



HASH TRICK BASED DEEP NEURAL NETWORK IMPLEMENTATION USING OPTIMIZED STOCHASTIC MULTIPLIER

¹K.Saranya, ²Chitravalavan

¹P.G.Scholar, ²Professor

¹Electronics and Communication Engineering,

¹A.V.C.College of Engineering, Mayiladuthurai, India.

Abstract—Deep neural networks (NN) have shown a significant promise in difficult tasks like image classification or speech recognition. Even well-optimized hardware implementations of digital NNs show significant power consumption. It is mainly due to non-uniform pipeline structures and inherent redundancy of numerous arithmetic operations that have to be performed to produce each single output vector. This paper provides a methodology for the design of well-optimized power-efficient NNs with a uniform structure suitable for hardware implementation with hash tricking. An error resilience analysis was performed in order to determine key constraints for the design of approximate multipliers that are employed in the resulting structure of NN. By means of a search based approximation method, approximate multipliers showing desired tradeoffs between the accuracy and implementation cost were created. Significant improvement in area efficiency was obtained in both cases with respect to regular NNs.

Keywords- Deep Neural Networks, stochastic computing, hashed neural nets, VLSI.

I. INTRODUCTION

Deep Neural Networks (DNNs) are amazingly ground-breaking AI models that accomplish fantastic execution on troublesome issues, for example, discourse acknowledgment and visual article acknowledgment. DNNs are ground-breaking since they can perform discretionary equal calculation for an unassuming number of steps. An astonishing illustration of the intensity of DNNs is their capacity to sort N-digit numbers utilizing just 2 shrouded layers of quadratic size. Thus, while neural organizations are identified with ordinary factual models, they get familiar with a complicated calculation. Besides, huge DNNs can be prepared with managed back proliferation at whatever point the named preparing set has enough data to indicate the organization's boundaries. Subsequently, if there exists a boundary setting of an enormous DNN that accomplishes great outcomes for instance, since people can tackle the undertaking quickly, directed back engendering will discover these boundaries and take care of the issue. Application zones of DNNs incorporate framework recognizable proof and control, game playing and dynamic, design acknowledgment, radar frameworks, face ID, object acknowledgment, arrangement acknowledgment, clinical determination, monetary applications, data mining perception and email spam sifting. Genuine applications the errands to which counterfeit neural organizations are applied will in general fall inside the accompanying general classifications: Function guess, or relapse investigation, including time arrangement expectation and displaying. Grouping, including example and succession acknowledgment, curiosity identification and consecutive dynamic information preparing that involves sifting, bunching, dazzle signal division and pressure.

Regardless of their adaptability and force, DNNs must be applied to issues whose information sources and targets can be reasonably encoded with vectors of fixed dimensionality. It is a huge constraint, since numerous significant issues are best communicated with arrangements whose lengths are not known from the earlier. For instance, discourse acknowledgment and machine interpretation are successive issues. Moreover, question noting can likewise be viewed as planning a succession of words representing the inquiry to an arrangement of words speaking to the appropriate response. It is thusly evident that an area free strategy that figures out how to plan arrangements to successions would be valuable. The term Neural Network was

generally used to allude to an organization or circuit of natural neurons. The cutting edge utilization of the term regularly alludes to fake neural organizations, which are made out of fake neurons or hubs.

Biological neural networks (BNN) are comprised of genuine organic neurons that are associated or practically related in the fringe sensory system or the focal sensory system. In the field of neuroscience, they are regularly distinguished as gatherings of neurons that play out a particular physiological capacity in research facility examination. Artificial neural networks (ANN) are comprised of interconnecting fake neurons programming develops that copy the properties of organic neurons. ANN may either be utilized to pick up a comprehension of natural neural organizations, or for taking care of computerized reasoning issues without fundamentally making a model of a genuine natural framework. The genuine, natural sensory system is profoundly intricate and incorporates a few highlights that may appear to be pointless dependent on a comprehension of counterfeit networks. The genuine structure of the organization and the techniques used to set the interconnection loads change starting with one neural system then onto the next, each with its favorable circumstances and hindrances. Organizations can proliferate data in one way just, or they can ricochet to and fro until self-initiation at a hub happens and the organization chooses a last state. The capacity for bi-directional progression of contributions between neurons/hubs was delivered with the Hopfield's organization, and specialization of these hub layers for explicit designs was presented through the main mixture organization. The utility of fake neural organization models lies in the way that they can be utilized to deduce a capacity from perceptions and furthermore to utilize it. This is especially valuable in applications where the unpredictability of the information or undertaking makes the plan of such a capacity by hand illogical.

II. RELATED WORK

Two plans of less error fixed-width sign equal multipliers and two's-supplement equal multipliers for computerized signal handling applications are introduced. Given two n -cycle inputs, the fixed-width multipliers create n -bit rather than $2n$ -piece items with low error, yet utilize just about a large portion of the territory and less defer when contrasted and a standard equal multiplier. Among those, the convey creating circuits are planned, separately, to make the items produced all the more precisely and rapidly[1]. Applying a similar methodology, a low-blunder diminished width multiplier with yield bit-width among n and $2n$ has likewise been planned. Exploratory outcomes show that the proposed fixed-width and diminished width multipliers have lower blunder than any remaining fixed-width multipliers are as yet savvy. Because of these properties, they are entirely reasonable for use in numerous media and advanced sign preparing applications, for example, computerized sifting, number-crunching coding, wavelet change, reverberation crossing out, and so on Pooling layers are down samplers of 2-D pictures. Max pooling layers give a spatial greatest capacity, which separates an information picture into little subtiles of a given window size and afterward replaces these with the most extreme incentive in the subtile [2]. A normal pooling layer is comparable; nonetheless, it finds the normal in the subtile as opposed to the greatest. Profound neural organization (DNN) preparing is an iterative cycle which has a feedforward way to register the yield and a backpropagation way for realizing, which includes figuring inclinations and update the organization loads. Preparing of low exactness networks commonly includes keeping a bunch of single-or twofold accuracy skimming point loads W which are quantized to a portrayal before induction [3].

As the quantization capacities utilized are piecewise and steady, the inclinations of quantized loads can be determined and applied to refresh their relating full-exactness loads [4]. A quantization work which diminishes confuse in forward and in reverse ways is pivotal for high precision. To additionally improve exactness, a choice to low-accuracy networks is to utilize a weight-sharing methodology [5]. Weight sharing includes picking a limited arrangement of full-accuracy loads ordered by a codebook. Commonly, these loads are picked to coordinate the ideal dissemination to decrease data misfortune, not at all like conventional fixed-point quantization where loads are consistently appropriated. Keeping the quantity of various loads in the codebook little diminishes the word size of the records prompting a little memory impression. Nonetheless, weight sharing is typically not applied in FPGA usage as the weight planning measure presents extra postponements in the basic way of the circuit and requires additional equipment[6].

Besides, higher exactness math units additionally burn-through more territory. For the proposed RCCM, an understood weight sharing is used, lessening coefficient memory without requiring any planning equipment. In the interim, our RCCMs are advanced for FPGA equipment, implying that they burn-through fewer regions than fixed-point counterparts. Numerous quantization techniques for neural organizations have been investigated with the point of accomplishing proficient surmising in equipment. A proficient method of preparing networks with various forward and in reverse capacities was presented in [7]. This prompted new inferences of uniform quantization capacities for low-exactness neural organizations in [8]. Learning symmetric quantization (SYQ) further investigated the significance of introductions and planning a quantization work which lessens the forward and confounds. They accomplished cutting edge exactnesses under low-accuracy loads and initiations. This propelled the deduction of the dissemination coordinating instatement technique for effective quantization. Viable nonuniform quantization structures were likewise investigated as log portrayals [9]. This structure can likewise process multiplier-less duplicate and gather (MACs) activities; in any case, the dispersion of the portrayals is confined to the log space. There have been a few quickening agent engineering plans for low-exactness CNNs with uniform quantization number-crunching.

Other FPGA models have been actualized to use the profoundly manageable nature of CNNs which oblige weight boundaries to be just paired or ternary portrayals [10]. With limitations in the productivity of both programming and equipment usage of neural organizations, programming equipment code sign is viewed as a powerful way to deal with accomplish ideal execution [11]. A strategy for planning a quantization work for both expanding precision of binarized CNNs while keeping up proficient multiplier less equipment was proposed in [12]. Furthermore, an effective long transient memory (LSTM) usage in [13] used burden balance-mindful pruning to accomplish both organization pressure and high equipment use. Additionally, preparing profoundly inadequate ternary organizations and planning proficient CNN equipment for misuse was portrayed in [14]. Apparently, AddNet is the primary quantization plot which inserts reconfigurability straightforwardly into its portrayals.

III. PROPOSED SYSTEM

The goal of the proposed system is to develop an optimal multiplier to give an actually smaller fast and low force utilization unit. Being a center part of math preparing unit multipliers are in very popularity on its speed and low force utilization. To diminish critical force utilization of multiplier plan it is a decent course to diminish number of activities along these lines lessening a unique force which is a significant piece of absolute force dissemination. In the past significant exertion were placed into planning multiplier in VLSI in this bearing.

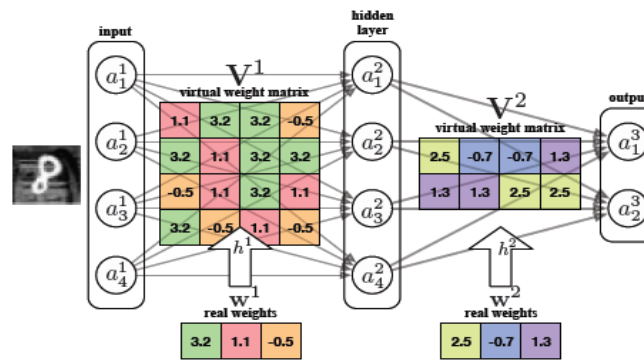


Fig 1 Architecture of Hashed Neural Network

A multiplier is one of the key equipment blocks in most of the digital signal processing frameworks. Common DSP applications where a multiplier plays an significant job incorporate advanced separating, computerized correspondences and examination. Numerous current DSP applications are focused at convenient, battery-worked frameworks, so that power dispersal gets one of the essential plan imperatives. Since multipliers are somewhat intricate circuits and should commonly work at a high framework clock rate, lessening the deferral of a multiplier is a fundamental piece of fulfilling the general plan. The hashed neural network is shown in Fig 1.

A. Stochastic Computing based Multiplier

Stochastic Computing (SC) measures information as (pseudo) arbitrary arrangements of 0 s and 1 s . These spot streams are alluded to as stochastic numbers (SNs) and are deciphered as probabilities. In its fundamental "unipolar" structure, a bitstream X of length n containing $n1$ 1 s indicates the likelihood $pX = n1/n$, I. e., the likelihood of a haphazardly picked piece of X being 1. A similar worth pX is spoken to by a wide range of touch streams, so this number arrangement is profoundly repetitive. For instance, 0010, 01000001 and 00001100 all indicate the equivalent number $1/4$. nipolar SNs straightforwardly rough numbers lying in the span $[0,1]$, however they can be scaled to different ranges, for example, the stretch $[-1, 1]$ utilized for marked numbers (the purported "bipolar" SNs) . The vital favorable position of SC is that math tasks can be performed by amazingly basic and standard rationale circuits utilizing any ideal equipment usage innovation. The above figure shows stochastic circuits for augmentation, expansion, and number configuration change. For instance, the two-information AND door plays out the augmentation $pX \times pY$ in n clock cycles, with the proviso that the info bit-streams X and Y should be factually free (uncorrelated) and adequately long to give worthy exactness. Expansion is executed by the two-path multiplexer in the scaled structure $0.5(pX + pY)$, which guarantees that the entirety consistently lies in the likelihood stretch $[0, 1]$.

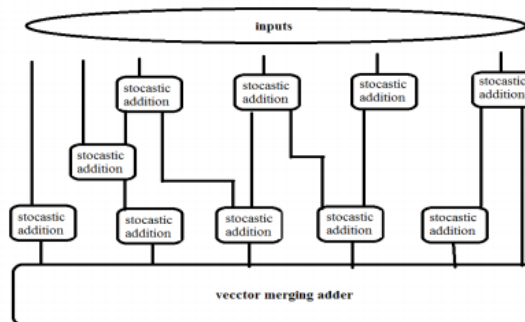


Fig 2 Block diagram of the Stochastic multiplier special adder

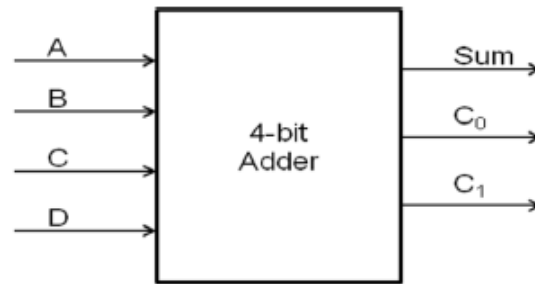


Fig 3 Architecture of 4 bit special adder

See that the scaling factor 0.5 is provided by a stochastic number R of steady worth 0.5, I. e., a simply arbitrary succession of 0 s and 1 s. The leftover circuits convert numbers between regular parallel and stochastic structures. They serve to interface standard and stochastic computerized circuits, and to re-randomize SNs that have gotten unduly corresponded. The (pseudo) irregular number generator is commonly executed by a straight criticism move register (LFSR). Number transformation circuits now and then establish the biggest piece of a SC framework. SC has a few highlights that limit its handiness, including long bitstreams and processing times, errors brought about by arbitrary cycle vacillations added to hidden layer) without increasing the actual number of parameters of the neural network. Connections are randomly grouped into three categories per layer and their weights are shown in the virtual weight matrices $V1$ and $V2$. Connections belonging to the same color share the same weight value, which are stored in $w1$ and $w2$, respectively. Overall, the entire network is compressed by a factor $1=4$, i.e. the 24 weights stored in the virtual matrices $V1$ and $V2$ are reduced to only six real values in $w1$ and $w2$. On data with four input dimensions and two output dimensions, a conventional neural network with six weights would be restricted to a single (trivial) hidden unit.

B. Multiplication Logic

$X1$ and $X2$ are the two independent bit streams of AND gate whose output is a bit stream Y , where $y = P(Y = 1) = P(X1 = 1 \text{ and } X2 = 1) = P(X1 = 1)P(X2 = 1) = x1x2$. The probability values of the two inputs $X1$ and $X2$ are multiplied and returned as the result of AND gate.

C. Scaled Addition

Consider a two-input multiplexer Suppose that its inputs are two independent stochastic bit streams $X1$ and $X2$ and its selecting input is a stochastic bit stream S . Its output is a bit stream Y , where $y = P(Y = 1) = P(S = 1)P(X1 = 1) + P(S = 0)P(X2 = 1) = sx1 + (1 - s)x2$. (Note that throughout the paper, multiplication and addition represent arithmetic operations, not Boolean AND and OR.) Thus, the multiplexer computes the scaled addition of the two input probability values

The benefit of the stochastic engineering as far as assets is that it endures blames effortlessly. Contrast a stochastic encoding with a standard double radix encoding, state with M pieces speaking to fragmentary qualities somewhere in the range of 0 and 1.0 Assume that the climate is boisterous; piece flips happen and these burden all the pieces with equivalent likelihood. With a double radix encoding, assume that the most huge piece of the information gets flipped. This causes an overall mistake of $2M-1/2M = 1/2$. Conversely, with a stochastic encoding, the information is spoken to as the fragmentary load on a piece stream of length $2M$. Accordingly, a solitary piece flip as it were changes the info esteem by $1/2M$, which is little in examination scaling necessities. Notwithstanding, in view of their intrinsically repetitive encoding design, blunders of the touch flip sort have little impact on the estimation of a SN. For instance, a solitary piece flip happening in a n -digit SN changes the yield an incentive by $1/n$, a generally little mistake whose noteworthiness reduces as n increments. Moreover, on the off chance that the blunders are bidirectional, I. e., if 0-to-1 and 1-to-0 digit flips are similarly likely, at that point the blunders will in general drop each other. This recommends that stochastic figuring can outflank paired in applications.

D. Hashed Nets

HashedNets, a novel variety of neural organizations with definitely decreased model sizes (and memory requests). We initially present our methodology as a method of irregular weight sharing across the organization connections and then depict how to encourage it with the hashing trick to dodge any extra memory overhead.

Random weight sharing

In a standard completely associated neural organization, there are $(n+1)_{n+1}$ weighted associations between a couple of layers, each with a comparing free boundary in the weight matrix V' . We expect a limited memory financial plan per layer, $K'_{(n+1)_{n+1}}$, that can't be surpassed. The obvious solution is to fit the neural organization inside spending plan by reducing the quantity of hubs n'_{n+1} in layers $'_{n+1}$ or by lessening the touch exactness of the weight lattices. Notwithstanding if K' is adequately small, both approaches essentially lessen the capacity of the neural network to sum up. All things being equal, we propose another option: we keep the size of V' immaculate but reduce its powerful memory impression through weight sharing. We just permit precisely K' various loads to occur within V' , which we store in a weight vector $w'_{2RK'}$. The loads inside w' are shared across different randomly chosen associations inside V' . We allude 0to the resulting matrix V' as virtual, as its size could be expanded (i.e. nodes are 52 added to concealed layer) without expanding the actual number of boundaries of the neural organization. Fig 1 shows a neural organization with one covered up layer, four input units and two yield units. Associations are randomly assembled into three classes for every layer also, their weights are appeared in the virtual weight grids $V1$ and $V2$. Associations having a place with a similar shading share the same weight esteem, which are put away in $w1$ and $w2$, separately.

Generally, the whole organization is compacted by a factor $1=4$, for example the 24 loads put away in the virtual matrices V1 and V2 are diminished to just six genuine qualities in w1 and w2. On information with four information measurements and two yield dimensions, a traditional neural organization with six weights would be confined to a solitary (insignificant) concealed unit.

IV. PERFORMANCE ANALYSIS

The performance of the proposed system is evaluated using the report generated by XILINK.

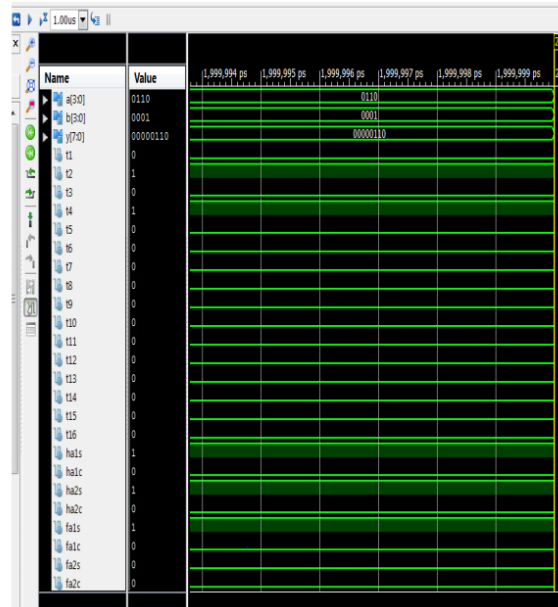


Fig 4 Multiplier Output

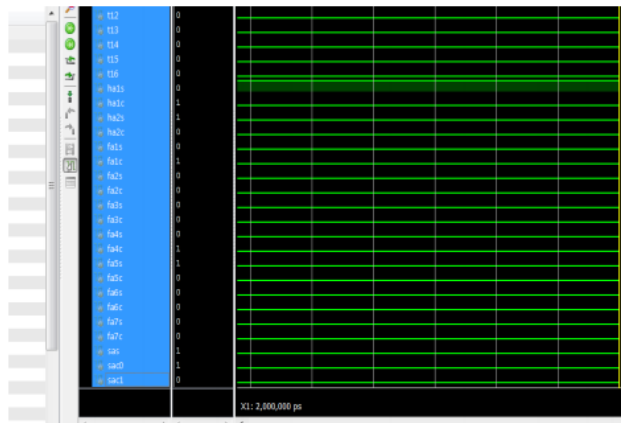


Fig 5 Adder Output

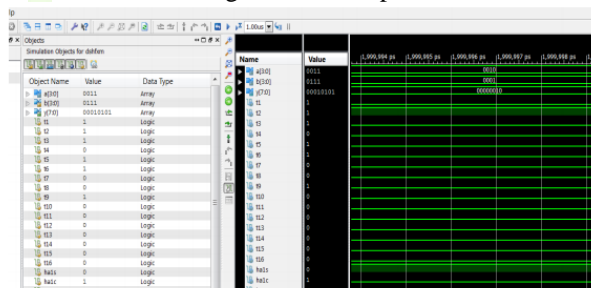


Fig 6 Stochastic Multiplication Output

Table 1 Parameters Used

s.no	Parameter	Existing	Proposed
1	Slices	78	68
2	LUT	143	125

The Figure given below is shown that there is a considerable reduction in time and area based on the implementation results which have been done by using Spartan-3 processor. The proposed algorithm significantly reduces area consumption when compared to the existing system.

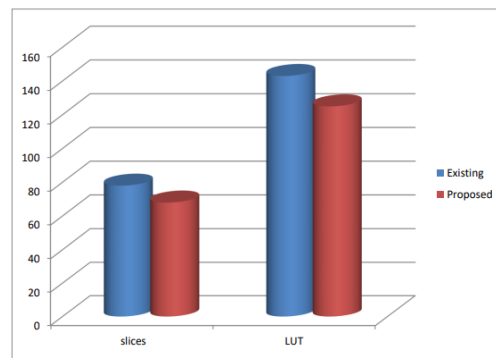


Fig 7 Comparison in the performance of proposed and existing system

V. CONCLUSION

Deep neural networks are state-of-the-art machine learning techniques for a wide range of applications. Many of these applications are error-resilient, where approximate computing is naturally used to achieve great energy savings with minor quality degradation. In this project, Approx ANN approximates the computation partial products in weight product stages, with the proposed solution, Approx ANN achieves energy benefit with less than 5% quality loss for various multiplication techniques used in our experiments. The proposed system for stochastic and the approximate multiplier using Artificial Neural Network achieved power savings of 72% and 38% respectively and, energy benefit and low quality loss for various multiplication techniques are used. In phase II, will stochastic computation of adder has been proposed and will be implemented in neural network addition stages.

VI. REFERENCES

- [1] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015. doi: 10.1038/nature14539.
- [2] C. Zhang, P. Li, G. Sun, Y. Guan, B. Xiao, and J. Cong, "Optimizing FPGA-based accelerator design for deep convolutional neural networks," in *Proc. ACM/SIGDA Int. Symp. Field-Program. Gate Arrays*, New York, NY, USA, Feb. 2015, pp. 161–170. doi: 10.1145/2684746.2689060.
- [3] Y. Umuroglu *et al.*, "FINN: A framework for fast, scalable binarized neural network inference," in *Proc. 2ACM/SIGDA Int. Symp. Field-Program. Gate Arrays*, New York, NY, USA, Feb. 2017, pp. 65–74. doi: 10.1145/3020078.3021744.
- [4] S. Han, H. Mao, and W. J. Dally, "Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding," in *Proc. 4th Int. Conf. Learn. Represent. (ICLR)*, San Juan, Puerto Rico, May 2016, pp. 1–14. [Online]. Available: <http://arxiv.org/abs/1510.00149>
- [5] A. G. Howard *et al.*, "MobileNets: Efficient convolutional neural networks for mobile vision applications," *CoRR*, 2017. [Online]. Available: <http://arxiv.org/abs/1704.04861>
- [6] M. Courbariaux and Y. Bengio, "Binarized neural networks: Training deep neural networks with weights and activations constrained to +1 or -1," *CoRR*, 2016. [Online]. Available: <http://arxiv.org/abs/1602.02830>
- [7] D. Miyashita, E. H. Lee, and B. Murmann, "Convolutional neural networks using logarithmic data representation," *CoRR*, 2016. [Online]. Available: <http://arxiv.org/abs/1603.01025>
- [8] D. D. Lin, S. S. Talathi, and V. S. Annapureddy, "Fixed point quantization of deep convolutional networks," in *Proc. 33rd Int. Conf. Int. Conf. Mach. Learn.*, vol. 48, Jun. 2016, pp. 2849–2858. [Online]. Available: <http://dl.acm.org/citation.cfm?id=3045390.3045690>
- [9] S. S. Demirsoy, A. G. Dempster, and I. Kale, "Design guidelines for reconfigurable multiplier blocks," in *Proc. Int. Symp. Circuits Syst.*, May 2003, p. 4.
- [10] R. I. Hartley, "Subexpression sharing in filters using canonic signed digit multipliers," *IEEE Trans. Circuits Syst. II, Analog Digit. Signal Process.*, vol. 43, no. 10, pp. 677–688, Oct. 1996.
- [11] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *CoRR*, 2015. [Online]. Available: <http://arxiv.org/abs/1512.03385>
- [12] Y. Bengio, N. Léonard, and A. C. Courville, "Estimating or propagating gradients through stochastic neurons for conditional computation," *CoRR*, 2013. [Online]. Available: <http://arxiv.org/abs/1308.3432>
- [13] W. Chen, J. T. Wilson, S. Tyree, K. Q. Weinberger, and Y. Chen, "Compressing neural networks with the hashing trick," in *Proc. 32nd Int. Conf. Int. Conf. Mach. Learn.*, vol. 37, Jul. 2015, pp. 2285–2294. [Online]. Available: <http://dl.acm.org/citation.cfm?id=3045118.3045361>
- [14] J. Wu, C. Leng, Y. Wang, Q. Hu, and J. Cheng, "Quantized convolutional neural networks for mobile devices," in *IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016, pp. 4820–4828.