlexNet-level accuracy architecture which reduces dramatically the number of MACs and network parameters, requiring half of the MACs and fifty times less parameters compared to AlexNet. Even though the SqN DCNN architecture is more suitable than others for use in embedded mobile applications, it is still computationally

¹ In this work, SqueezeNet refers to SqueezeNet v1.1.

very demanding and cannot be used in applications running on an embedded mobile processor.

The contribution of this work is the design of SqueezeJet (SqJ), a small FPGA convolutional (conv) layer accelerator for SqN, that can be used as a coprocessor to an embedded mobile processor and enable the development of mobile computer vision (CV) applications. Specifically, the SqJ design: (1) deals with the challenge of the implementation of a single accelerator for multiple conv layers with variable input arguments, (2) implements streaming input/output (I/O) interfaces which, after the initialization phase, consume and produce data pixel-by-pixel², (3) uses a sophisticated hardware (HW) mechanism, which mimics software (SW) pointers to the rows of a two-dimensional array, taking advantage of the spatial locality of data and minimizing unnecessary data movement, (4) presents the possibilities of high-level synthesis (HLS) design by using the Xilinx Vivado HLS (VHLS) tool, (5) is implemented on a low-end FPGA system on chip (SoC) device, the Xilinx XC7Z020, using the Xilinx SDSoC tool, and (6) it achieves 80.29% ILSVRC12 top-5 accuracy when it is used for the inference phase of SqN. To the best of the authors' knowledge, the current work presents the first low-end FPGA SoC (XC7Z020) DCNN implementation which achieves 80.29% ILSVRC12 top-5 accuracy.

The rest of this paper is organized as follows: Sect. 2 presents related work. Section 3 is an introduction to the conv layer's operation. Section 4 presents the architecture, the HLS design, and the implementation of the SqJ accelerator. Section 5 shows results related to the performance, the accuracy, and the power consumption of SqJ. Finally, Sect. 6 concludes the paper and proposes future work.

2 Related Work

Works related to DCNN FPGA accelerators can be classified into two main categories; those which accelerate only the conv layer and those which accelerate two or more layer types of a DCNN.

Conv layer accelerators: *Zhang et al.* [4] designed an architecture template for the conv layer using loop tiling, loop arrangement based on data dependencies, computation optimizations (loop unrolling and pipelining), and optimizations for efficient data reuse. Using the parameters of the template and the roofline model, they performed design space exploration (DSE) and found the optimal solution which defined the parameters of their accelerator. A similar approach is followed by *Motamedi et al.* [5] starting with a completely different architectural template. Specifically, they designed their template to take advantage of all the possible forms of parallelism; intra/inter-kernel and inter-output. They eventually used the design parameters and proceeded as in the aforementioned work. Both of these works use DSE to minimize the execution time of the accelerator and 32-bit floating-point arithmetic.

² A pixel is comprised by all the channels at a specific (x, y) location in the future map volume (see Fig. 1).

Multi-layer accelerators: *Qiu et al.* [8] developed a dynamic-precision data quantization flow and designed a dynamic-precision 16-bit fixed-point accelerator which is capable of accelerating conv, fully connected (FC), and pooling layers. Their implementation is used to accelerate the VGG16-SVD DCNN, which is the VGG16 DCNN with reduced weight matrices for the FC layers; SVD is used for the weight matrix reduction. This accelerator also uses a huge amount of FPGA resources to accelerate one of the most computational demanding DCNNs, requiring 15470 million MACs for a single forward pass. *Gschwend* [9] converted all the layers, except the last global pooling layer, of the SqueezeNet v1.0 DCNN architecture to conv layers and accelerated, using floating-point arithmetic, the new DCNN, called ZynqNet, using VHLS. *Gokhale et al.* [10] designed and implemented nn-X, a complete low-power system for DCNN acceleration composed from a host processor, a coprocessor, and external memory. The coprocessor consists of an array of processing elements which can perform convolution, sub-sampling, and non-linear functions. *Ma et al.* [11] designed an accelerator that supports conv, pooling and fully-connected layers by following a strategy that minimizes computing latency, partial sum storage, access of on-chip buffer, access of external memory, and uses loop optimization techniques. Their accelerator uses 8–16 bit dynamic fixed point arithmetic and it is evaluated by accelerating the VGG-16 DCNN.

SqJ is a conv layer accelerator and it uses fixed-point arithmetic for both parameters (8 bits) and activations (16 bits), which results in considerable savings in both the resources and the power consumption compared to floating-point implementations [4,5,9]. Furthermore, even though works in [8,10,11] use fixed-point arithmetic, they require large costly FPGA devices for their implementation.

3 Convolutional Layer Basics

The conv layer of a DCNN can be described by:

$$\begin{aligned}
 FM_o(y_o, x_o, c_o) = & \\
 \sum_{k_h=0}^{K_h-1} \sum_{k_w=0}^{K_w-1} \sum_{c_i=0}^{C_i-1} & FM_i((y_o \cdot S + k_h), (x_o \cdot S + k_w), c_i) \cdot W(c_o, k_h, k_w, c_i) \quad (1) \\
 & + B(c_o),
 \end{aligned}$$

where FM_o , FM_i are the output and the input feature maps (fmaps) respectively, and W , B are the weight and bias parameters respectively. The y , x , c , represent the vertical, the horizontal, and the channel dimensions of the fmaps, S is the stride, and k_h , k_w are the vertical and horizontal dimensions of the kernel³.

The second line in Eq.1 represents a 3D convolution between FM_i , and C_o number of 3D kernels, the weight parameters. To calculate the first output

³ In this work, kernel has the same meaning as filter.

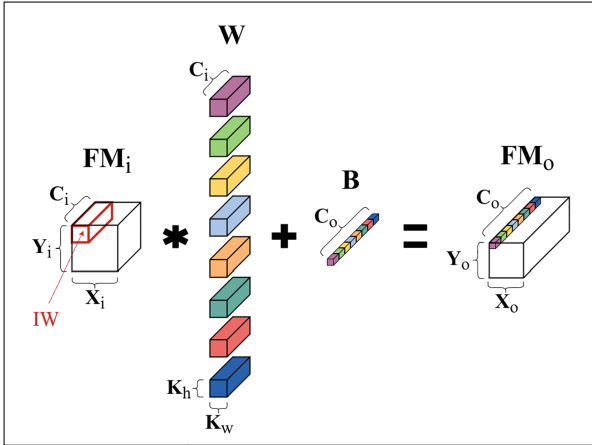


Fig. 1. Calculation of the channels of the first pixel of FM_o . The number of 3D kernels is equal to C_o , the number of output channels.

channel of the first output pixel of FM_o , an input window IW of the FM_i , of size $IW[K_h][K_w][C_i]$ is multiplied element-wise with kernel $W[0][K_h][K_w][C_i]$ and all partial results are accumulated to a single value. This value is then added to the respective bias term, $B(0)$, to produce the first output channel of the first output pixel of FM_o . This procedure is depicted in Fig. 1, which shows the calculation of all the channels of the first pixel of FM_o . To calculate all the elements of FM_o , the IW is moved vertically and horizontally by $y_o \cdot S$ and $x_o \cdot S$ respectively, and the above procedure is repeated. The resulted size of each Y , X dimension of the FM_o is calculated by:

$$\begin{aligned} Y_o &= (Y_i - K_h + 2 \cdot P) / S + 1 \\ X_o &= (X_i - K_w + 2 \cdot P) / S + 1, \end{aligned} \tag{2}$$

where P denotes the number of pixels added for padding the FM_i . In all the practical cases $Y_i = X_i$ and $K_h = K_w$.

An activation function always follows a conv layer. Thus, it is convenient, from an implementation point of view, to include the activation layer in the conv layer. In this case, the output of the conv layer becomes:

$$FM_{o,a}(y_o, x_o, c_o) = f(FM_o(y_o, x_o, c_o)), \tag{3}$$

where $f()$ is the activation function used in the specific DCNN, e.g. the Rectified linear unit (ReLU) described by:

$$f(x) = \max(0, x) \tag{4}$$

The accelerator described in the next section, accelerates this fused convolution-activation layer with output given by Eq. 3.

4 The SqueezeJet Accelerator

SqN is a DCNN architecture focused in reducing the network parameter count for a given accuracy. Specifically, SqN achieves AlexNet-level accuracy with fifty times less parameters, making its model sufficiently small to be stored in on-chip FPGA memories and removing the need for off-chip memory access. For an FPGA accelerator, such as SqJ, implemented on a device with a few Mbits of block RAM (BRAM) resources, this means that the parameters (weight and bias values) of a single layer can fit in the BRAMs. Thus, for the calculation of the FM_o of a specific conv layer by an accelerator, the following procedure is required: the parameters are brought from off-chip memory and stored to BRAMs, the FM_i is streamed from off-chip memory in the accelerator, the calculation of FM_o pixel(s) takes place, and the resulting FM_o pixel(s) are streamed back to the off-chip memory. Having the layer’s parameters stored on-chip is a big advantage as they will be reused for the calculation of each pixel of FM_o .

Following the architecture principle “make the common case fast”, SqJ is designed to accelerate conv layers described by Eq. 3 with stride limited to one; it can be used for the acceleration of all the SqN conv layers except the first one, which can be implemented as a distinct module. All SqN conv layers, except the first one, share the following common characteristics: (1) a stride equal to 1, (2) an input channel dimension with a greatest common divisor (GCD) equal to 16 and (3) an output channel dimension which is divisible by a power of 2. SqJ uses all these three characteristics to accelerate a conv layer; the first SqN conv layer does not have characteristics (1) and (2). Implementing SqJ to support the first SqN conv layer would significantly degrade the acceleration of the other 17 conv layers (25 conv modules) of SqN.

This section describes the architecture, the high-level synthesis design, and the implementation of SqJ.

4.1 Architecture

Data organization: The data organization of all the convolution array arguments is shown in Eq. 1. This data organization is imposed by the 3D convolution operation; it is necessary to read all the input channels of the IW pixels in order to be able to calculate a single output channel. Because SqJ accelerates 3D convolutions, the design of a streaming architecture is not possible, but it is possible to design the accelerator to use streaming I/O interfaces.

Buffering: The implementation of the 3×3 convolution introduces an input data access pattern which requires multiple lines of the input. Because FM_i data is streamed in the accelerator, FM_i data lines must be buffered. In the general case, the size of the input tile buffer ITB is:

$$ITB = K \cdot Y_i \cdot C_i, \quad (5)$$

where K denotes the kernel size (considering that $K = K_h = K_w$, see Fig. 1), and Y_i and C_i denote the width and the channels of FM_i respectively. In the

SqJ case, where support for up to $3 \times 3 \times C_i$ 3D kernels is required, $K = 3$ and $ITB_{3 \times 3}$ is implemented as a set of 3 line buffers whose access is determined by a pointer array. In this way, $ITB_{3 \times 3}$ shifts down the FM_i without the need for any data shift to take place; only the lowest, as defined by the pointer array, line buffer gets updated. This shift mechanism is also used by the input tile window buffer $ITWB$ (depicted as IW in Fig. 1) to update only one of its columns as it shifts horizontally on the ITB , taking advantage of the spatial locality of the input data. Figure 2 shows the internal organization of ITB and the operation of pointer array for $ITB_{3 \times 3}$. Apart from the ITB and $ITWB$, buffers are used to store the weights, the bias, and one pixel of FM_o . The buffer used to store the FM_o pixel could be omitted if each output channel was calculated serially, but buffering is required to calculate multiple output channels in parallel and to stream them out of the accelerator in order.

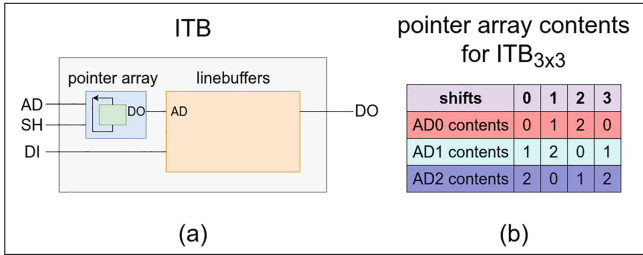


Fig. 2. ITB : (a) schematic and (b) pointer array content of $ITB_{3 \times 3}$ after a number of shifts. AD denotes memory address, SH denotes a shift signal, and DI, DO denote data input and output respectively.

Parallelism exploitation: SqJ takes advantage of the fact that SqN increases in the input channel dimension and, with the exception of the first conv layer, all conv layers' input channels have a GCD equal to $CI_{min} = 16$. The accelerator is designed to perform CI_{min} multiplications concurrently. These CI_{min} products are then fed to an accumulator unit which outputs a CI_{min} MAC result. The combination of the CI_{min} concurrent multiplications plus the accumulator unit forms a $MAC-CI_{min}$ unit which is pipelined. CI_{min} is a design parameter and can be easily modified according to the architecture of a different DCNN. This intra-kernel parallelism has the advantage that it exploits parallelism in the input channel dimension and it is independent from the kernel size K . Thus, SqJ can be easily modified to support kernel sizes larger than 3×3 . Another form of parallelism that is used is the concurrent calculation of multiple output channels for a specific output pixel. This is achieved by splitting the weights buffer in 2^n ($n = 1, 2, 3, \dots$) equal groups of 3D kernels and assigning them to 2^n $MAC-CI_{min}$ units.

Operation: First step in the operation of SqJ is the initialization of the input buffers. Weights and bias are brought from off-chip memory and, only in the

case where kernel $K = 3$, the ITB is initialized. After the initialization step, the convolution begins:

- For each row of FM_o : (1) only if $K = 3$, the ITB is shifted down (in FM_i) and two FM_i pixels are written in the empty line buffer, and (2) only if $K = 3$, the $ITWB$ is initialized with ITB data.
- For each column of each row of FM_o : (1) ITB is updated with a new FM_i pixel and $ITWB$ is updated with a new ITB column, (2) the weight buffers and the $ITWB$ are used to calculate one pixel of FM_o , and (3) the computed pixel is written back to off-chip memory.

4.2 Implementation

FPGA algorithm acceleration is not as trivial as implementing an algorithm in SW using a general purpose programming language such as C/C++. Even though HLS tools advertise the automatic generation of FPGA IP cores from C/C++ code, this process requires knowledge of the architecture of the FPGA device, knowledge of the internals of the HLS compiler [13], and use of a C/C++ coding style compatible with the HLS capabilities. This paragraph describes the process of generating an IP core for SqJ using the Xilinx VHLS tool and implementing it as a real application using the SDSoC tool.

Coding style: Hardware description languages (HDL) books warn the reader that if the designer cannot understand what logic circuit is described by the HDL code, then the design tool is not likely to synthesize the circuit that the designer is trying to model [14]. The same applies for the C/C++ code used as input to VHLS. A result of this coding style is the implementation of ITB shown in Fig. 2, which uses the HW model of pointers to the rows of a two-dimensional array. Even though VHLS simplifies the HW design of an algorithm, it doesn't provide a straightforward way for making a design scalable as it is the case with the combination of generate constructs and generics/parameters used in HDLs.

Interfaces: The SqJ IP core requires buffers for the weights, the bias, the ITB (FM_i), and the FM_o buffer for storing the output pixel. Three FIFO interfaces are used to stream data in and out of the IP core; one for streaming in the parameter (weights, bias) data, one for the FM_i data, and one for the output (FM_o) data. In addition, an AXI-Lite interface is used for acquiring the rest of the required HW function arguments. The SDSoC tool is used for interface synthesis.

Optimizations: VHLS provides many optimization possibilities both in terms of performance and resource usage [15].

- **Parallelism:** SqJ exploits parallelism in: (a) the input channel dimension (intra-kernel parallelism), by calculating the result of CI_{min} MACs every clock cycle of the operation of the pipelined MAC- CI_{min} unit, and (b) the output channel dimension, by calculating 2^n ($n = 1, 2, 3, \dots$) output channels concurrently. Parallelism in (a) requires a CI_{min} -wide data register and

partitioning the operand buffers (array partitioning) in a way which makes them able to provide CI_{min} outputs concurrently. Parallelism in (b) requires 2^n *ITWB* buffers and the same number of MAC-*CI* units.

- **Arbitrary precision types:** To further decrease the model size of SqN and reduce the amount of logic required by SqJ, fixed-point quantization in both the parameters and the FM_i is used. Specifically, Ristretto [16] is used to specify the proper quantization of the parameters (weights and bias) and the FM_i . Parameters are quantized at 8 bits (1 bit integer + 7 bits fractional) and FM_i at 16 bits (13 bits integer + 3 bits fractional), achieving 0.88% top-5 accuracy loss without performing any fine-tuning.

In Fig. 3, the block diagram of SqJ, implemented (for simplicity) with 4 MAC- CI_{min} units, is shown. Since the parallelization factor is equal to 4, the sizes of the buffers are: (1.179648/4) Mbits for the $weights_i$, (2048/4) bits for the $bias_i$, 344.064 Kbits for the ITB, 73.728 Kbits for each $ITWB_i$, and (4096/4) bits for the $fmap_{o_i}$. Table 1 presents the FPGA resources required for the implementation of **conv_10**, the accelerator of the first SqN conv layer, and **SqJ**, in an 8 MAC-16 unit configuration, on the XC7Z020 FPGA SoC. The **conv_10 + SqJ** implementation is the one used in the results of the next Section.

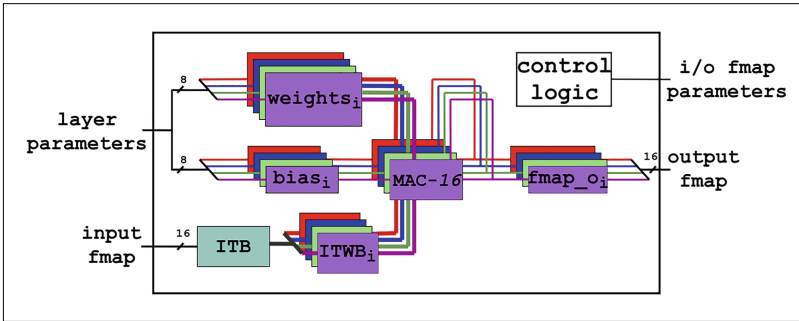


Fig. 3. SqJ block diagram implemented with 4 MAC- CI_{min} units. Bold lines denote $CI_{min} = 16$ times the data size shown at the left side of the figure.

Table 1. Resource utilization of conv_10 and SqJ on the XC7Z020 FPGA SoC

Resource	Available	conv_10		SqJ		conv_10 + SqJ	
		Util.	Util. %	Util.	Util. %	Util.	Util. %
LUT	53200	9405	17.678	12692	23.857	20631	38.780
LUTRAM	17400	707	4.063	726	4.172	1273	7.316
FF	106400	15459	14.529	18114	17.024	30554	28.716
BRAM	140	13	9.285	124	88.571	134.5	96.071
DSP	220	37	16.818	149	67.727	186	84.545

5 Performance Evaluation

Table 2 presents the per-layer execution times, the accuracy, and the chip power consumption⁴ of SqN implemented on 4 different processing unit configurations, an Intel Core i3-7100U@2.4 GHz core (Intel NUC), an ARM Cortex-A53@1.2 GHz core (Raspberry Pi 3 (RPI3) Model B V1.2), an ARM Cortex-A9@667 MHz core (Xilinx ZC702), and an ARM Cortex-A9@667 MHz core with the SqJ@100 MHz accelerator in an 8 MAC-16 unit configuration (Xilinx ZC702).

Table 2. SqN application execution time/accuracy/power results

Processing Unit	NUC	RPI3	ZC702	ZC702
	Intel i3@2.4GHz	ARM A53@1.2GHz	ARM A9@667MHz	ARM A9@667MHz conv_10@100MHz SqJ@100MHz
SqN Implementation Accuracy (bits)				
Activations	32			16
Weights, Bias	32			8
SqN Application Per-Layer Execution Time Results (ms)				
Load Image	0.1761	1.2137	21.3210	54.4263
0:Conv	25.3118	131.5186	297.2426	26.2756
1:Maxpool	2.0531	18.2868	28.7206	22.7574
2:Fire	16.1473	142.7623	446.1214	32.6526
3:Fire	17.0744	150.7194	474.1013	34.7981
4:Maxpool	1.3333	13.4446	27.3646	18.0916
5:Fire	13.5606	124.2315	450.0168	17.7738
6:Fire	14.5805	135.3108	482.2875	18.9882
7:Maxpool	0.6023	7.1370	14.4114	9.4158
8:Fire	7.4712	69.1218	257.9832	8.6426
9:Fire	7.8755	72.4013	273.4599	8.8704
10:Fire	13.1197	125.8514	497.6390	12.2322
11:Fire	13.6331	132.5349	517.09514	12.7946
12:Conv	34.7681	324.9181	1257.4682	33.9618
13:Fixed2float	0.0001	0.0004	0.0003	15.4479
13:Avgpool	1.5295	4.3149	5.7796	5.7085
14:Softmax	0.0260	0.1528	0.2212	0.2220
Total Conv	162.4322	1395.4275	4892.9386	174.9867
Total Merge	13.74	13.89	60.43	31.97
Total Maxpool	3.99	38.87	70.49	50.26
Total	169.2627	1453.9202	5051.2337	333.0595
FPS	5.907	0.687	0.198	3.002
SqN ILSVRC12 Accuracy Results (%)				
Top-1	58.38			57.46
Top-5	81.01			80.29
SqN Application CPU/SoC Power Consumption Results (Watts)				
Technology	14nm	n/a	28nm	28nm
Chip Power	5.3253	2.9	1.569	2.275
FPS/W	1.109	0.237	0.126	1.319

SqN is a single floating point precision C/C++ Linux application accelerated with single-instruction multiple-data (SIMD) instruction set extensions (Intel AVX, ARM NEON) and executed on a single core of the target CPU-only processing systems. In the case where the SqJ accelerator is used, the implementation uses 16 bits for the activations and 8 bits for the weights and bias. GCC

⁴ In the case of the ARM Cortex-A53, we measure RPI3 board power consumption, because there is no way to acquire power consumption measurements or estimations for the Broadcom 2837 SoC.

(version 6.3.0 for the Intel (64-bit) and RPI3 (32-bit) configurations, and version 6.2.1 for the ZC702 (32-bit) configuration) with the -O3 flag is used to build the SqN Linux application. Execution times are an average of 1000 inference iterations. Power consumption is acquired: (1) using Intel PCM⁵ while the processing system executes 1000 SqN iterations, in the case of the Intel i3 CPU, (2) using a power plug and measuring board power consumption, in the case of RPI3, and (3) using Xilinx XPE⁶ in the case of Xilinx ZC702. Accuracy is evaluated using the Ristretto⁷ tool.

Results show that the SqJ configuration achieves an 15.16x execution time speedup in SqN inference when compared to the ARM A9 core configuration, 4.36x execution time speedup in SqN inference when compared to the ARM A53 core, and similar convolution performance (see **Total Conv** in Table 2) to the Intel i3 core configuration, with less than 1% top-5 accuracy loss. In terms of performance per Watt, frames per second per Watt (FPS/W), the SqJ implementation is 10.46 times better than the ARM A9 core configuration; again, with less than 1% accuracy loss. The **Load Image** execution time in the SqJ implementation includes the conversion of the image from 32-bit floating point to 16-bit fixed point; that’s why it takes more than double of the ARM A9 corresponding time. Because of the use of lower precision for the activations, **Total Merge** (merge operations are included in the Fire layers) and **Total Maxpool** operations require much less time than the ARM A9 implementation. Furthermore, the Maxpool layers require 15% of the **Total SqJ** implementation time and could be incorporated in a future SqJ implementation. Table 3 summarizes the characteristics of the SqJ implementation.

Table 3. SqJ (conv_l0+SqJ) implementation summary

	SqueezeNet v1.1
FPGA	Zynq XC7Z020
Frequency (MHz)	100
Design Tool	Vivado HLS
DCNN Ops (GOPs)	0.7755
Precision	8-16 bits
DSP (Util.)	186 (84.5%)
BRAM (Util.)	134.5 (96%)
LUT (Util.)	20631 (38.8%)
LUTRAM (Util.)	1273 (7.3%)
FF (Util.)	30554 (28.7%)
Conv Latency/Image (ms)	175
Throughput (GOPs)	4.43
Top-5 ILSVRC12 Accuracy	80.29%

⁵ <https://www.intel.com/software/pcm>.

⁶ <https://www.xilinx.com/products/technology/power/xpe.html>.

⁷ <https://github.com/pmgysel/caffe>.

6 Conclusion

In this paper, we present the design and the implementation of SqJ, an FPGA-based convolution layer accelerator which can be used to boost the performance of an embedded mobile processor running a CV task. The accelerator, consisting of a buffering architecture and multiple computational units, is designed using the Xilinx Vivado HLS tool. The Ristretto tool is used to squeeze the SqN DCNN in the Xilinx XC7Z020 FPGA SoC, and the Xilinx SDSoC tool is used to deploy SqJ accelerated SqN to the XC7Z020 device. To the best of our knowledge, our work is the first one which implements the SqN DCNN in a small FPGA SoC device, such as the XC7Z020, and achieves 80.29% top-5 ILSVRC12 accuracy (using XC7Z020). Results show that SqJ accelerates by 15.16 times the SqN inference execution time of an embedded mobile processor while being 10.46 times more power efficient with less than 1% top-5 accuracy drop. Improvements to the HLS SqJ design could include: (1) Maxpool layer support, since they require considerable amount (15%) of the total inference time on a mobile ARM core, and (2) streaming execution, to avoid memory accesses for fmaps (requires additional BRAM resources). Future work could use an enhanced version of SqJ as a template and perform multiobjective optimization for finding the best solution in terms of performance, resources, accuracy, power, and cost.

References

1. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: *Advances in Neural Information Processing Systems*, pp. 1097–1105 (2012)
2. Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., Rabinovich, A.: Going deeper with convolutions. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1–9 (2015)
3. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778 (2016)
4. Zhang, C., Li, P., Sun, G., Guan, Y., Xiao, B., Cong, J.: Optimizing FPGA-based accelerator design for deep convolutional neural networks. In: *Proceedings of the 2015 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, pp. 161–170. ACM, 2015 February
5. Motamedi, M., Gysel, P., Akella, V., Ghiasi, S.: Design space exploration of FPGA-based deep convolutional neural networks. In: *2016 21st Asia and South Pacific, Design Automation Conference (ASP-DAC)*, pp. 575–580. IEEE, January 2016
6. Ovtcharov, K., Ruwase, O., Kim, J.Y., Fowers, J., Strauss, K., Chung, E.S.: Accelerating deep convolutional neural networks using specialized hardware. *Microsoft Res. Whitepaper* 2(11) (2015)
7. Iandola, F.N., Han, S., Moskewicz, M.W., Ashraf, K., Dally, W.J., Keutzer, K.: SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and <0.5 MB model size. *arXiv preprint* (2016). [arXiv:1602.07360](https://arxiv.org/abs/1602.07360)

8. Qiu, J., Wang, J., Yao, S., Guo, K., Li, B., Zhou, E., Yu, J., Tang, T., Xu, N., Song, S., Wang, Y.: Going deeper with embedded FPGA platform for convolutional neural network. In: Proceedings of the 2016 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays, pp. 26–35. ACM, February 2016
9. Gschwend, D.: Zynqnet: an FPGA-accelerated embedded convolutional neural network. Masters thesis, Swiss Federal Institute of Technology Zurich (ETH-Zurich) (2016)
10. Gokhale, V., Jin, J., Dundar, A., Martini, B., Culurciello, E.: A 240 G-ops/s mobile coprocessor for deep neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, pp. 682–687 (2014)
11. Ma, Y., Cao, Y., Vrudhula, S., Seo, J.S.: Optimizing loop operation and dataflow in FPGA acceleration of deep convolutional neural networks. In: Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays, pp. 45–54. ACM, February 2017
12. Iandola, F.: SqueezeNet/SqueezeNet_v1.1 at master. DeepScale/SqueezeNet (2017). https://github.com/DeepScale/SqueezeNet/tree/master/SqueezeNet_v1.1
13. Xilinx Inc.: High-Level Synthesis. Vivado Design Suite User Guide. UG902 (2017). https://www.xilinx.com/support/documentation/sw_manuals/xilinx2017_2/ug902-vivado-high-level-synthesis.pdf
14. Vranesic, Z., Brown, S.: Fundamentals of Digital Logic with Verilog Design, 3rd edn. McGraw-Hill Education, New York (2014)
15. Ali, K.M.A., Ben Atitallah, R., Fakhfakh, N., Dekeyser, J.-L.: Exploring HLS optimizations for efficient stereo matching hardware implementation. In: Wong, S., Beck, A.C., Bertels, K., Carro, L. (eds.) ARC 2017. LNCS, vol. 10216, pp. 168–176. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-56258-2_15
16. Gysel, P., Motamedi, M., Ghiasi, S.: Hardware-oriented approximation of convolutional neural networks. arXiv preprint (2016). [arXiv:1604.03168](https://arxiv.org/abs/1604.03168)