

Enabling Capsule Networks at the Edge through Approximate Softmax and Squash Operations

Original

Enabling Capsule Networks at the Edge through Approximate Softmax and Squash Operations / Marchisio, Alberto; Bussolino, Beatrice; Salvati, Edoardo; Martina, Maurizio; Masera, Guido; Shafique, Muhammad. - ELETTRONICO. - (2022). (Intervento presentato al convegno ISLPED '22: ACM/IEEE International Symposium on Low Power Electronics and Design tenutosi a Boston (USA) nel 1-2 August 2022) [10.1145/3531437.3539717].

Availability:

This version is available at: 11583/2968137 since: 2022-09-11T09:42:47Z

Publisher:

Association for Computing Machinery

Published

DOI:10.1145/3531437.3539717

Terms of use:

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

Publisher copyright

(Article begins on next page)

Enabling Capsule Networks at the Edge through Approximate Softmax and Squash Operations

Alberto Marchisio^{1,*}, Beatrice Bussolino^{2,*}, Edoardo Salvati^{2,*}, Maurizio Martina², Guido Masera², Muhammad Shafique³

¹ Institute of Computer Engineering, Technische Universität Wien (TU Wien), Vienna, Austria

² Department of Electronics and Telecommunications, Politecnico di Torino, Turin, Italy

³ eBrain Lab, Division of Engineering, New York University Abu Dhabi, UAE

alberto.marchisio@tuwien.ac.at, beatrice.bussolino@polito.it, edoardo.salvati@studenti.polito.it, maurizio.martina@polito.it, guido.masera@polito.it, muhammad.shafique@nyu.edu

ABSTRACT

Complex Deep Neural Networks such as Capsule Networks (CapsNets) exhibit high learning capabilities at the cost of compute-intensive operations. To enable their deployment on edge devices, we propose to leverage approximate computing for designing approximate variants of the complex operations like softmax and squash. In our experiments, we evaluate tradeoffs between area, power consumption, and critical path delay of the designs implemented with the ASIC design flow, and the accuracy of the quantized CapsNets, compared to the exact functions.

KEYWORDS

Deep Neural Networks, Capsule Networks, Approximate Computing, Nonlinear Functions, Squash, Softmax.

ACM Reference Format:

Alberto Marchisio^{1,*}, Beatrice Bussolino^{2,*}, Edoardo Salvati^{2,*}, Maurizio Martina², Guido Masera², Muhammad Shafique³. 2022. Enabling Capsule Networks at the Edge through Approximate Softmax and Squash Operations. In *ACM/IEEE International Symposium on Low Power Electronics and Design (ISLPED '22)*, August 1–3, 2022, Boston, MA, USA. ACM, New York, NY, USA, 6 pages. <https://doi.org/10.1145/3531437.3539717>

1 INTRODUCTION

In recent years, Deep Neural Networks (DNNs) have achieved outstanding performance in a wide range of applications [7][12]. Among the latest DNN models, Capsule Networks (CapsNets) [20] enable high learning capabilities due to the capsules, which add a layer of abstraction compared to the traditional neurons of DNNs. Despite their groundbreaking success, the most advanced DNNs, such as CapsNets, exhibit high complexity due to their compute-intensive operations, which hinders their deployments on energy-constrained edge devices. Therefore, several optimizations have been proposed to increase the performance and reduce the energy consumption of complex DNNs on edge devices, such as network pruning [4][8][14], and quantization [13][24]. In this work, we

*These authors contributed equally to this work.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).
ISLPED '22, August 1–3, 2022, Boston, MA, USA
© 2022 Copyright held by the owner/author(s).
ACM ISBN 978-1-4503-9354-6/22/08.
<https://doi.org/10.1145/3531437.3539717>

focus on leveraging approximate computing for CapsNet optimization. Therefore, our approach is orthogonal to other optimization techniques (e.g., quantization), as we directly perform experiments on the quantized CapsNets.

1.1 Motivations and Research Challenges

While for generic matrix multiplications (that are used on convolution operations) a common approach is to use approximate adders and multipliers [16][17], there are other complex operations (i.e., squash and softmax) that need more specialized designs to be computed in approximate form. Indeed, as shown in Fig. 1, the squash and softmax operations are the most compute-intensive operations of the CapsNets. More precisely, ① the squash constitutes the performance bottleneck of their execution on GPUs, and ② the softmax has a high execution time on a Capsule Network Hardware Accelerator (CapsAcc [15]). Hence, these results motivate our research to focus on the design of approximate squash and softmax units. The research challenges tackled in this work are to find tradeoffs for area, power, and delay of the approximate softmax and squash units, without reducing the complete CapsNets' inference accuracy much.

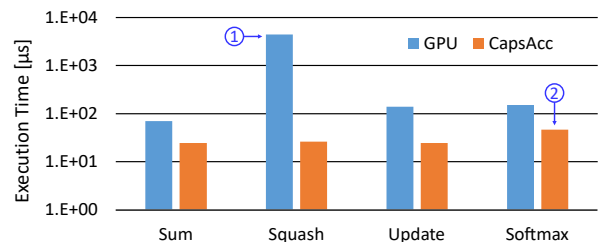


Figure 1: Execution time breakdown for the Dynamic Routing operations of the ShallowCaps [20] on the Nvidia GeForce RTX 2080 Ti GPU and the CapsAcc hardware accelerator [15].

1.2 Novel Contributions

Our novel contributions are:

- We analyze the state-of-the-art CapsNets models and the most advanced designs of approximate squash and approximate softmax (Section 2).
- We design specialized approximate softmax units using domain transformations (Section 3).
- We design approximate squash units with piecewise approximations (Section 4).

- We implement the approximate softmax and squash architectures in VHDL, synthesize them in a 45nm technology node with the ASIC design flow, and perform gate-level simulations to evaluate the area, power consumption, and critical path delay.
- We also integrate the functional approximations into the open-source Q-CapsNets framework to evaluate the inference accuracy of state-of-the-art CapsNet models using the proposed approximate units (Section 5).
- Our proposed approximate *softmax-b2* design outperforms the related works, having -11% area, -8% power, and -19% critical path delay, and comparable accuracy results.
- Our proposed approximate *squash-exp* and *squash-pow2* have up to -6% power consumption and up to -36% critical path delay compared to the state-of-the-art, while showing similar accuracy as having the exact squash function.

2 BACKGROUND AND RELATED WORKS

2.1 Capsule Networks

Capsule Networks (CapsNets) are introduced to improve the generalization ability of Convolutional Neural Networks (CNNs) in image classification tasks. CapsNets include multi-dimensional *capsules*, i.e., groups of neurons that encode the existence probability and spatial properties of a specific feature, and exploit the dynamic *routing-by-agreement* algorithm to detect entities that are consistent with lower-level features.

As a limitation, CapsNets show a higher computational complexity than traditional CNNs, as indicated by the MACs per memory ratio [13], due to the vectors of neurons and routing algorithm involving the iterative computation of complex operations, i.e., *softmax* and *squash*. The softmax function is a nonlinear function used to compute the routing coefficients connecting a lower-level capsule to the higher-level capsules. It normalizes its input values into a probability distribution. The squash function is the nonlinear activation function applied to produce the activity vector of the capsules. It ensures that the vector norm is below 1 to represent the existence probability of the associated entity and preserves the vector orientation in agreement with the lower-level capsules predictions.

The ShallowCaps model proposed by Hinton *et al.* [20] is designed for image classification on the *MNIST* dataset [11] with greyscale images of handwritten digits. The architecture consists of three layers for the inference pass, which are a convolutional layer, a convolutional capsule layer, and a fully-connected capsule layer. The first layer has $256, 9 \times 9 \times 1$ kernels with ReLU activation. The second layer applies $256, 9 \times 9 \times 256$ kernels with a stride of 2 and ReLU activation, and the output feature maps are reshaped in 32 channels of 8-dimensional capsules with squash activation. The final layer that performs the dynamic routing algorithm consists of 10 16-dimensional capsules, one for each dataset class.

The *DeepCaps* model [18] is introduced to improve the classification accuracy on complex image datasets like *CIFAR-10* [9] with color images of animals and vehicles. The architecture consists of a convolutional layer, four middle stages (*CapsCells*) including convolutional capsule layers (*ConvCaps*) and a final fully-connected capsule layer. The two main architectural improvements w.r.t. the

ShallowCaps are the skip connections that enable an efficient gradient flow during training and the routing-by-agreement algorithm based on 3D convolution to avoid the computational bottleneck that would occur by stacking multiple fully-connected capsule layers.

2.2 Approximate Computing for DNNs Nonlinear Operations

Approximate computing is an effective design methodology that aims to achieve low power consumption, high performance, and reduced circuit area by relaxing the accuracy requirement in error-tolerant applications [3]. Extensive research efforts have been dedicated to the optimization of matrix multiplications in DNNs by proposing approximate designs for adders [22] and multipliers [10]. However, a key factor for achieving high computing efficiency in DNNs and CapsNets is represented by the implementation of nonlinear functions, including nonlinearities such as sigmoid, hyperbolic tangent, softmax, and squash.

Various techniques have been proposed to compute nonlinear functions in an approximate form and enable an efficient hardware implementation with limited accuracy loss. The work in [1] proposed a piecewise linear approximation of the sigmoid function by storing the curve breakpoints in a look-up table and applying linear interpolation.

As for the softmax function, the work in [5] proposed an approximate softmax design where the exponential function is evaluated by using Taylor series expansion and a look-up table method. The work in [21] presented a hardware architecture that exploits a mathematical transformation into the logarithmic domain to simplify the division operation and approximates the exponential and logarithmic functions by using linear fitting within a specific range.

As regards the squash function, the work in [2] described a set of approximations of the Euclidean norm that avoid the computation of square and square-root operations. The work in [6] introduced an approximate square-accumulate architecture with a self-healing mechanism that is suitable for computing the sum of squared components in the Euclidean norm. *However, the previous works did not consider advanced methods, like piecewise approximations and domain transformations that are possible due to the error tolerance of these functions inserted in the CapsNets computations, that we indeed exploit in this work.*

3 APPROXIMATE SOFTMAX DESIGN

In the following, we present three approximate softmax designs describing the algorithmic approximations and the RTL architectures. The proposed softmax approximations are referred to as *softmax-taylor*, *softmax-lnu* and *softmax-b2*, with names enclosing their key features.

The softmax function shown in Eq. 1 is a probabilistic version of the argmax function, which returns 1 for the highest input value and 0 for all the other values.

$$y_i = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}} \quad (i = 1, \dots, n) \quad (1)$$

The softmax computation involves three fundamental operations: natural exponential, sum, and division. In the following approximate softmax designs, we mainly focus on the approximate

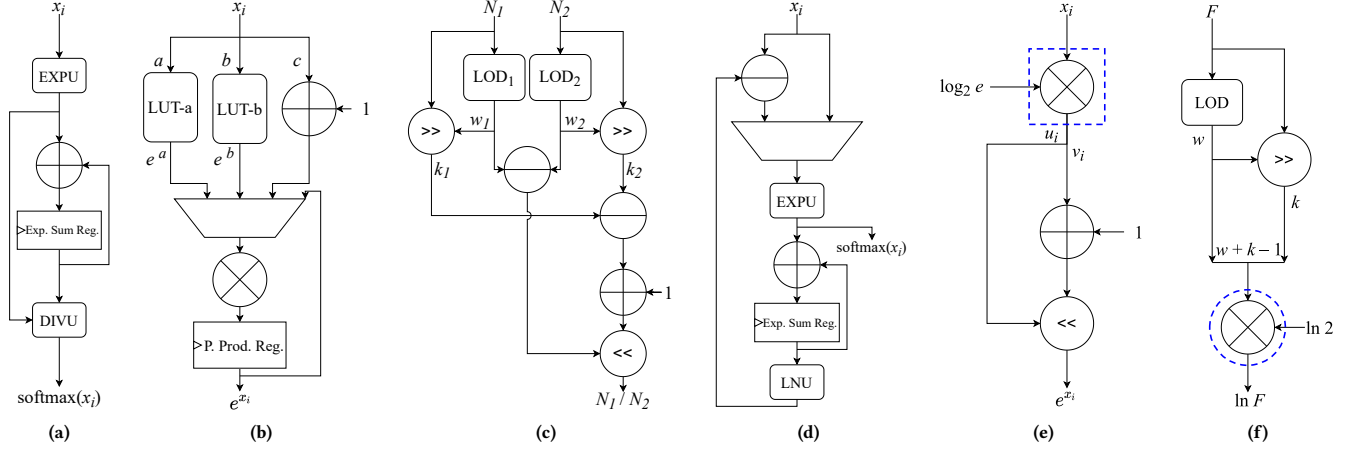


Figure 2: Architectures of the approximate softmax designs: (a) Softmax-taylor unit. (b) Softmax-taylor exponent unit. (c) Softmax-taylor division unit. (d) Softmax-lnu unit. (e) Softmax-lnu exp unit. (f) Softmax-lnu natural log unit. Softmax-b2 unit replaces expu and lnu in Softmax-lnu with pow2u and log2u by removing the multipliers $\log_2 e$ and $\ln 2$.

computation of the exponentiation and division, which are the most complex operations of the softmax function.

The **softmax-taylor** design is based on a specific softmax approximation [5] which exploits the Taylor series expansion method for exponential computation and performs divisions in the logarithmic domain.

The natural exponent operation is simplified as in Eq. 2 using the first-order Taylor polynomial approximation. At the architecture level, the exponent unit consists of 2 look-up tables to implement the first two exponent contributions, a specific bus arrangement to get $1 + c$ and a multiplier to compute the final product iteratively (see Figure 2a and Figure 2b).

$$e^{x_i} = e^{a+b+c} \approx e^a \cdot e^b \cdot (1 + c) \quad (2)$$

The division operation is performed in the logarithmic domain by exploiting the mathematical transformation in Eq. 3

$$\begin{aligned} \text{pow2}(\log_2(e^{x_i} / \sum_{j=1}^n e^{x_j})) &= \text{pow2}(w_1 + \log_2 k_1 - (w_2 + \log_2 k_2)) \\ &\approx \text{pow2}(w_1 - w_2 + k_1 - k_2) = 2^{u_i + v_i} \approx 2^{u_i} \cdot (1 + v_i) \end{aligned} \quad (3)$$

First of all, $N_1 = e^{x_i}$ and $N_2 = \sum_{j=1}^n e^{x_j}$ are expressed as $2^{w_l} \cdot k_l$, with $w_l \in \mathbb{Z}$ and $k_l \in [1, 2)$ for $l = 1, 2$ and the base-2 logarithm of k_l is approximated by the linear fitting function $k_l - 1$. Secondly, the argument of the power-2 operation is split into its integer and fractional parts, u_i and v_i , with $u_i \in \mathbb{Z}$ and $v_i \in [0, 1)$ and 2^{v_i} is estimated as $(1 + v_i)$.

The division unit is composed of 2 base-2 logarithm units, a leading one detector (LOD) and shift unit that compute the logarithm of the dividend and divisor, a subtraction unit that performs the division in the log domain, and a power-2 unit (bus arrangement and shift unit) that computes the softmax output value (see Figure 2c).

To be compliant with the capsule network models [20] [18] that we use in our experiments, the softmax architecture is able to process 10, 32 or 128 inputs.

The **softmax-lnu** design builds on a peculiar softmax approximation [21] which adopts a mathematical domain transformation

involving natural logarithm and natural exponential operations (see Eq. 4).

$$\exp(\ln(e^{x_i} / \sum_{j=1}^n e^{x_j})) = \exp(x_i - \ln(\sum_{j=1}^n e^{x_j})) \quad (4)$$

The transformation into the logarithm domain allows to perform the division by using a more straightforward subtraction. On the other hand, the exponentiation is required to convert the softmax output values from the logarithmic domain to the linear one. The design architecture mainly consists of three computational units to compute the natural exponential of the softmax inputs (EXPU), sum up the exponentials, and evaluate the natural logarithm of the sum (LNU) required for the division (see Figure 2d).

The natural exponential operation is performed by using the mathematical transformation in Eq. 5, with $u_i \in \mathbb{Z}$ and $v_i \in [0, 1)$. At the architecture level, the natural exponential unit is composed of a constant multiplier by $\log_2 e$, a specific bus arrangement to implement $1 + v_i$ and a shift unit to compute the result (see Figure 2e).

$$e^{x_i} = 2^{x_i \cdot \log_2 e} = 2^{u_i + v_i} = 2^{u_i} \cdot 2^{v_i} \approx 2^{u_i} \cdot (1 + v_i) \quad (5)$$

The natural logarithm is computed as in Eq. 6, where $F = \sum_{j=1}^n e^{x_j}$ is expressed as $2^w \cdot k$, with $w \in \mathbb{Z}$ and $k \in [1, 2)$ and the base-2 logarithm of k is approximated by the linear fitting function $k - 1$. The natural logarithm unit consists of 4 main subunits: a leading one detector to determine w , a shift unit to compute k , a specific bus arrangement to get the base-2 logarithm of F and a constant multiplier by $\ln 2$ (see Figure 2f).

$$\ln F = \ln 2 \cdot \log_2 F = \ln 2 \cdot (w + \log_2 k) \approx \ln 2 \cdot (w + k - 1) \quad (6)$$

The architecture includes other hardware units to compute the maximum input value, scale the inputs, execute the division in the log domain, and allow for the processing of a variable number of softmax inputs.

The **softmax-b2** design implements the idea of computing a softmax-like function with powers of 2 in place of natural exponentials, and it exploits a domain transformation with base-2 logarithm

and power-2 operations (see Eq. 7).

$$\text{pow2}(\log_2(2^{x_i} / \sum_{j=1}^n 2^{x_j})) = \text{pow2}(x_i - \log_2(\sum_{j=1}^n 2^{x_j})) \quad (7)$$

The proposed approximation allows for a complexity reduction of the hardware implementation of the *softmax-lnu* design, thanks to the removal of two constant multipliers.

Compared to the *softmax-lnu* design, the *softmax-b2* architecture avoids the preliminary multiplication by $\log_2 e$ in the exponential unit (see dashed square in Figure 2e) and the final multiplication by $\ln 2$ in the logarithmic unit (see dashed circle in Figure 2f), by implementing the power-2 and base-2 logarithm unit, respectively.

4 APPROXIMATE SQUASH DESIGN

The proposed approximate squash designs are called *squash-norm*, *squash-exp* and *squash-pow2*.

The squash function requires to compute the norm of the input vector and the squashing coefficient that multiplies the input vector to produce the output vector, as shown in Eq. 8.

$$\mathbf{y} = \frac{\|\mathbf{x}\|^2}{1 + \|\mathbf{x}\|^2} \frac{\mathbf{x}}{\|\mathbf{x}\|} \quad (8)$$

The first design exploits a specific norm approximation [2], while the remaining two techniques introduce novel solutions to approximate the squashing coefficient.

The **squash-norm** design is inspired by the specific Euclidean norm approximation proposed by Chaudhuri *et al.* [2], which is shown in Eq. 9.

$$\|\mathbf{x}\| \approx D_\lambda(\mathbf{x}) = |x_{i_{max}}| + \lambda \sum_{\substack{i=1 \\ i \neq i_{max}}}^n |x_i| \quad (9)$$

This architecture does not require the square root operator and the multiplications needed to square the vector components, but it involves the computation of the absolute values and the maximum absolute value components. The parameter λ depends on the number of vector components and it is selected accordingly [19].

The designed architecture is composed of two main units. The norm unit computes the approximate vector norm, and the squashing unit produces the squash outputs (see Figure 3a).

The norm unit implements the Chaudhuri approximation [2] in Eq. 9. It consists of multiple arithmetic modules. A dedicated component computes the absolute value of the inputs, an accumulator sums up the absolute values, a unit determines the maximum absolute value, a subtractor gets the second term of the formula, a multiplier scales the sum by λ and an adder adds the maximum value to the sum (see Figure 3b).

The squashing unit consists of two look-up tables to implement the squashing coefficient and a multiplier to compute the squash outputs as the product between the inputs and the squashing coefficient (see Figure 3c).

To be compliant with the two capsule network models employed in our experiments, the squash architecture is able to process 4, 8, 16, or 32 inputs.

The **squash-exp** design exploits a piecewise approximation of the squashing coefficient $\|\mathbf{x}\|/(1 + \|\mathbf{x}\|^2)$ in two ranges of norm values. The coefficient is approximated by the nonlinear function

$1 - e^{-\|\mathbf{x}\|}$ in the first range and by a direct mapping method in the second range. The range of norm values is derived experimentally by executing inference steps with two capsule network models on two image datasets.

At the architecture level, the design mainly consists of two computational units: the norm unit and the squashing unit.

The norm unit computes the Euclidean norm of the input vector. It is composed of a multiplier to square the input components, an accumulator to sum up the squared inputs, and two look-up tables to implement the square root function over two specific ranges of squared norm values (see Figure 3d).

The squashing unit implements the piecewise approximation of the squashing coefficient and computes the output values. The nonlinear function in the first range is implemented by a component composed of a 2's complement of the norm value, a natural exponential unit, and a subtractor. The second-range approximation is performed with a look-up table. The final multiplier is used to compute the squash outputs (see Figure 3e).

The **squash-pow2** design builds on the piecewise approximation of the squashing coefficient used in the *squash-exp* architecture, but the approximating nonlinear function used in the first range of norm values is $1 - 2^{-\|\mathbf{x}\|}$.

At the architecture level, in the exponential unit the constant multiplication by $\log_2 e$ is removed to implement the power-2 unit (see Figure 3f). The hardware cost reduction is obtained at the expense of a higher worst-case approximation error of the squashing coefficient in the range of low norm values (see Figure 4).

5 EVALUATION OF OUR DESIGNS

In the following, the approximate softmax and squash designs are evaluated in terms of inference accuracy loss and hardware implementation metrics.

First, we explore the inference accuracy degradation induced by the proposed softmax and squash approximations in 4 case studies, with two capsule network models on two image classification datasets. Secondly, we synthesize the complete architectures and analyze our designs' area usage, power consumption, and timing performance.

The objective of the evaluation is to explore possible *trade-offs* between the classification accuracy loss of a CapsNet using the approximations and the hardware implementation cost of the approximate designs.

5.1 Experimental Setup

We implement the approximate softmax and squash algorithms in Python and perform extensive software simulations to evaluate the quality of each approximation w.r.t. the exact function. The experiments are conducted for over 1,000 input vectors in a specific range. We analyze the *Mean Error Distance* on the maximum and average component errors, in absolute and relative terms.

To assess how the softmax and squash approximations affect the inference accuracy of the complete capsule networks, we include the approximate functions in a Python-based CapsNet model provided by the open-source framework *Q-CapsNets* [13] and we perform an image classification task with two CapsNet models, ShallowCaps [20] and DeepCaps [18], on two image datasets, MNIST [11] and Fashion-MNIST [23].

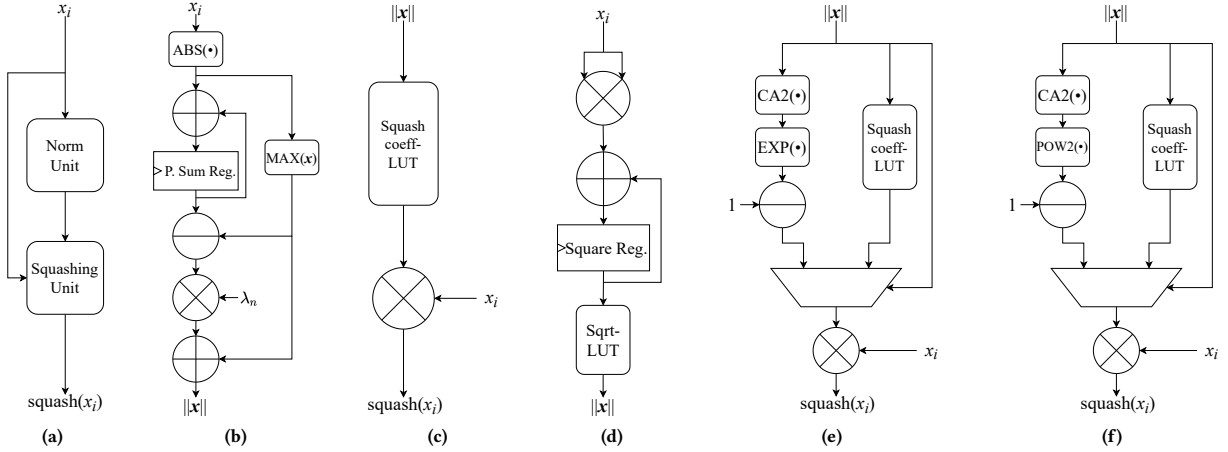


Figure 3: Architectures of the approximate squash designs: (a) Squash function unit. (b) Squash-norm norm unit. (c) Squash-norm squashing unit. (d) Squash-exp and -pow2 norm unit. (e) Squash-exp squashing unit. (f) Squash-pow2 squashing unit.

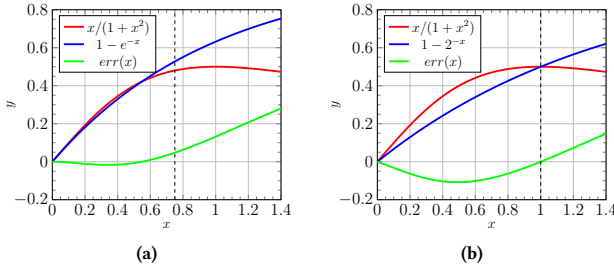


Figure 4: Behavior of (a) squash-exp and (b) squash-pow2 approximation with $x := ||x||$.

As shown in Figure 5, our experimental setup consists of both software and hardware components. We use a software environment with PyTorch library and Nvidia CUDA Toolkit and execute the inference passes on an Nvidia GeForce RTX 2080 Ti GPU. To comply with the hardware implementation, we perform the quantization of the approximate softmax and squash data, and we test the quantized approximate functions in quantized CapsNet models (see Table 1). Using the Q-CapsNets framework, we quantize weights and activations of the CapsNet models on the image datasets and input data of the softmax and squash functions.

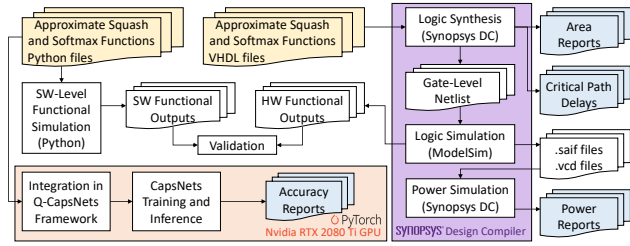


Figure 5: Setup and tool-flow for conducting our experiments.

We implement the design architectures in VHDL and perform functional simulation using ModelSim to check the results against the Python model outputs. We synthesize the architectures in a 45nm academic technology library, *Nangate OCL*, by using the

Table 1: Percentage of quantized inference accuracy.

	MNIST		Fashion-MNIST	
	ShallowCaps	DeepCaps	ShallowCaps	DeepCaps
<i>exact functions</i>	99.44	99.35	92.42	94.69
softmax-lnu	99.46	99.42	92.37	94.71
softmax-b2 (ours)	99.49	99.33	92.33	94.64
softmax-taylor	99.42	99.41	92.47	94.69
squash-exp (ours)	99.18	98.79	91.32	94.76
squash-pow2 (ours)	99.00	98.58	89.05	94.62
squash-norm	99.26	99.23	92.51	94.70

ASIC design flow with Synopsys Design Compiler and obtain area usage, power consumption, and maximum path delay of each design (see Table 2). Finally, we conduct post-synthesis functional and timing validation of the gate-level netlist.

Table 2: Hardware characteristics with clock frequency 100 MHz.

	Area usage (μm^2)	Power consumption (μW)	Critical path delay (ns)
softmax-lnu	12,511	2,572	6.46
softmax-b2 (ours)	11,169	2,244	4.22
softmax-taylor	14,944	2,430	5.24
squash-exp (ours)	7,937	1,414	5.64
squash-pow2 (ours)	7,543	1,340	4.17
squash-norm	6,806	1,431	6.53

5.2 Evaluating the Softmax

From the experimental results, we derive the following key observations regarding the approximate softmax designs.

The **softmax-b2** design is the best solution in terms of hardware metrics but it implies the highest CapsNet accuracy loss in all the case studies except for the ShallowCaps on MNIST. Actually, the *b2* design consumes less area (-11% and -25%) and power (-13% and -8%) than the *lnu* and *taylor* designs. Moreover, it has the lowest critical path delay (-35% and -19% w.r.t. *lnu* and *taylor*).

The **softmax-taylor** design is the best choice in terms of inference accuracy loss since it outperforms the other designs in the ShallowCaps for Fashion-MNIST. However, it is characterised

by the worst area usage (+20% and +35% w.r.t. *lnu* and *b2*) and intermediate power consumption and critical path delay.

The **softmax-*lnu*** design shows the highest power consumption (+15% and +5% w.r.t. *b2* and *taylor*) and maximum path delay (+53% and +23%) but intermediate area usage. Its performance in inference accuracy loss is similar to the *taylor* design in all the case studies, except for the ShallowCaps for Fashion-MNIST, where the *lnu* performs worse (+0.1% loss).

5.3 Evaluating the Squash

The **squash-norm** design is the best approximate squash solution in terms of CapsNet accuracy loss. It also has the benefit of having the best area usage (−13% and −8% w.r.t. *exp* and *pow2*), but as a drawback, it shows the worst power (+1% and +7%) and delay metrics (+15% and +56%).

The **squash-pow2** design is the best option in terms of power consumption (−5% and −6% w.r.t. *exp* and *norm*) and critical path delay (−25% and −36%), and it has intermediate area usage. However, it implies the highest CapsNet accuracy loss among all the case studies.

The **squash-exp** design is characterized by an accuracy loss similar to the *norm* design in two case studies and significantly worse accuracy in the other two cases. In exchange for the reduced accuracy, it has intermediate power and delay metrics, but as a downside, it shows the worst area usage (+5% and +17% w.r.t. *pow2* and *norm*).

6 CONCLUSION

To enable efficient CapsNets inference on edge devices, we propose approximate designs for the most compute-intensive CapsNets operations, which are the softmax and squash. Our **softmax-*b2*** design based on approximating the natural exponential with powers of 2 significantly reduces the hardware complexity, with limited accuracy drop. Our squash designs based on piecewise approximations show interesting tradeoffs between accuracy, area, power consumption, and critical path delay. We believe that our findings will contribute to the deployment of CapsNets and other complex DNN models on resource-constrained devices.

ACKNOWLEDGMENTS

This work has been supported in part by the Doctoral College Resilient Embedded Systems, which is run jointly by the TU Wien's Faculty of Informatics and the UAS Technikum Wien.

REFERENCES

- [1] H. Amin, K.M. Curtis, and B.R. Hayes-Gill. 1997. Piecewise linear approximation applied to nonlinear function of a neural network. *IEE Proceedings - Circuits, Devices and Systems* 144, 6 (1997), 313–317.
- [2] M. Emre Celebi, Fatih Celiker, and Hassan A. Kingravi. 2011. On Euclidean norm approximations. *Pattern Recognit.* 44, 2 (2011), 278–283. <https://doi.org/10.1016/j.patcog.2010.08.028>
- [3] Vinay K. Chippa, Swagath Venkataramani, Srimat T. Chakradhar, Kaushik Roy, and Anand Raghunathan. 2013. Approximate computing: An integrated hardware approach. In *2013 Asilomar Conference on Signals, Systems and Computers*. IEEE, 111–117. <https://doi.org/10.1109/ACSSC.2013.6810241>
- [4] Jonathan Frankle and Michael Carbin. 2019. The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks. In *7th International Conference on Learning Representations, ICLR 2019*. OpenReview.net.
- [5] Yue Gao, Weiqiang Liu, and Fabrizio Lombardi. 2020. Design and Implementation of an Approximate Softmax Layer for Deep Neural Networks. In *IEEE International Symposium on Circuits and Systems, ISCAS 2020*. IEEE, 1–5. <https://doi.org/10.1109/ISCAS45731.2020.9180870>
- [6] G. A. Gillani, Muhammad Abdullah Hanif, M. Krone, Sabih H. Gerez, Muhammad Shafique, and André B. J. Kokkeler. 2018. Squash: Approximate Square-Accumulate With Self-Healing. *IEEE Access* 6 (2018), 49112–49128. <https://doi.org/10.1109/ACCESS.2018.2868036>
- [7] Sorin Mihai Grigorescu, Bogdan Trasnea, Tiberiu T. Cocias, and Gigel Macesanu. 2019. A Survey of Deep Learning Techniques for Autonomous Driving. *arXiv abs/1910.07738* (2019). [arXiv:1910.07738](https://arxiv.org/abs/1910.07738)
- [8] Song Han, Jeff Pool, John Tran, and William J. Dally. 2015. Learning both Weights and Connections for Efficient Neural Network. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015*. 1135–1143.
- [9] Alex Krizhevsky. 2012. Learning Multiple Layers of Features from Tiny Images. *University of Toronto* (05 2012).
- [10] Parag Kulkarni, Puneet Gupta, and Milos D. Ercegovac. 2011. Trading Accuracy for Power with an Underdesigned Multiplier Architecture. In *VLSI Design 2011: 24th International Conference on VLSI Design*. IEEE Computer Society, 346–351. <https://doi.org/10.1109/VLSID.2011.51>
- [11] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. 1998. Gradient-based learning applied to document recognition. *Proc. IEEE* (1998).
- [12] Zewen Li, Fan Liu, Wenjie Yang, Shouheng Peng, and Jun Zhou. 2021. A Survey of Convolutional Neural Networks: Analysis, Applications, and Prospects. *IEEE Transactions on Neural Networks and Learning Systems* (2021), 1–21. <https://doi.org/10.1109/TNNLS.2021.3084827>
- [13] Alberto Marchisio, Beatrice Bussolino, Alessio Colucci, Maurizio Martina, Guido Masera, and Muhammad Shafique. 2020. Q-CapsNets: A Specialized Framework for Quantizing Capsule Networks. In *57th ACM/IEEE Design Automation Conference, DAC 2020*. IEEE, 1–6. <https://doi.org/10.1109/DAC18072.2020.9218746>
- [14] Alberto Marchisio, Muhammad Abdullah Hanif, Maurizio Martina, and Muhammad Shafique. 2018. PruNet: Class-Blind Pruning Method For Deep Neural Networks. In *2018 International Joint Conference on Neural Networks, IJCNN 2018*. IEEE, 1–8. <https://doi.org/10.1109/IJCNN.2018.8489764>
- [15] Alberto Marchisio, Muhammad Abdullah Hanif, and Muhammad Shafique. 2019. CapsAcc: An Efficient Hardware Accelerator for CapsuleNets with Data Reuse. In *In Design, Automation & Test in Europe Conference & Exhibition, DATE 2019*. IEEE, 964–967. <https://doi.org/10.23919/DATE.2019.8714922>
- [16] Alberto Marchisio, Vojtech Mrazek, Muhammad Abdullah Hanif, and Muhammad Shafique. 2020. ReD-CaNe: A Systematic Methodology for Resilience Analysis and Design of Capsule Networks under Approximations. In *2020 Design, Automation & Test in Europe Conference & Exhibition, DATE 2020*. IEEE, 1205–1210. <https://doi.org/10.23919/DATE48585.2020.9116393>
- [17] Vojtech Mrazek, Zdenek Vasicek, Lukás Sekanina, Muhammad Abdullah Hanif, and Muhammad Shafique. 2019. ALWANN: Automatic Layer-Wise Approximation of Deep Neural Network Accelerators without Retraining. In *Proceedings of the International Conference on Computer-Aided Design, ICCAD 2019*. 1–8. <https://doi.org/10.1109/ICCAD45719.2019.8942068>
- [18] Jathushan Rajasegaran, Vinoj Jayasundara, Sandaru Jayasekara, Hirunima Jayasekara, Suranga Seneviratne, and Ranga Rodrigo. 2019. DeepCaps: Going Deeper With Capsule Networks. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019*. 10725–10733. <https://doi.org/10.1109/CVPR.2019.01098>
- [19] Frank Rhodes. 1995. On the metrics of Chaudhuri, Murthy and Chaudhuri. *Pattern Recognit.* 28, 5 (1995), 745–752. [https://doi.org/10.1016/0031-3203\(94\)00134-8](https://doi.org/10.1016/0031-3203(94)00134-8)
- [20] Sara Sabour, Nicholas Frosst, and Geoffrey E. Hinton. 2017. Dynamic Routing Between Capsules. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017*. 3856–3866.
- [21] Meiqi Wang, Siyuan Lu, Danyang Zhu, Jun Lin, and Zhongfeng Wang. 2018. A High-Speed and Low-Complexity Architecture for Softmax Function in Deep Learning. In *2018 IEEE Asia Pacific Conference on Circuits and Systems, APCCAS 2018*. IEEE, 223–226. <https://doi.org/10.1109/APCCAS.2018.8605654>
- [22] Yi Wu, You Li, Xiangxuan Ge, Yuan Gao, and Weikang Qian. 2019. An Efficient Method for Calculating the Error Statistics of Block-Based Approximate Adders. *IEEE Trans. Computers* 68, 1 (2019), 21–38. <https://doi.org/10.1109/TC.2018.2859960>
- [23] Han Xiao, Kashif Rasul, and Roland Vollgraf. 2017. Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms. *CoRR abs/1708.07747* (2017). [arXiv:1708.07747](https://arxiv.org/abs/1708.07747)
- [24] Yiren Zhou, Seyed-Mohsen Moosavi-Dezfooli, Ngai-Man Cheung, and Pascal Frossard. 2018. Adaptive Quantization for Deep Neural Network. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th Innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18)*. AAAI Press, 4596–4604.