

## Enhancing the Performance of Convolutional Neural Networks on FPGAs

Syed Roshan<sup>1</sup>, Amarendra Jadda<sup>2</sup>

<sup>1</sup>M.Tech., Scholar, <sup>2</sup>Associate Professor

Department of Electronics and Communication Engineering.

Audisankara College of Engineering & Technology, Gudur, Andhra Pradesh, India

### Abstract:

We presents a novel strategy to twofold the calculation pace of convolutional neural network (CNN) quickening agents by pressing two duplicate and-gather (MAC) tasks into one DSP square of off-the-rack FPGAs (called Double MAC). While an overall SIMD MAC utilizing a solitary DSP block appears to be unthinkable, our answer is custom-made for the sort of MAC activities needed for a convolution layer. Our fundamental assessment shows that not exclusively can our Double MAC approach increment the calculation throughput of a CNN layer by twice with basically a similar asset use, the network level execution can likewise be improved by 14~84% over an exceptionally enhanced cutting edge quickening agent arrangement relying upon the CNN hyper-boundaries.

### Introduction

An interesting choice accessible just to equipment quickening agents for convolutional neural networks (CNNs) [1]–[3] is the adaptability in number-crunching exactness. Specifically, FPGAs have the adaptability of picking whatever exactness adequate for the objective CNN application. Practically speaking, notwithstanding, FPGA mappings commonly depend on DSP blocks for most elevated productivity, and DSP blocks accompany a fixed exactness in particular. For example, a DSP48E1 block in Xilinx FPGAs can play out a 25x18-piece duplication followed by a 48-piece gathering. Regardless, with some change we can pack two diminished piece width augmentations into one DSP block, basically multiplying the calculation pace of CNN calculations than is right now conceivable

with the best FPGA planning procedures today. This paper is about how to turn a conventional DSP square of an off-the rack FPGA into a Double MAC, or a 2-way SIMD (Single-Instruction Multiple-Data) MAC (Multiply-and-Accumulate) unit, that can convey 4 operations/cycle by performing two duplicate and-add activities at the same time with decreased bit width. Despite the fact that we assess our procedure for a Xilinx Virtex-7 FPGA in particular, our technique is conventional and relevant to different FPGAs with comparable equipment DSP blocks. Our method doesn't need any adjustment in the FPGA texture itself. Acknowledging SIMD expansion on a FPGA is inconsequential and as of now upheld. A SIMD MAC additionally is insignificant whenever actualized with LUTs (Look-Up Tables). Anyway the test is the manner by which to acknowledge SIMD

duplication on a DSP square of a FPGA. Utilizing DSP blocks is basic, since most CNN usage on FPGAs depend on DSP blocks for MAC tasks [1], and thusly having the option to perform two MACs with one DSP block basically implies free 2X improvement in calculation throughput. Further, however it is conceivable to build throughput by different methods (e.g., actualizing extra MACs with LUTs [2]), our technique is symmetrical to them. Without local SIMD augmentation uphold on DSP blocks, we should make virtual SIMD paths inside one DSP block, ensuring that information on the paths don't slam into one another. We should likewise deal with sign pieces and flood recognition, all inside one DSP block. In addition to the fact that this is a long way from direct, a straightforward examination uncovers that an overall SIMD multiplier on one DSP block is almost incomprehensible. In this paper we characterize a unique class of SIMD multipliers, customized for the sort of MAC tasks found in convolutional layers of CNNs. In particular, our Double MAC necessitates that the increases of a SIMD duplicate offer a typical operand, viz.,  $A \times C$  and  $B \times C$ , and that the basic operand,  $C$ , be an unsigned whole number.

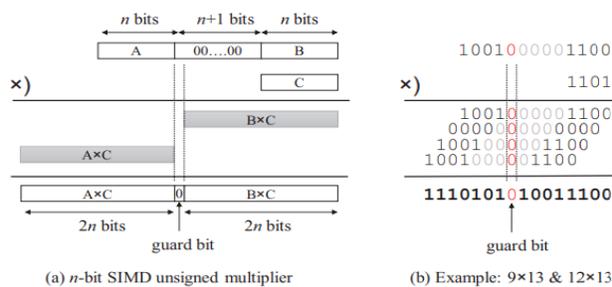


Fig. 1. How 2-way SIMD multiplication works (unsigned case).

We accept that lone testing, instead of preparing, of a CNN is done on a FPGA. We

likewise accept that the CNN quickening agent quickens convolution layers just, which represent by far most of calculation. We show that our Double MAC can twofold the calculation throughput of a MAC cluster when given a similar number of DSP blocks, when contrasted with the past best in class [1]. On the network level, i.e., with all convolution layers consolidated, our technique can produce execution upgrades that range, contingent upon the CNN hyper-boundaries, from 14% to over 80% over the exceptionally streamlined best in class quickening agents, intended for two enormous, genuine CNNs [4], [5], without relinquishing the yield quality by over 1% (top-5 precision) in the wake of applying our scaling plan.

### Literature Survey

#### Boosting Convolutional Neural Networks Performance Based on FPGA Accelerator

Convolutional Neural Network (CNN) has been widely utilized for picture acknowledgment because of its incredible precision. This precision is accomplished through imitating the optic nerves conduct in living individuals. The rapid advancement of the momentum applications got from profound learning calculations has extra improved exploration and improvements. All the more explicitly, a few profound CNN quickening agents have been moved toward FPGA-based stage, because of its quick advancement round, reconfigure-capacity, and elite. The FPGA is amazingly quicker than the CPU since it dependent on equal component, just as, devours low energy. This paper utilizes the FPGA in setting up the CNN engineering of type

VGG16 model. The FPGA explains the convolutional calculations for quickening the calculation time by 11%, without losing a lot of constancy when utilizing 16-digit fixed-point information design instead of 32-bit skimming point information design.

### **Boosting the Performance of CNN Accelerators with Dynamic Fine-Grained Channel Gating**

This paper proposes another fine-grained dynamic pruning method for CNN deduction, named channel gating, and presents a quickening agent engineering that can successfully misuse the dynamic sparsity. Naturally, channel gating distinguishes the districts in the component guide of each CNN layer that offer less to the characterization result and turns off a subset of channels for registering the initiations in these less significant areas. Dissimilar to static network pruning, which eliminates excess loads or neurons preceding surmising, channel gating misuses dynamic sparsity explicit to each contribution at run time and in an organized way. To expand register reserve funds while limiting precision misfortune, channel gating learns the gating edges along with loads consequently through preparing. Trial results show that the proposed approach can essentially accelerate best in class networks with a minor precision misfortune, and empower a compromise among execution and exactness. This paper likewise shows that channel gating can be upheld with a little arrangement of expansions to a CNN quickening agent, and executes a model for quantized ResNet-18 models. The quickening agent shows a normal accelerate

of  $2.3\times$  for ImageNet when the hypothetical FLOP decrease is  $2.8\times$ , demonstrating that the equipment can viably abuse the dynamic sparsity uncovered by channel gating.

### **Double MAC on a DSP: Boosting the Performance of Convolutional Neural Networks on FPGAs**

Profound learning, for example, Convolutional Neural Networks (CNNs) are a significant remaining burden progressively requesting superior equipment increasing speed. One separating highlight of profound learning remaining burden is that it is innately versatile to little mathematical blunders and functions admirably with low accuracy equipment. In this manner we propose a novel technique, called Double MAC, to hypothetically twofold the calculation pace of CNN quickening agents by pressing two duplicate and-amass (MAC) activities into one DSP square of off-the-rack FPGAs. There are a few specialized difficulties, which we defeat by misusing the method of activity in the CNN quickening agent. We have approved our technique through FPGA union and Verilog reproduction, and assessed our strategy by applying it to the condition of-the-art CNN quickening agent. We find that our Double MAC approach can build the calculation throughput of a CNN layer by twice. On the network level (all convolution layers consolidated), the exhibition improvement differs relying upon the CNN application and FPGA size, from 14% to over 80% over an exceptionally advanced cutting edge quickening agent arrangement, without

relinquishing the yield quality fundamentally.

### Related Work

Early work on CNN quickening utilized drifting point precisions (e.g., [2]), however late plans utilize fixed-point precisions (e.g., [4], [7]). For deduction instead of preparing, even 8-cycle fixed-point is demonstrated to be sufficient for some huge CNNs [17]. As of late, diminished exactness CNN models have become a subject of extraordinary interest because of their execution neighborly nature, for example, low force and ease. For instance, in the IBM TrueNorth processor [18], weight boundaries are allowed to have one of four potential qualities. Creators in [19] presents a learning technique for TrueNorth. Scientists are investigating even binarized neural networks (e.g., [20]) where weight boundaries are either 1, 0, or - 1. Creators of [21] endeavor to break down the effect of quantization in CNNs. A somewhat unique methodology is to utilize surmised multipliers [22], [23]. Then again, our CNN quickening agent is not quite the same as inexact figuring based ones (e.g., [23]–[26]). Inexact figuring based quickening agents exploit the versatility of AI calculations by utilizing surmised yet more energy-proficient useful units, for example, rough adders and estimated multipliers, and subsequently gain in cost and energy effectiveness. Our MAC exhibit produces accurate outcomes for the given info, and there is no truncation or adjusting presented by our MAC cluster, which is the reason we can utilize Caffe [27] to reenact our CNN

quickening agent. There are a few methodologies that actualize FPGA-based SIMD processors [28]–[31]. A large portion of those processors comprise of a variety of preparing components (PEs) that works in a SIMD style, consequently unique in relation to our meaning of SIMD, which is truly multi-word guidance. The work in [31] proposes a 2D convolution processor which executes a SIMD convolver. A 16-bit convolution activity is part into two at the same time 8-cycle convolutions by working on the operands subwords. Our methodology is totally extraordinary, since we execute two concurrent procedure on a solitary process unit. Our plan upholds two simultaneous marked unsigned duplication with only one marked multiplier. Some amendment steps are needed since the second path of our Double MAC can't perceive marked activity. A comparative technique to accomplish right marked increase with unsigned multiplier is referenced.

### Methodology

#### SIMD Multiplication of Unsigned Numbers

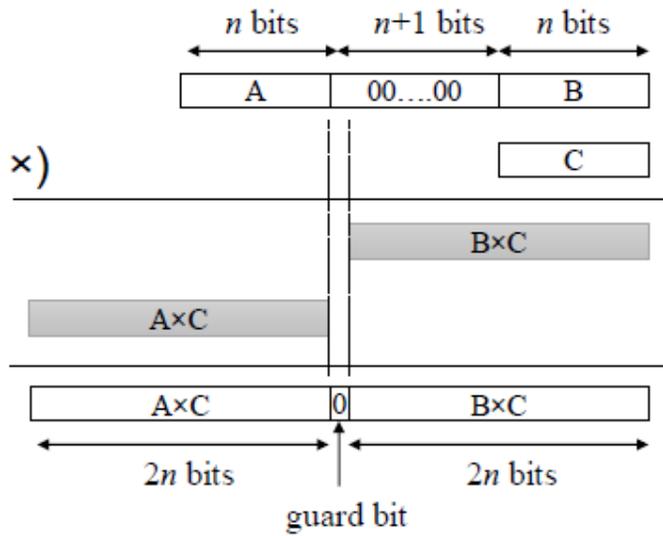
Let us initially consider the situation where all the operands are unsigned. a solitary multiplier can perform two unsigned duplications with a typical operand all the while, i.e., A\_C and B\_C. For this to work, two conditions must be fulfilled. Leave n alone the width of every operand. To begin with, the yield register must be in any event 4n-chomped wide. Second, the two operands in one of the sources of info must be isolated by at any rate n bits.

For amassing to work without flood, we need in any event one watchman chomped as represented in the figure. During the aggregation a similar gatekeeper spot can be utilized to distinguish any complete of the lower  $2n$  bits. On location, the convey bit is cleared quickly and the quantity of complete functions is tallied independently utilizing a little counter. The estimation of the counter is added independently once all aggregation is done. Hence to perform two  $n$ -bit duplications in a SIMD design we need one  $(3n + 1)$   $n$ -digit multiplier. For instance, with the  $25 \times 18$ -piece multiplier of a DSP48E in Xilinx FPGAs,  $n$  can be all things considered 8.

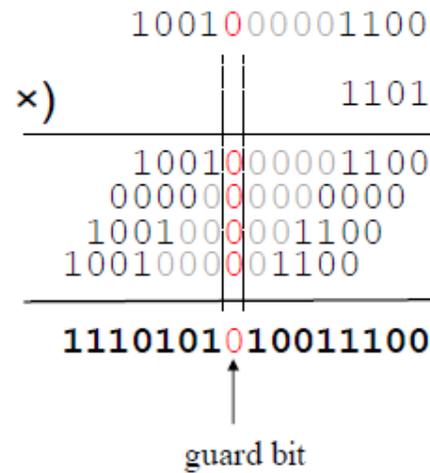
bit unsigned integer. Recalling  $[-2^n, 2^n] = 2n+1$  bits, the product  $B \cdot C$  can be computed as follows, where  $\hat{B}$  represents the value of  $B$  interpreted as an unsigned number.

$$\begin{aligned}
 B \times C &= (-b_{n-1} \cdot 2^{n-1} + \sum_{i=0}^{n-2} b_i \cdot 2^i) \cdot C \\
 &= (b_{n-1} \cdot 2^{n-1} + \sum_{i=0}^{n-2} b_i \cdot 2^i) \cdot C - b_{n-1} \cdot 2^n \cdot C \\
 &= \hat{B} \cdot C - b_{n-1} \cdot 2^n \cdot C
 \end{aligned}$$

In other words, we can compute  $B \cdot C$  by



(a)  $n$ -bit SIMD unsigned multiplier



(b) Example:  $9 \times 13$  &  $12 \times 13$

How 2-way SIMD multiplication works (unsigned case).

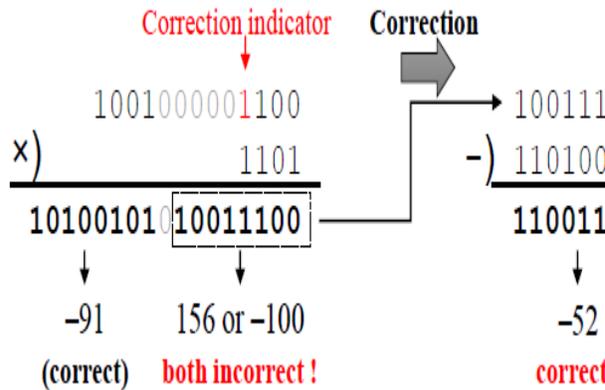
**SIMD Signed-Unsigned Multiplication**

Let us consider how to perform a signed-unsigned multiplication using an unsigned multiplier. Let  $B = b_n \dots b_1 \dots b_0$  be an  $n$ -bit signed integer, and  $C$  an  $n$ -

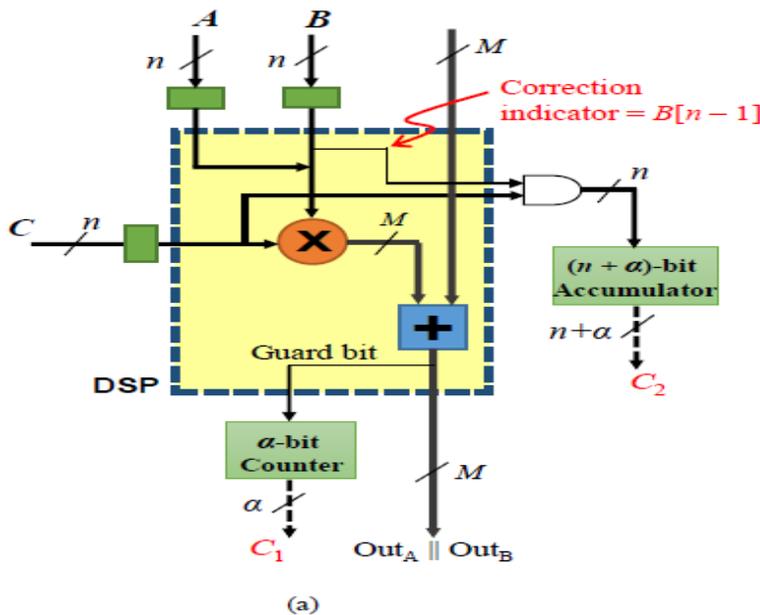
performing an unsigned multiplication on  $B$  and  $C$  followed by a subtraction of  $n$ -bit left-shifted  $C$  if  $B$  is negative.

Extending this to 2-way SIMD multiplication is straightforward. We can perform two  $n$ -bit signed-unsigned multiplications in a SIMD fashion by performing one  $(3n + 1)$   $n$ -bit unsigned multiplication, followed by at most two

subtractions. However we can reduce the number of subtractions to one by performing a signed multiplication for the  $(3n + 1) \_ n$ -bit multiplication. In this case the upper  $2n$ -bit (i.e.,  $A \_ C$ )



Example SIMD execution:



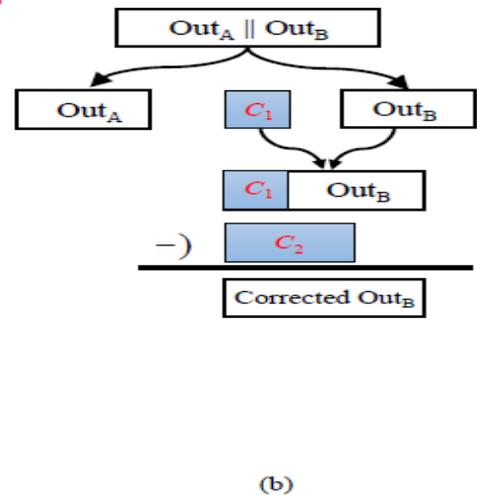
(a)Datapath of our Double MAC architecture, where  $n = 8$  and  $M = 48$ , and (b) how it works.

It is as of now right (there is no requirement for remedy), and just the lower  $2n$ -chomped

needs an adjustment if B is negative, as exhibited by the model in Fig. 2.

The amendment term could be added set up on the off chance that we do only one duplicate add activity. For collection of countless increase results, be that as it may, we should defer doing the adjustments until all aggregation is done on the grounds that we have just one watchman bit.1 Instead the remedy terms are gathered independently in a little collector, whose worth is later deducted from the primary gatherer. We decide the size of additional collectors in the following segment.

**Accumulation and Our Double MAC**



**Architecture**

Above Fig. represents our Double MAC engineering. One DSP square can execute both a 2-way SIMD multiplier and a 2-way SIMD snake. Both the multiplier and the

snake have an adjustment yield signal, which is appeared as slim bolts in the figure. The amendment yields are collected into a different register or counter to be utilized later for conclusive change. In the last change two terms should be added/deducted. The snake rectification term C1, which originates from the complete counter, has no cover with the fundamental gatherer yield, and accordingly can be simply connected. Hence we need just a single expansion—the deduction of the duplicate remedy term C2 from the link result.

The width  $\_$  of the complete counter is resolved from the quantity of qualities collected,  $V$ , as follows:  $\_ \log_2 V$ . On account of our benchmark CNN quickening agent (see Section II-A),  $V = (K2N=TN)$ , accordingly  $\_$  ought to be no under  $\log_2(K2N=TN)$  for each layer of the CNN, where  $K$  is the size of convolution channel in one measurement. Additionally the width of the collector for multiplier rectification terms is  $n + \_$ . The post-amassing change should be possible effectively by utilizing an additional viper of  $(n + \_)$ - bit width that lives outside the DSP block. There is no runtime overhead because of this change, since the change is done at the same time while new qualities are stacked.

## Results

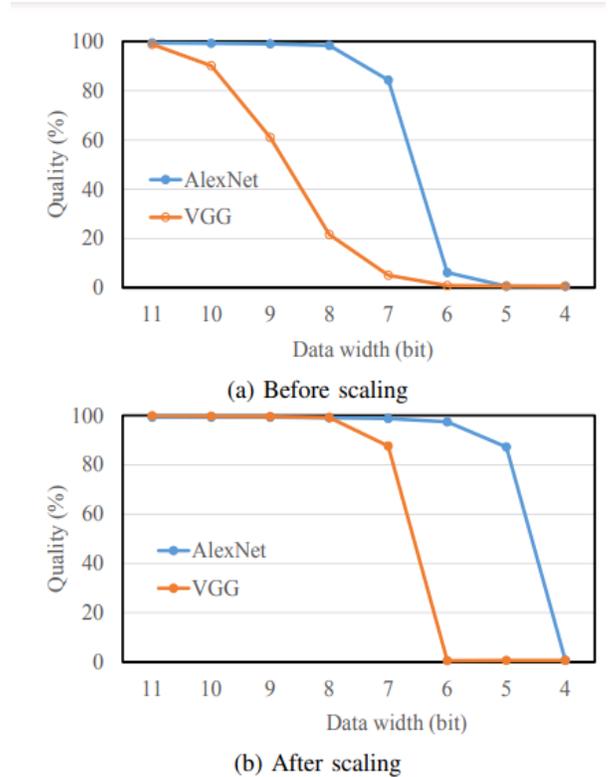
Table looks at the asset necessities of one MAC unit for various precisions. The information for coasting point and 32-cycle

fixed-point are from [2], appeared here for examination. We have actualized our Double MAC exhibit in Verilog RTL and approved its usefulness utilizing Vivado test system. The MAC cluster is incorporated utilizing Vivado 2015.2 focusing on the Xilinx Virtex7 485T FPGA, which has 2,800 DSP blocks. The MAC exhibit has two key boundaries, TM and TN, likewise called tile boundaries, that sway the throughput of the quickening agent. Their ideal qualities, appeared in Table II, are discovered utilizing comprehensive quest for AlexNet. True to form, with a similar degree of DSP usage our SIMD engineering can execute a MAC exhibit twice huge than the non-SIMD variant. The table likewise reports the territory and most extreme recurrence of the integrated circuit. Despite the fact that our Double MAC cluster burns-through more LUT and FF, which is because of the revision circuit, the extra asset utilization isn't high. Then again, having the option to transform the additional assets into execution improvement can be a bit of leeway, particularly when the assets are regularly underutilized in CNN quickening agents except if a gliding point information type is utilized.

Data type	DSP	LUT	FF
Floating-point 32-bit [2]	5	349	355
Fixed-point 32-bit [2]	4	0	0
Fixed-point 16-bit	1	0	0
Fixed-point 8-bit	1	0	0
Fixed-point 8-bit, Double MAC	0.5	11	12

## Performance Scalability

The underneath figure shows how the exhibition of two CNNs on the two CNN quickening agents (8-bit fixed-point benchmark versus our Double-MAC cluster) differs as the quantity of DSPs is changed. Different conditions, for example, memory data transfer capacity are kept unaltered. The quantity of-DSPs esteem is taken from the DSP checks of existing Xilinx FPGA gadgets [33]. We initially see that distinctive CNNs have totally different execution scaling designs. VGG's presentation on the gauge cluster scales straightly as the quantity of DSPs increments while AlexNet's exhibition is almost soaked at around 400 GOPS. Subsequently when our Double MAC is applied, the presentation improvement by our plan is definitely extraordinary between the two applications. While our Double MAC accomplishes critical speedup on VGG (84%), its presentation with AlexNet isn't as amazing (14%). In any case, we watch a typical example that our Double MAC cluster can accomplish the presentation level of the standard 0 2 4 8 16 32 0 20 40 60 80 100 5 10 15 20 LUT Usage (%) Performance (GOPS) 8b-fixed 8b-fixed-D (a) LUT use 0 2 4 8 16 32 0 5 10 15 20 25 5 10 15 20 FF Usage (%) Performance (GOPS) 8b-fixed 8b-fixed-D (b) FF utilization For the situation of 8b-Fix, the number on the left half of the imprints speaks to the quantity of additional columns of MACs. 0 500 1000 1500 2000 2500 0 500 1000 1500 2000 2500 3000 Performance (GOPS) Number of DSPs AlexNet, 8b-fixed AlexNet, 8b-fixed-D VGG, 8b-fixed VGG, 8b-fixed-D Fig. 6: Performance versus Number of DSPs. at double the quantity of DSPs. As such our



Double MAC cluster can effectively twofold the successful number of DSPs for both CNNs, in any event, when decided by the calculation rate improvement.

### CONCLUSION

We introduced how to expand the calculation pace of CNN quickening agents on FPGAs by pressing various MAC tasks into one DSP squares of off-the-rack FPGAs. By misusing the setting where MAC activities are utilized, our strategy can find some kind of harmony among ease of use and execution overhead. We have approved our proposed engineering through Verilog recreation and FPGA union, and assessed it utilizing a cutting edge CNN quickening agent, which shows that our Double MAC can build calculation throughput of a CNN layer frequently by

twice, and accomplish significant execution enhancements for the network level going from 14% to over 80% over a previously improved quickening agent arrangement relying upon the hyper-boundaries of the CNN and the FPGA size. All these are managed without relinquishing the yield quality essentially. This improvement in execution can legitimately convert into energy sparing, which can make quickening agent arrangements much additionally engaging when contrasted with GP-GPU arrangements.

## References

1. M. Courbariaux, Y. Bengio and J.-P. David, "BinaryConnect: Training Deep Neural Networks with binary weights during propagations", NIPS, vol. abs/1511.00363, 2015.
2. Sugil Lee; Daewoo Kim; Dong Nguyen; Jongeun Lee "Double MAC on a DSP: Boosting the Performance of Convolutional Neural Networks on FPGAs"
3. W. Sung, S. Shin and K. Hwang, Resiliency of Deep Neural Networks under Quantization, CoRR, vol. abs/1511.06488, 2015.
4. D.D Lin, S.S. Talathi and V.S. Annapureddy, "Fixed Point Quantization of Deep Convolutional Networks", ICML, vol. abs/1511.06393, 2016.
5. A. Sajid, H. Kyuyeon and W. Sung, "Fixed point optimization of deep convolutional neural networks for object recognition", IEEE International Conference on Acoustic Speech and Signal Processing (ICASSP), 2015.
6. P. Gysel, Ristretto: Hardware-Oriented Approximation of Convolutional Neural Networks, CoRR, vol. abs/1605.06402, 2016.
7. Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, et al., Caffe: Convolutional Architecture for Fast Feature Embedding, CoRR, vol. abs/1408.5093, 2014.
8. C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, et al., "Going Deeper with Convolutions", CVPR, 2015.
9. S. Han, H. Mao and W.J. Dally, "Deep Compression: Compressing Deep Neural Networks with Pruning Trained Quantization and Huffman Coding", ICLR, 2016.
10. S. Han, X. Liu, H. Mao, J. Pu, A. Pedram, M.A. Horowitz, et al., "EIE: Efficient Inference Engine on Compressed Deep Neural Network", ISCA, 2016.
11. S. Shin, K. Hwang and W. Sung, "Fixed Point Performance Analysis of Recurrent Neural Networks", ICASSP, vol. abs/1512.01322, 2016.
12. O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, et al., ImageNet Large Scale Visual Recognition Challenge, CoRR, vol. abs/1409.0575, 2014.
13. A. Krizhevsky and G. Hinton, Learning multiple layers of features from tiny images, 2009.

14. Randy Yates, Practical considerations in fixed-point fir filter implementations, 2000.
15. K. He, X. Zhang, S. Ren and J. Sun, "Deep residual learning for image recognition", CVPR, 2016.