# CapsAcc: An Efficient Hardware Accelerator for CapsuleNets with Data Reuse

Alberto Marchisio, Muhammad Abdullah Hanif, and Muhammad Shafique

*Vienna University of Technology, Vienna, Austria*

{alberto.marchisio,muhammad.hanif,muhammad.shafique}@tuwien.ac.at

*Abstract*—**Deep Neural Networks (DNNs) have been widely deployed for many Machine Learning applications. Recently, *CapsuleNets* have overtaken traditional DNNs, because of their improved generalization ability due to the multi-dimensional capsules, in contrast to the single-dimensional neurons. Consequently, CapsuleNets also require extremely intense matrix computations, making it a gigantic challenge to achieve high performance. In this paper, we propose *CapsAcc*, the first specialized CMOS-based hardware architecture to perform CapsuleNets inference with high performance and energy efficiency. State-of-the-art convolutional DNN accelerators would not work efficiently for CapsuleNets, as their designs do not account for key operations involved in CapsuleNets, like squashing and dynamic routing, as well as multi-dimensional matrix processing. Our CapsAcc architecture targets this problem and achieves significant improvements, when compared to an optimized GPU implementation. Our architecture exploits the massive parallelism by flexibly feeding the data to a specialized systolic array according to the operations required in different layers. It also avoids extensive load and store operations on the on-chip memory, by reusing the data when possible. We further optimize the routing algorithm to reduce the computations needed at this stage. We synthesized the complete CapsAcc architecture in a 32nm CMOS technology using Synopsys design tools, and evaluated it for the MNIST benchmark (as also done by the original CapsuleNet paper) to ensure consistent and fair comparisons. This work enables highly-efficient CapsuleNets inference on embedded platforms.**

## I. Introduction

Machine Learning (ML) algorithms are widely used for Internet of Things and Artificial Intelligence applications, such as computer vision [7], speech recognition [3] and natural language processing [2]. Deep Neural Networks (DNNs) have reached state-of-the-art results in terms of accuracy, compared to other ML algorithms. Recently, Sabour and Hinton et al. [12] proposed the Dynamic Routing algorithm to efficiently perform training and inference on CapsuleNets [5]. Such CapsuleNets are able to encapsulate multi-dimensional features across the layers, while traditional Convolutional Neural Networks (CNNs) do not. Thus, *CapsuleNets can beat traditional CNNs in multiple tasks*, like image classification, as shown in [12]. The most evident difference is that the CapsuleNets are deeper in width than in height, when compared to DNNs, as each capsule incorporates the information hierarchically, thus preserving other features like position, orientation and scaling (see an overview of CapsuleNets in Section II). The data is propagated towards the output using the so-called routing-by-agreement algorithm.

Current state-of-the-art DNN accelerators [4] [1] [6] [11] [9] proposed energy-aware solutions for inference using traditional CNNs. As for our knowledge, **we are the first to propose a hardware accelerator-based architecture for the complete CapsuleNets inference**. Although systolic array based designs like [9] perform parallel matrix multiply-and-accumulate (MAC) operations with good efficiency, the existing CNN accelerators cannot compute several key operations of the CapsuleNets (i.e., squashing and routing-by-agreement) with high performance. An efficient data-flow mapping requires a direct feedback connection from the outputs coming from the activation unit back to the inputs of the processing element. Thus, such key optimizations

can highly increase the performance and reduce the memory accesses.

**Our Novel Contributions:**

1) We analyze the memory requirements and the performance in the forward pass of CapsuleNets, through experiments on a high-end GPU, which allows to identify the corresponding bottlenecks.
2) We propose CapsAcc, an accelerator that can perform inference on CapsuleNets with an efficient data reuse based mapping policy.
3) We optimize the routing-by-agreement process at algorithm level, by skipping the first step and directly initializing the coupling coefficients.
4) We implement and synthesize the complete CapsAcc architecture for a 32nm technology using the ASIC design flow, and perform evaluations for performance, area and power consumption. We performed the functional and timing validation through gate-level simulations. Our results demonstrate a speed-up of $12\times$ in the ClassCaps layer, of $172\times$ in the Squashing and of $6\times$ in the overall CapsuleNet inference, compared to a highly optimized GPU implementation.

**Paper Organization:** Section II summarizes the fundamental theory behind CapsuleNets and highlights the differences with traditional DNNs. In Section III, we systematically analyze the forward pass of the CapsuleNets executing on a GPU, to identify the potential bottlenecks. Section IV describes the architectural design of our CapsuleNet accelerator, for which the data-flow mapping is presented in Section V. The results are presented in Section VI.

## II. Background: An Overview of CapsuleNets

Sabour and Hinton et al. [12] introduced many novelties compared to CNNs, such as the concept of capsules, the squashing activation function, and the routing-by-agreement algorithm. In this paper, since we analyze the inference process, the layers and the algorithms that are involved in the training process *only* (e.g., decoder, margin loss and reconstruction loss) are not discussed.

### A. CapsuleNet Architecture

Figure 1 illustrates the CapsuleNet architecture [12] designed for the MNIST [8] dataset. It consists of 3 layers:

- **Conv1:** traditional convolutional layer, with 256 channels, with filter size of 9x9, stride=1, and ReLU activations.
- **PrimaryCaps:** first capsule layer, with 32 channels. Each eight-dimensional (8D) capsule has 9x9 convolutional filters with stride=2.
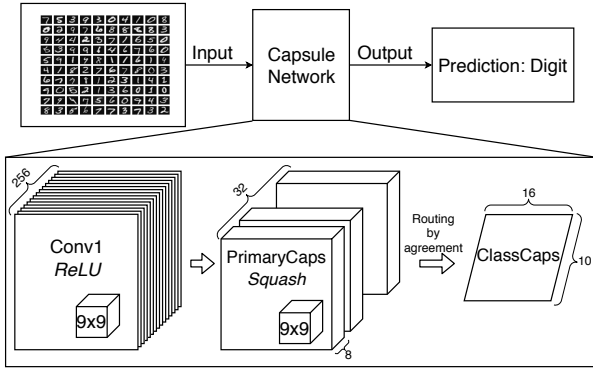- **ClassCaps:** last capsule layer, with 16D capsules for each output class.

Fig. 1: An overview of the CapsuleNet architecture, based on the design of [12] for the MNIST dataset.
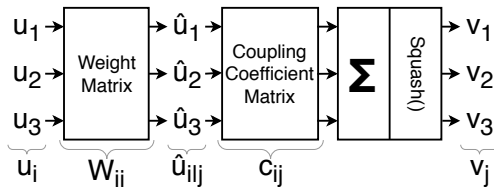


Fig. 2: Simple representation of how a CapsuleNet works.

At the first glance, it is evident that a capsule layer contains multi-dimensional capsules, which are groups of neurons nested inside a layer. One of the main advantages of CapsuleNets over traditional CNNs is the ability to learn the hierarchy between layers, because *each capsule element is able to learn different types of information* (e.g., position, orientation and scaling). Indeed, CNNs have limited model capabilities, which they try to compensate by increasing the amount of training data (with more samples and/or data augmentation) and by applying pooling to select the most important information that will be propagated to the following layers. In capsule layers, however, the outputs are propagated towards the following layers in form of a prediction vector, whose size is defined by the capsule dimension. A simple visualization of how a CapsuleNet works is presented in Figure 2. After the weight matrix multiplication $W_{ij}$, the values $\hat{u}_{i|j}$ are multiplied by the coupling coefficients $c_{ij}$, before summing together the contributions and applying the squash function. The coupling coefficients are computed and updated at run-time during each inference pass, using the *routing-by-agreement* algorithm (Figure 4).

### B. Squashing

The squashing is an activation function designed to efficiently fit for the prediction vector. It introduces the nonlinearity into an array and normalizes the outputs to values between 0 and 1. Given $s_j$ as the input of the capsule $j$ (or, from another perspective, the sum of the weighted prediction vector) and $v_j$ as its respective output, the squashing function is defined by the Equation (1).

The behaviors of the squashing function and its first derivative are shown in Figure 3. Note that we have plotted the single-dimensional input function, since a multi-dimensional input version cannot be visualized in a chart. The squashing function produces an output bounded between 0 and 1, while its first derivative follows the behavior of the red line, with a peak at the point $(0.5767, 0.6495)$.
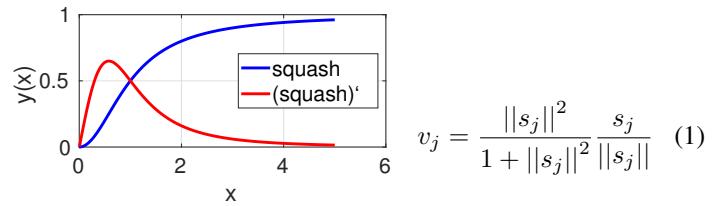


$$v_j = \frac{||s_j||^2}{1 + ||s_j||^2} \frac{s_j}{||s_j||} \quad (1)$$

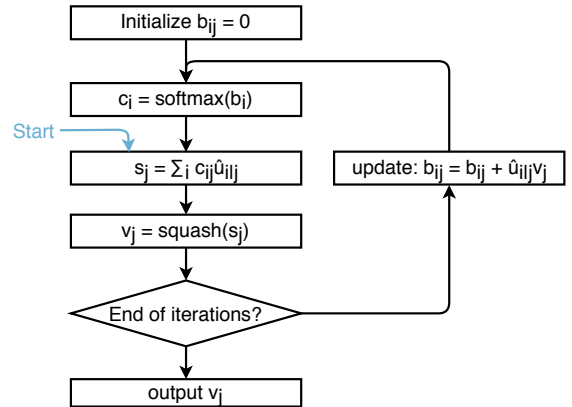Fig. 3: Squashing function and its first derivative, considering single-dimensional input.



Fig. 4: Flow of the routing-by-agreement algorithm.

### C. Routing-by-Agreement Algorithm

The predictions are propagated across two consecutive capsule layers through the routing-by-agreement algorithm. It is an iterative process, that introduces a feedback path in the inference pass. For clarity, we present the flow diagram (Figure 4) of the routing-by-agreement at software level. Note, this algorithm introduces a loop in the forward pass, because the coupling coefficients $c_{ij}$ are learned during the routing, as their values depend on the current data. Thus, they cannot be considered as constant parameters, learned during the training process. Intuitively, this step can cause a computational bottleneck, as demonstrated in Section III.

### III. MOTIVATIONAL ANALYSIS OF CAPSULENET COMPLEXITY

In the following, we perform a comprehensive analysis to identify how CapsuleNet inference is performed on a standard GPU platform, like the one used in our experiments, i.e., the Nvidia Ge-Force GTX1070 GPU (see Figure 6). First, in Section III-A we quantitatively analyze how many trainable parameters per layer must be fed from the memory. Then, in Section III-B we benchmark our pyTorch [13] based CapsuleNet implementation for the MNIST dataset to measure the performance of the inference process on our GPU.

### A. Trainable parameters of the CapsuleNet

Figure 5 shows quantitatively how many parameters are needed for each layer. As evident, the majority of the weights belong to the PrimaryCaps layer, due to its 256 channels and 8D capsules. Even if the ClassCaps layer has fully-connected behavior, it counts just for less than 25% of the total parameters of the CapsuleNet. Finally, Conv1 and the coupling coefficients counts for a very small percentage of the parameters. The detailed computation of the parameters is reported in Table I. Based
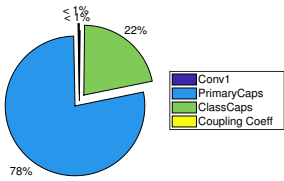
**Fig. 5:** Distribution of trainable parameters on the CapsuleNet across different layers.

**TABLE I:** Input size, number of trainable parameters and output size of each layer of the CapsuleNet.

| | Inputs | # parameters | Outputs |
|---|---|---|---|
| Conv1 | 784 | 20992 | 102400 |
| PrimaryCaps | 102400 | 5308672 | 102400 |
| ClassCaps | 102400 | 1474560 | 160 |
| Coupling Coeff | 160 | 11520 | 160 |

**Fig. 8:** Layer-wise performance of the inference pass of the CapsuleNet.

**Fig. 9:** Performance of the inference pass on each step of the routing-by-agreement algorithm.

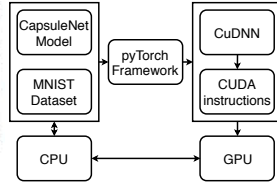**Fig. 6:** Nvidia Ge-Force GTX1070.

**Fig. 7:** Experimental setup for GPU analyses.

on that, we make an observation valuable for designing our hardware accelerator: *by considering an 8-bit fixed point weight representation, we can estimate that an on-chip memory size of 8MB is large enough to contain every parameter of the CapsuleNet.*

### B. Performance Analysis on a GPU

At this stage, we measure the time required for an inference pass on the GPU. The experimental setup is shown in Figure 7. Figure 8 shows the measurements for each layer. The ClassCaps layer is the computational bottleneck, because it is around $10\times$ slower than the previous layers. To obtain more detailed results, a further analysis has been performed, regarding the performance for each step of the routing-by-agreement (Figure 9). It is evident that *the Squashing operation inside the ClassCaps layer represents the most compute-intensive operation.* This analysis gives us the motivation to *spend more effort in optimizing routing-by-agreement and squashing* in our CapsuleNet accelerator.

### C. Summary of Key Observations from our Analyses

From the analyses performed in Sections III-A and III-B, we derive the following key observations:

- The CapsuleNet inference performed on GPU is more compute-intensive than memory-intensive, because the bottleneck is represented by the squashing operation.
- A massive parallel computation capability in the hardware accelerator is desirable to achieve the same or a better level of performance than the GPU for Conv1 and ClassCaps layers.
- Since the overall memory required to store all the weights is quite high, the buffers located in between the on-chip memory and the processing elements are beneficial to maintain high throughput and to mitigate the latency due to on-chip memory reads.

### IV. DESIGNING THE CAPSACC ARCHITECTURE

Following the above observations, we designed the complete CapsAcc accelerator and implemented it in hardware (RTL). The top-level architecture is shown in Figure 10, where the
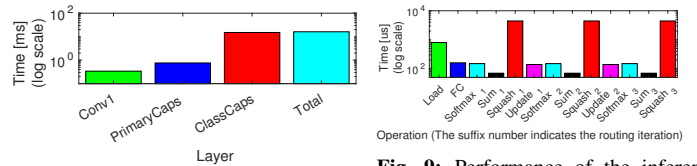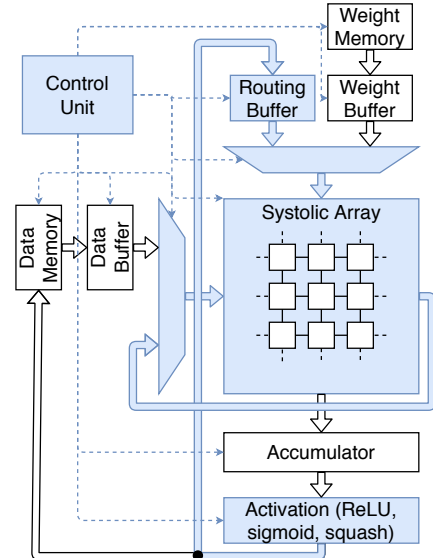


**Fig. 10:** Overview of our CapsAcc Architecture.

blue-colored blocks highlight our novel contributions over other existing accelerators for CNNs. The detailed architectures of different components of our accelerator are shown in Figure 11. Our CapsAcc architecture has a systolic array supporting a specialized data-flow mapping (see Section V), which allows to exploit the computational parallelism for multi-dimensional matrix operations. The partial sums are stored and properly added together by the accumulator unit. The activation unit performs different activation functions, according to the requirements for each stage. The buffers (Data, Routing and Weight Buffers) are essential to temporarily store the information to feed the systolic array without accessing every time to the data and weight memories. The two multiplexers in front of the systolic array introduce the flexibility to process new data or reuse them, according to the data-flow mapping. The control unit coordinates all the accelerator operations, at each stage of the inference.

### A. Systolic Array

The systolic array of our CapsAcc architecture is shown in Figure 11a. It is composed of a 2D array of Processing Elements (PEs), with $n$ rows and $m$ columns. For illustration and space reasons, Figure 11a presents the $4\times4$ version, while in our actual CapsAcc design we use a $16\times16$ systolic array. The inputs are propagated towards the outputs of the systolic array both horizontally (Data) and vertically (Weight, Partial sum). In the first row, the inputs corresponding to the Partial sums are zero-valued, because each sum at this stage is equal to $0$. Meanwhile, the Weight outputs in the last row are not connected, because they are not used in the following stages.
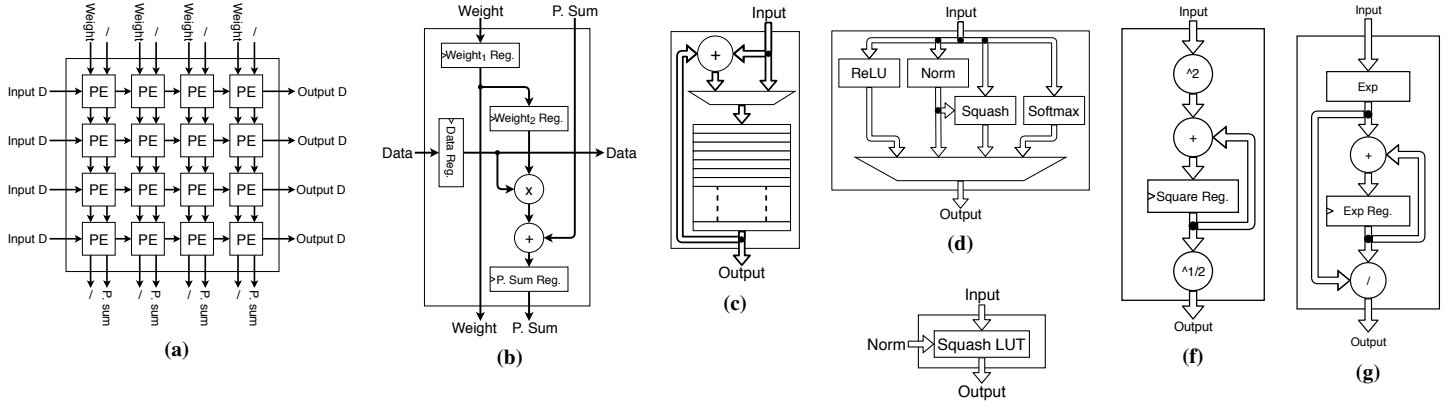
**Fig. 11:** Architecture of Different Components of our CapsAcc Accelerator: (a) Systolic Array. (b) A Processing Element of the Systolic Array. (c) Accumulator. (d) Activation Unit. (e) Squashing Function Unit. (f) Norm Function Unit. (g) Softmax Function Unit.

Figure 11b shows the data path of a single Processing Element (PE). It has 3 inputs and 3 outputs: Data, Weight and Partial sum, respectively. The core of the PE is composed of the sequence of a multiplier and an adder. As shown in Figure 11b, it has 4 internal registers: (1) *Data Reg.* to store and synchronize the Data value coming from the left; (2) *Sum Reg.* to store the Partial sum before sending it to the neighbor PE below; (3) *Weight$_1$ Reg.* synchronizes the vertical transfer; (4) *Weight$_2$ Reg.* stores the value for data reuse. The latter is particularly useful for convolutional layers, where the same weight of the filter must be convolved across different data. For fully-connected computations, the second weight register introduces just one clock cycle latency, without affecting the throughput. The bit-widths of each element have been designed as follows: (1) each PE computes the product between an 8-bit fixed-point Data and an 8-bit fixed-point Weight; and (2) the sum is designed as a 25-bit fixed-point value. At full throttle, each PE produces one output-per-clock cycle, which also implies one output-per-clock cycle for every column of the systolic array.

### B. Accumulator

The Accumulator unit consists of a FIFO buffer to store the Partial sums coming from the systolic array, and sum them together when needed. The multiplexer allows the choice to feed the buffer with the data coming from the systolic array or with the one coming from the internal adder of the Accumulator. We designed the Accumulator to have 25-bit fixed-point data. Figure 11c shows the data path of our Accumulator. In the overall CapsAcc there are as many Accumulators as the number of columns of the systolic array.

### C. Activation Unit

The Activation Unit follows the Accumulators. As shown in Figure 11d, it performs different functions in parallel, while the multiplexer (placed at the bottom of the figure) selects the path to propagate the information towards the output. As for the case of the Accumulator, the figure shows only one unit, while in the complete CapsAcc architecture there is one Activation Unit per each column of the systolic array. The 25-bits data values coming from the Accumulators are reduced to an 8-bit fixed-point value, to reduce the computations at this stage.

Note: the Rectified Linear Unit (ReLU) [10] is a very simple function and its implementation description is omitted, since it

is straightforward. This function is used for every feature of the first two layers of the CapsuleNet.

We designed the **Normalization operator (Norm)** with a structure similar to the Multiply-and-Accumulate operator, where, instead of a traditional multiplier, there is the *Power2* operator. Its data path is shown in Figure 11f. A register stores the partial sum and the *Sqrt* operator produces the output. We designed the square operator as a Look Up Table with 12-bit input and 8-bit output. It produces a valid output every $n+1$ clock cycles, where $n$ is the size of the array for which we want to compute the Norm. This operator is used either as it is to compute the classification prediction, or as an input for the Squashing function.

We designed and implemented the **Squashing function** as a Look Up Table, as shown in Figure 11e. Looking at Equation (1), the function takes an input $s_j$ and its norm $||s_j||$. The Norm input is coming from its respective unit. Hence, this Norm operation is not implemented again inside the Squash unit. The LUT takes as input a 6-bit fixed-point data and a 5-bit fixed-point norm to produce an 8-bit output. We decided to limit the bit-width to reduce the computational requirements at this stage, following the analysis performed in Section III that shows the highest computational load for this operation. A valid output is produced with just one additional clock cycle compared to the Norm.

The **Softmax function** design is shown in Figure 11g. First, it computes the exponential function (8-bit Look Up Table) and accumulates the sum in a register, followed by division. Overall, having an array of $n$ elements, this block is able to compute the softmax function of the whole array in $2n$ clock cycles.

### D. Control Unit

At each stage of the inference process, it generates different control signals for all the components of the accelerator architecture, according to the operations needed. It is essential for correct operation of the accelerator.

## V. DATA-FLOW MAPPING

In this section, we provide the details on how to map the processing of different types of layers and operations onto our CapsAcc accelerator, in a step-by-step fashion. To feed the systolic array, we adopt the mapping policy described in Figure 13. For the ease of understanding, we illustrate the process with the help of an example performing MNIST classification on our CapsAcc accelerator, which also represents our case study. Note,
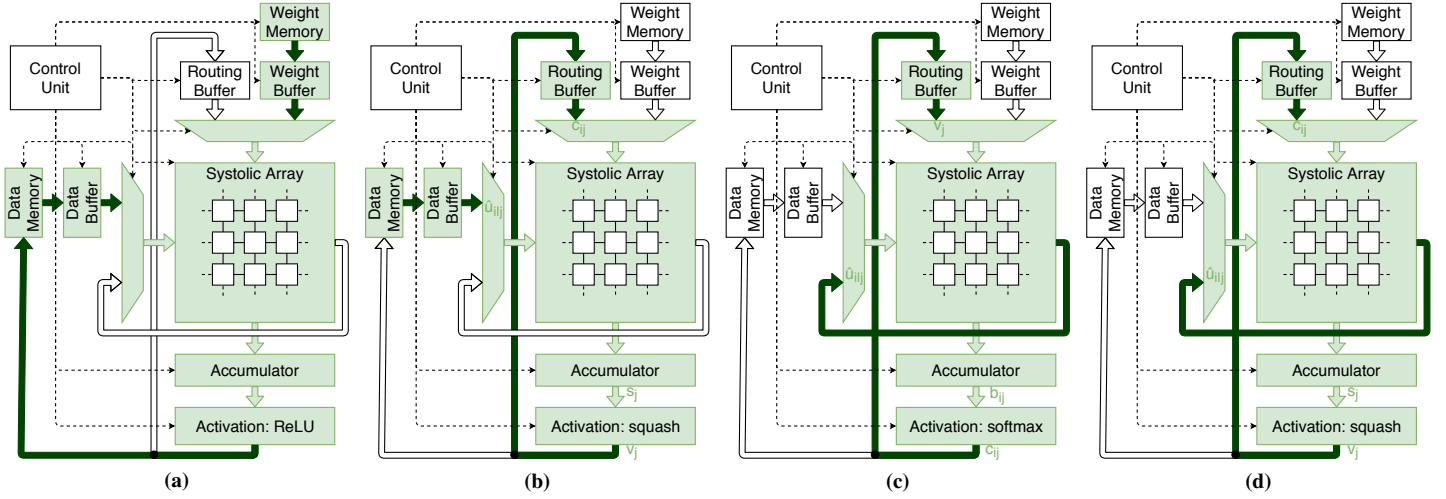
**Fig. 12:** Data-flow mapping onto our CapsAcc accelerator for different scenarios of the case study. (a) Convolutional layer mapping. (b) First sum generation & squashing operation mapping. (c) Update and softmax operation mapping. (d) The sum generation & squashing operation mapping other than the first routing iteration.

```
1: for(l=0; l<L; l++) //output capsules
2:   for(k=0; k<K; k++) //output channels
3:     for(j=0; j<J; j++) //input capsules
4:       for(i=0; i<I; i++) //input channels
5:         for(g=0; g<G; G++) //output columns in a feature map
6:           for(f=0; f<F; f++) //output rows in a feature map
7:             for(c=0; c<C; c++) //kernel/input columns
8:               for(r=0; r<R; r++) //kernel/input rows
9:                 Sum += Weight·Data //multiply and accumulate
```

**Fig. 13:** Mapping algorithm for CapsuleNet operations onto the systolic array.
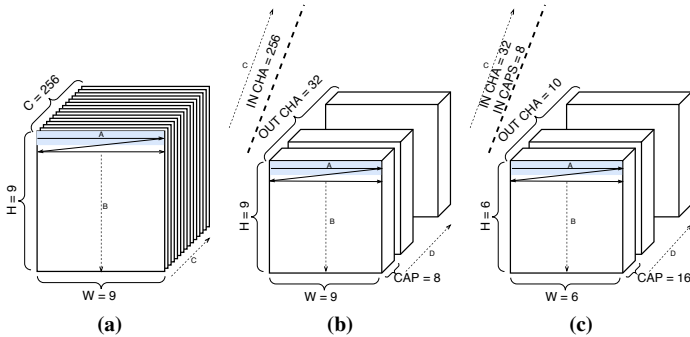


**Fig. 14:** Overview of the process of mapping different layers onto our CapsAcc accelerator. (a) Conv1 layer. (b) PrimaryCaps layer. (c) ClassCaps layer.

each stage of the CapsuleNet inference requires its own mapping scheme.

### A. Conv1 mapping

The Conv1 layer has filters of size 9×9 and 256 channels. As shown in Figure 14a, we designed the mapping row by row (A,B), and after the last row we move to the next channel (C). Figure 12a shows how the data-flow is mapped onto our CapsAcc accelerator. To perform the convolution efficiently, we hold the weight values into the systolic array to reuse the filter across different input data.

### B. PrimaryCaps mapping

Compared to the Conv1 layer, the PrimaryCaps layer has one more dimension, which is the capsule size (i.e., 8). However, we treat the 8D capsule as a convolutional layer with 8 output

channels. Thus, Figure 14b shows that we map the parameters row-by-row (A,B), then moving through different input channels (C), and only at the third stage we move on to the next output channel (D). This mapping procedure allows us to minimize the accumulator size, because our CapsAcc accelerator computes first the output features for the same output channel. Since the type of this layer is convolutional, the data-flow is the same as the one in the previous layer, as reported in Figure 12a.

### C. ClassCaps mapping

The mapping of the ClassCaps layer is shown in Figure 14c. After mapping row by row (A,B), we consider input capsules and input channels as the third dimension (C), and output capsules and output channels as the fourth dimension (D).

Then, for each step of the routing-by-agreement process, we design the corresponding data-flow mapping. It is a critical phase, because a less efficient mapping can potentially have a huge impact on the overall performance.

First, *we apply an algorithmic optimization on the routing-by-agreement algorithm.* During the first operation, instead of initializing $b_{ij}$ to 0 and computing the *softmax* on them, we directly initialize the coupling coefficients $c_{ij}$. The starting point is indicated with the blue arrow in Figure 4. With this optimization, we can skip the *softmax* computation at the first routing iteration. In fact, this operation is dummy, because all the inputs are equal to 0, thus they do not depend on the current data.

Regarding the data-flow mapping in our CapsAcc accelerator, we can identify three different data-flow scenarios during the routing-by-agreement algorithm:

1) **First sum generation and squash:** The predictions $\hat{u}_{j|i}$ are loaded from the Data Buffer, the coupling coefficients $c_{ij}$ are coming from the Routing Buffer, the systolic array computes the sums $s_j$, the Activation Unit computes and selects Squash, and the outputs $v_j$ are stored back in the Routing Buffer. This data-flow is shown in Figure 12b.

2) **Update and softmax:** The predictions $\hat{u}_{j|i}$ are reused through the horizontal feedback of the architecture, $v_j$ are coming from the Routing Buffer, the systolic array computes the updates for $b_{ij}$, and the Softmax at the Activation Unit produces $c_{ij}$ that are stored back in the Routing Buffer. Figure 12c shows the data-flow described above.
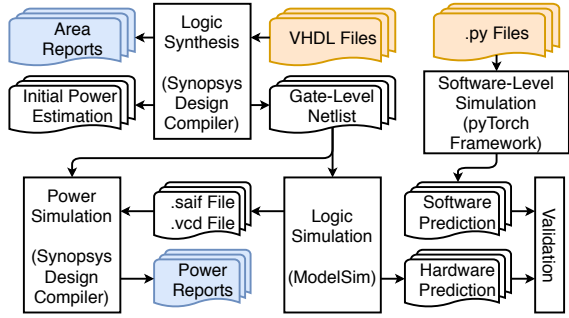
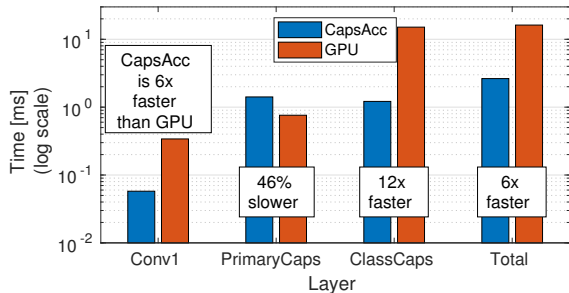**Fig. 15:** Synthesis flow and tool chain of our experimental setup.



**Fig. 17:** Performance of the inference pass on each step of the routing-by-agreement algorithm on our CapsAcc accelerator, compared to the GPU.



**Fig. 16:** Layer-wise performance of the inference pass on the CapsuleNet on our CapsAcc accelerator, compared to the GPU.

| Tech. node [nm] | 32 |
|---|---|
| Voltage [V] | 1.05 |
| Area [mm$^2$] | 2.90 |
| Power [mW] | 202 |
| Clk Freq. [MHz] | 250 |
| Bit width | 8 |
| On-Chip Mem. [MB] | 8 |

**TABLE II:** Parameters of our synthesized CapsAcc accelerator

| Component | Area [um^2] | Power [mW] |
|---|---|---|
| Accumulator | 311961 | 22.80 |
| Activation | 143045 | 5.94 |
| Data Buffer | 1332349 | 95.96 |
| Routing Buffer | 316226 | 22.78 |
| Weight Buffer | 115643 | 8.34 |
| Systolic Array | 680525 | 46.09 |
| Other | 4330 | 0.13 |

**TABLE III:** Area and power, measured for the different components of our CapsAcc accelerator.

**3) Sum generation and Squash:** Figure 12d shows the data-flow mapping for this scenario. Compared to the Figure 12b, the predictions $\hat{u}_{j|i}$ are coming from the horizontal feedback link, thus exploiting data reuse also in this stage.

## VI. RESULTS AND DISCUSSION

### A. Experimental Setup

We implemented the complete design of our CapsAcc architecture in RTL (VHDL), and evaluated it for the MNIST dataset *(to stay consistent with the original CapsuleNet paper)*. We synthesized the complete architecture in a 32nm CMOS technology library using the ASIC design flow with the Synopsys Design Compiler. We did functional and timing validation through the gate-level simulations using ModelSim, and obtained the precise area, power and performance of our design. The complete synthesis flow is shown in Figure 15, where the orange and blue colored boxes represent the inputs and the output results of our experiments, respectively.

**Important Note:** since our hardware design is fully functionally compliant with the original CapsuleNet design of the work of [12], we observed the same accuracy of classification. Therefore, we do not present any classification results in this paper, and only focus on the performance, area and power results, which are more relevant for an optimized hardware architecture.

### B. Discussion on Comparative Results

The graphs shown in Figure 16 report the performance (execution time) results of the different layers of CapsuleNet inference on our CapsAcc accelerator, while Figure 17 shows the performance of every sequence of the routing process. Compared with the GPU performance (see Figures 8 and 9), we obtained a significant speed-up for the overall computation time of a CapsuleNet inference pass (6×). *The main notable improvements*
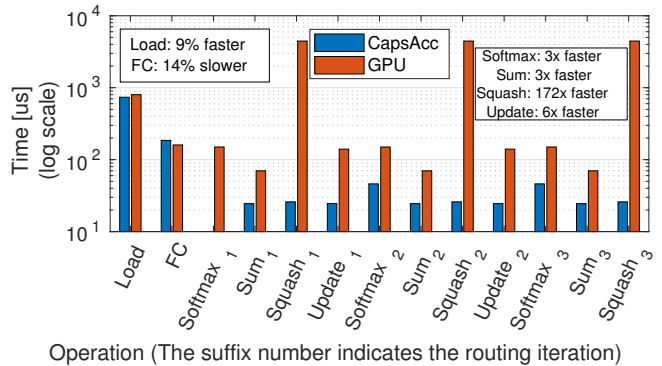
*are witnessed in the ClassCaps layer (12×) and in the Squashing operation (172×).*

### C. Detailed Area and Power Breakdown

The details and synthesis parameters for our design are reported in Table II. Table III shows the absolute values for the area and power consumption of all the components of the synthesized accelerator. Figures 18a and 18b show the area and power breakdowns, respectively, of our CapsAcc architecture. These figures show that the area and power contributions are dominated by the buffers, and the systolic array is just 1/4 of the total budget.

## VII. CONCLUSIONS

We presented the first CMOS-based hardware accelerator for the complete CapsuleNet inference. To achieve high performance, our CapsAcc architecture employs a flexible systolic array with several optimized data-flow patterns that enable it to fully exploit a high level of parallelism for diverse operations of the CapsuleNet processing. To efficiently use the proposed hardware design, we also optimized the routing-by-agreement algorithm without changing its functionality and thereby preserving the classification accuracy of the original CapsuleNet design of [12]. Our results show a significant speedup compared to an optimized GPU implementation. We also presented power and area breakdown of our hardware design. Our CapsAcc provides the first proof-of-concept for realizing CapsuleNet hardware, and opens new avenues for its high-performance inference deployments.

## REFERENCES

[1] Y. H. Chen, J. Emer, and V. Sze. Eyeriss: A spatial architecture for energy efficient dataflow for convolutional neural networks. In *ISCA*, 2016.

[2] R. Collobert, J. Weston, L. Bottou, M. Karlen, K. Kavukcuoglu, and P. Kuksa. Natural language processing (almost) from scratch. In *JMLR*, 2011.
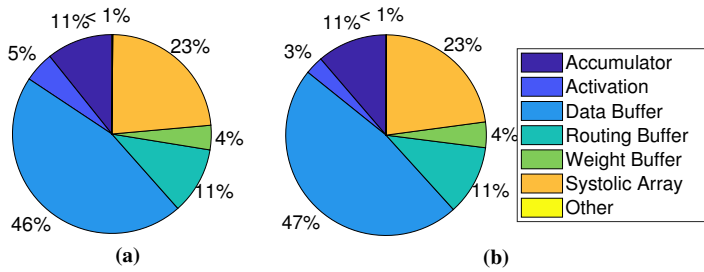
**Fig. 18:** (a) Area and (b) Power Breakdown of our CapsAcc Accelerator.

[3] A. Graves and J. Schmidhuber. Framewise phoneme classification with bidirectional lstm and other neural network architectures. In *Neural Networks*, 2005.

[4] S. Han, X. Liu, H. Mao, J. Pu, A. Pedram, M. A. Horowitz, and W. J. Dally. EIE: Efcient Inference Engine on Compressed Deep Neural Network. In *ISCA*, 2016.

[5] G. E. Hinton, A. Krizhevsky, and S. D. Wang. Transforming auto-encoders. In *ICANN*, 2011.

[6] N. P. Jouppi et al. In-datacenter performance analysis of a tensor processing unit. In *ISCA*, 2017.

[7] A. Krizhevsky, I. Sutskever, and G. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012.

[8] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. In *Proceedings of the IEEE*, 1998.

[9] W. Lu, G. Yan, J. Li, S. Gong, Y. Han, and X. Li. Flexflow: A flexible dataflow accelerator architecture for convolutional neural networks. In *HPCA*, 2017.

[10] V. Nair and G. E. Hinton. Rectified linear units improve restricted boltzmann machines. In *ICML*, 2010.

[11] A. Parashar, M. Rhu, A. Mukkara, A. Puglielli, R. Venkatesan, B. Khailany, J. Emer, S. W. Keckler, and W. J. Dally. SCNN: An accelerator for compressed-sparse convolutional neural networks. In *ISCA*, 2017.

[12] S. Sabour, N. Frosst, and G. E. Hinton. Dynamic routing between capsules. In *NIPS*, 2017.

[13] pyTorch framework: https://github.com/pytorch/pytorch