# Enabling Efficient Processing of Spiking Neural Networks with On-Chip Learning on Commodity Neuromorphic Processors for Edge AI Systems

Rachmad Vidya Wicaksana Putra, Pasindu Wickramasinghe, Muhammad Shafique eBrain Lab, New York University (NYU) Abu Dhabi, Abu Dhabi, UAE {rachmad.putra, pmw6287, muhammad.shafique}@nyu.edu

Abstract—The rising demand for energy-efficient edge AI systems (e.g., mobile agents/robots) has increased the interest in neuromorphic computing, since it offers ultra-low power/energy AI computation through spiking neural network (SNN) algorithms on neuromorphic processors. However, their efficient implementation strategy has not been comprehensively studied, hence limiting SNN deployments for edge AI systems. Toward this, we propose a design methodology to enable efficient SNN processing on commodity neuromorphic processors. To do this, we first study the key characteristics of targeted neuromorphic hardware (e.g., memory and compute budgets), and leverage this information to perform compatibility analysis for network selection. Afterward, we employ a mapping strategy for efficient SNN implementation on the targeted processor. Furthermore, we incorporate an efficient on-chip learning mechanism to update the systems' knowledge for adapting to new input classes and dynamic environments. The experimental results show that the proposed methodology leads the system to achieve low latency of inference (i.e., less than 50ms for image classification, less than 200ms for real-time object detection in video streaming, and less than 1ms in keyword recognition) and low latency of on-chip learning (i.e., less than 2ms for keyword recognition), while incurring less than 250mW of processing power and less than 15mJ of energy consumption across the respective different applications and scenarios. These results show the potential of the proposed methodology in enabling efficient edge AI systems for diverse application use-cases.

*Index Terms*—Neuromorphic computing, spiking neural networks, neuromorphic processors, event-based processing, on-chip learning, edge AI systems, real-world workloads.

## I. Introduction

In recent years, the demand for employing energy-efficient edge AI systems (e.g., mobile agents/robots and IoT devices) has increased rapidly due to their advantages in improving quality of services (QoS) and human productivity. However, these systems are typically powered by portable batteries with limited capacity [1]. Therefore, they usually suffer from short battery lifespan. To address this limitation, employing larger battery capacity may not be a scalable solution, since larger battery capacity means heavier mass, which typically require more power/energy consumption to operate or mobilize the system. Current trends also show that, real-world autonomous systems with heavier mass typically consume higher power/energy than the smaller ones, which leads to shorter battery lifespan; see Fig. 1(a). Therefore, the potential solution is making these systems to employ ultra-low power/energy AI algorithms, thus optimizing the overall systems' energy.

Toward this, neuromorphic computing (NC) with spiking neural networks (SNN) algorithms has emerged as a potential

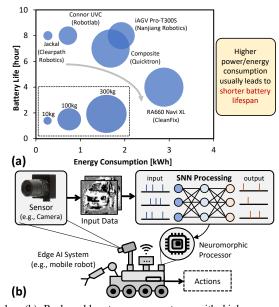


Fig. 1. (b) Real-world autonomous systems with higher power/energy consumption usually have heavier mass and shorter battery lifespan; adapted from studies in [1]. (b) Overview of SNN-based computation for edge AI systems, considering an example from a mobile robot use-case.

solution due to its highly sparse spike-driven operations [2] [3]; see an illustration in Fig. 1(b). To maximize the energy-efficiency benefits from NC, SNN processing needs to be performed on specialized neuromorphic hardware processors, which accommodate spike-driven operations [4]–[6].

Currently, most of the existing neuromorphic processors (e.g., SpiNNAker, NeuroGrid, IBM's TrueNorth, and Intel's Loihi) [4] were developed mainly for research purpose and not commercially available, thereby making it difficult to use them in SNN-based edge AI systems for real-world application usecases. Recently, several neuromorphic processors are released and available commercially in the market (such as BrainChip's Akida [7] and SynSense's DYNAP [8]), which can be used for developing real-world SNN-based edge AI systems. However, their efficient implementation strategy has not been studied, hence limiting the systems from achieving further efficiency gains considering different SNNs and workloads. Therefore, in this paper, the targeted research problem is: How can we enable efficient execution of SNN models on commodity neuromorphic processors, while achieving good trade-offs between accuracy, memory, and power/energy consumption? A solution to this problem may enable energy-efficient SNN deployments for diverse application use-cases at the edge.

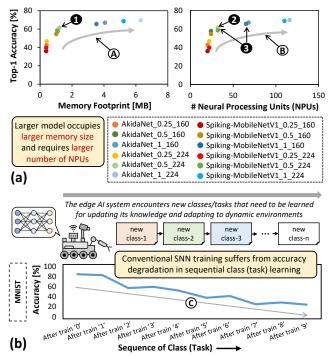


Fig. 2. (a) Results of running different SNN models considering different input resolutions on the same commodity neuromorphic processor (i.e., Akida); based on data from [17]. Here, each network name denotes "Network\_Alpha\_InputResolution", where Alpha represents the width multiplier that shrinks the network uniformly from the original size. (b) An edge AI system may encounter new classes at run time that need to be learned for updating its knowledge. Conventional SNN training suffers from accuracy degradation in sequential class learning; based on studies in [16].

#### A. State-of-the-Art and Their Limitations

State-of-the-art works in employing commodity neuromorphic processors for edge AI systems typically focus on the implementation of a specific application use-case, such as system control [9] [10], tactile sensing [11], gait analysis [12], object detection [13]–[15]. Moreover, the state-of-the-art works have not studied the on-chip learning aspect, which is important for updating the systems' knowledge and adapting to changing environments [16]. This condition shows that, a design methodology for enabling efficient SNN execution considering different workloads on commodity neuromorphic processors has not been explored. To show the importance of such a methodology, we conduct a case study in Section I-B.

# B. Case Study and Research Challenges

We aim to observe the compute and memory requirements for running different SNN models on the same commodity neuromorphic processor. Here, we consider different SNN models (i.e., AkidaNet and Spiking-MobileNetV1¹) with different input resolutions (i.e., 160x160 and 224x224) that run on the Akida neuromorphic processor². The investigation results are shown in Fig. 2(a). They show that, running larger SNN model on the neuromorphic processor potentially offer

higher accuracy due to higher feature extraction capabilities, but at the cost of larger memory footprint and larger compute resource (i.e., number of Neural Processing Units), and hence higher power/energy consumption; see (A)-(B). Such increased compute, memory, and power/energy requirements can reduce the efficiency gains of SNN-based edge AI systems.

Edge AI systems may also encounter new class at run time, that need to be learned for updating the systems' knowledge. Otherwise, these systems can suffer from accuracy degradation when their knowledge becomes obsolete over time [18], as the offline-trained SNN may struggle in training the new classes while preserving the old ones (i.e., previously learned classes) [19]; see © in Fig. 2(b)

**Required:** A design methodology that identifies the compatibility of the selected network for the targeted neuromorphic processor, and facilitates on-chip learning mechanism for knowledge updates. However, this requirement exposes several research challenges, as described in the following.

- The memory and compute costs of the selected network should be efficiently accommodated by the processor.
- The selected SNN model needs to be efficiently mapped on the processor to ensure energy-efficient SNN execution.
- The system demands an efficient on-chip learning capability to learn new classes, while preserving the knowledge of previously learned ones.

#### C. Our Novel Contributions

To address the targeted problem and associated challenges, we propose a novel design methodology that ensures efficient execution of SNNs on commodity neuromorphic processors, thereby enabling energy-efficient deployments of SNN-based edge AI systems for diverse applications. Our novel contributions include the following points (an overview in Fig. 3).

- A design methodology (Section III) that employs several key steps, as summarized below.
- 1) **Network compatibility analysis** (Section III-A) that quickly evaluates whether the selected network can be efficiently executed, i.e., by leveraging the characteristics of neuromorphic processor (e.g., memory and compute budgets) and the proposed analytical model.
- 2) Efficient deployment on the processor (Section III-B) by employing an efficient mapping strategy considering good trade-off between hardware costs and efficiency through the processor runtime settings.
- 3) On-chip learning strategy (Section III-C) that facilitates learning new classes after SNN deployment on the processor through a last layer modification technique.
- Comprehensive evaluation (Section V) that covers multiple design metrics, e.g., accuracy, latency, throughput, and power/energy consumption of different SNNs in different applications (i.e., image classification, robject recognition in video streaming, and keyword recognition) and different scenarios (i.e., offline-based and on-chip learning settings).

**Key Results:** Our design methodology is evaluated through a real-world edge AI prototype using the Akida neuromorphic

<sup>&</sup>lt;sup>1</sup>AkidaNet is a MobileNetV1-inspired network optimized for deployment on Akida [17], while Spiking-MobileNetV1 is the converted MobileNetV1 in spiking domain.

<sup>&</sup>lt;sup>2</sup>The details of Akida hardware architecture are provided in Section II-B, while the details of experimental setup is presented in Section IV.

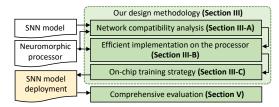


Fig. 3. Overview of our novel contributions, highlighted in green.

processor. The experimental results show that, our methodology leads the system to achieve low latency of inference (i.e., less than 50ms for image classification, less than 200ms for real-time object detection in video streaming, and less than 1ms in keyword recognition) and low latency of on-chip learning (i.e., less than 2ms for keyword recognition), while consuming less than 250mW of power.

#### II. BACKGROUND

# A. Spiking Neural Networks (SNNs)

**Overview:** SNNs are considered the bio-plausible neural network (NN) models [3], since they are modeled after the neural process observed in the human brain, specifically on how neurons utilize spikes for transferring and processing data. An SNN model mainly consists of several components: spiking neuron, synapses, network topology/architecture, and neural encoding [20]; see Fig. 4(a).

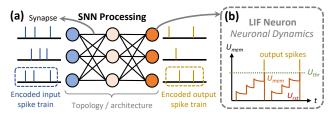


Fig. 4. (a) Illustration of an SNN and its components. (b) Overview of the neuronal dynamics of the widely-used spiking neuron model (i.e., LIF).

**Spiking Neuron:** The dynamics of spiking neuron depend on the neuron model, and the widely-used one is the Leaky Integrate-and-Fire (LIF) neuron [21]. The neuronal dynamics of LIF is illustrated in Fig. 4(b). Here,  $U_{mem}$ ,  $U_{thr}$ , and  $U_{rst}$  denote the neurons' membrane potential, threshold potential, and reset potential, respectively. When an incoming spike arrives in the LIF neuron, it triggers the increasing of  $U_{mem}$ ; otherwise,  $U_{mem}$  decays. If the  $U_{mem}$  reaches or surpasses the  $U_{thr}$ , then an output spike is generated.

# B. Neuromorphic Processors

1) Overview: The energy efficiency potentials offered by SNNs can be maximized by employing neuromorphic hardware processors [22]. In the literature, several processors have been proposed, and they can be categorized as research and commodity processors. Research processors refer to neuromorphic chips that are designed only for research and not commercially available, hence access to these processors is limited. Several examples in this category are SpiNNaker, NeuroGrid, IBM's TrueNorth, and Intel's Loihi [4]. Meanwhile, commodity processors refer to neuromorphic chips that are available commercially, such as BrainChip's Akida [7] and

SynSense's DYNAP-CNN [8]. In this work, we consider the Akida processor as it supports on-chip learning for SNN fine-tuning, which is beneficial for adaptive edge AI systems [7].

2) Akida Neuromorphic Processor System-on-Chip: Akida Neuromorphic SoC (NSoC) is designed by BrainChip, which aims at accelerating SNN processing for low-power application use-cases [7]. Its commercially available version is the Akida v1.0 (AKD1000) [7], which is fabricated using the TSMC's 28nm technology, and it can run at 300MHz clock frequency. The overview of Akida NSoC architecture is shown in Fig. 5. It mainly consists of a Cortex-M4 CPU as the SoC host processor and 80 Akida Neural Processing Units (NPUs) as the SNN processors/cores. Four NPUs form a single node, hence forming 20 nodes. Each NPU mainly consists of 8 Neural Processing Engines (NPEs) as the compute units for executing synaptic and neuronal operations (e.g., event-based convolutions), and 100KB SRAM buffer as the local memory for storing weights (40KB) and data spikes (60KB) [23]. Akida NSoC has a direct memory access (DMA) controller, a power management unit (PMU), several data interfaces (i.e., USB 3.0, PCIe 2.1, I2S, I3C, UART, and JTAG), two memory interfaces (i.e., SPI Flash and LPDDR4), and an interface for multi-chip expansion. For data encoding, Akida employs the data/pixelspike converter. Furthermore, to facilitate SNN developments and their Akida implementations, BrainChip provides MetaTF framework [17], which accommodates features like ANN-to-SNN conversion, SNN mapping, and on-chip learning.

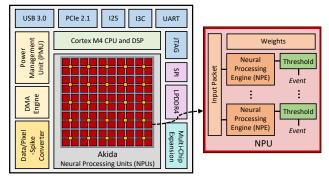


Fig. 5. Overview of the Akida architecture; adapted from [7] [23].

#### III. OUR DESIGN METHODOLOGY

We propose a design methodology to address the targeted problem and related challenges, whose key steps are shown in Fig. 6, and discussed in Sections III-A until III-C. For an overview, we describe the flow of our methodology as follows.

- 1) It starts with the network development using the existing environment (e.g., TensorFlow+Keras) to provide offline-trained NN models.
- 2) The offline-trained NN models and the processor configuration are leveraged to perform the network compatibility analysis by using our analytical model for selecting an NN model that meets the memory and compute budgets.
- 3) If the selected NN model is not in spiking domain, we convert the ANN model into an SNN model. Otherwise, the model is already in spiking domain and ready to use.

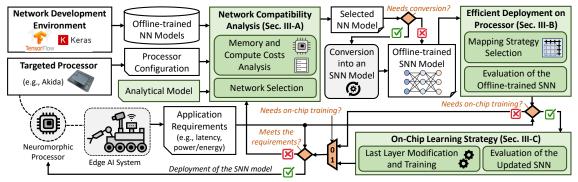


Fig. 6. Our proposed design methodology, showing the novel contributions highlighted in green boxes.

- 4) The SNN model is then deployed on the processor using a specific mapping strategy, and then evaluated with a specific workload (e.g., image classification).
- If a knowledge update is needed, we perform an on-chip learning on the deployed SNN model, and then evaluate it. Otherwise, the original SNN is evaluated.
- 6) If the evaluated SNN meets the application requirements (e.g., latency and power/energy), then it can be deployed on the processor. Otherwise, we can select a smaller NN model through the network compatibility analysis step.

## A. Network Compatibility Analysis

This step aims at analyzing whether the selected network can be executed efficiently in the targeted neuromorphic processors. Specifically, this refers to the condition where the network can be fully mapped and executed on the processor at one time, and hence no network partitioning and scheduling are required. In this manner, costly memory access and data movements can be minimized, as these operations typically dominate the neuromorphic systems' energy [24]. This is important because it evaluates the processing requirements in advance before the actual deployment on the hardware, hence guiding the users to better develop and/or select a suitable network to deploy on the targeted processor.

Memory and Compute Costs Analysis: We identify the characteristics of targeted processor that determine whether the selected network can be fully mapped and executed on the processor at one time, hence avoiding network partitioning and scheduling. Specifically, we investigate the memory and compute budgets, and how they are distributed in the targeted processor. The memory budget represents the maximum size of network and activation data that can be fully mapped at one time, while the compute budget represents the maximum event-based computations that can be executed at one time, hence they are both leveraged for the network compatibility analysis. To enable this analysis, we propose an analytical model to estimate the memory and compute costs for the given network, while considering the hardware architecture from the processor. Specifically, it investigates the total number of NPUs (cores) required to fully map and execute the given network and data (i.e., denoted as  $N_{NPU\ tot}$ ).

**Poposed Analytical Model:**  $N_{NPU\_tot}$  is defined as a total number of NPUs required across different layers of the given network; see Eq. 1. Here, L represents the number of layers in

the network, and  $N_{NPU\_mem}^l$  represents the number of NPUs required for storing network parameters and data in layer-l.

$$N_{NPU\_tot} = \sum_{l=1}^{L} N_{NPU\_mem}^{l} \tag{1}$$

We observe that, SNN parameters and data in the same layer may have different sizes, hence requiring different numbers of NPUs. To properly allocate hardware resources for such a condition, we select the bigger number of NPUs to ensure sufficient memory resource, whose function can be can stated as Eq. 2. Here,  $N_{NPU\_net}^l$  and  $N_{NPU\_dat}^l$  denote the number of NPUs for network parameters and data in layer-l, respectively.

$$N_{NPU\_mem}^l = \max(N_{NPU\_net}^l, N_{NPU\_dat}^l) \qquad (2)$$

We can obtain  $N_{NPU\_net}^l$  and  $N_{NPU\_dat}^l$  using Eq. 3. Here,  $M_{net}^l$  and  $M_{dat}^l$  denote the size of network parameters and activation data in layer-l, respectively. Meanwhile,  $B_{net}$  and  $B_{dat}$  denote the local memory (buffer) size in each NPU for network parameters and data, respectively.

$$N_{NPU\_net}^{l} = \left\lceil \frac{M_{net}^{l}}{B_{net}} \right\rceil$$
 and  $N_{NPU\_dat}^{l} = \left\lceil \frac{M_{dat}^{l}}{B_{dat}} \right\rceil$  (3)

Furthermore,  $M_{net}^l$  can be obtained by leveraging the number of parameters (i.e., weights  $N_w$  and bias  $N_b$ ) with their bit precision  $(bit_{par})$  in layer-l; see Eq. 4. Meanwhile,  $M_{dat}^l$  can be obtained by leveraging the number of feature maps with their bit precision  $(bit_{dat})$  in layer-l; see Eq. 5. Note,  $H^l$ ,  $W^l$ , and  $C^l$  denote the feature maps' dimension in layer-l for height, width, and channel, respectively.

$$M_{net}^l = (N_w^l + N_b^l) \cdot bit_{par}^l \tag{4}$$

$$M_{dat}^{l} = (H^{l} \cdot W^{l} \cdot C^{l}) \cdot bit_{dat}^{l} \tag{5}$$

**Network Selection:** We use the proposed analytical model to identify the network models that can be efficiently executed in the targeted processor, i.e., by selecting the network models whose memory and compute costs  $(N_{NPU\_tot})$  are less than the memory and compute budgets from the processor  $(N_{NPU\_proc})$ , while offering high accuracy.

# B. Efficient Deployment on the Processor

This step aims to enable efficient deployment of the selected network on the processor. It requires a mapping strategy that

leads the SNN processing to meet the application requirements (e.g., latency and power). In practice, the possible strategies also depend on the availability of related application programming interface (API) of the chip's implementation framework.

**Mapping Strategy:** In this work, we employ a mapping strategy from the Akida's MetaTF framework [17] that maximizes hardware resources for holding network parameters and data with minimum passes (i.e., sequential processing), thus providing a trade-off between performance and efficiency. This strategy aims at mapping the entire SNN model and data in the NPU local memories, while employing a sequential processing for synaptic and neuronal operations in each layer, minimizing parallel NPU processing. Therefore, the optimization objective is to minimize the number of NPUs for executing synaptic and neuronal operations for each layer  $(N_{NPU}^l_{exe})$ ; see Eq. 6.

Objective: minimize(
$$N_{NPU\ exe}^{l}$$
) (6)

The cost function for hardware mapping (C) is defined as the total NPU allocation for storing network parameters and data  $(N_{NPU\_mem})$  and executing synaptic and neuronal operations (optimized  $N_{NPU\_exe}$ ) across all layers (L). To allocate resources for such a condition, we select the bigger number of NPUs for facilitating each layer processing; see Eq. 7.

$$C = \sum_{l=1}^{L} \{ \max(N_{NPU\_mem}^{l}, N_{NPU\_exe}^{l}) \}$$
 (7)

Consequently, the hardware utilization (U) can be determined through the ratio between the mapping cost C and the total number of NPUs in the processor  $N_{NPU\_proc}$ ; see Eq. 8.

$$U = \frac{C}{N_{NPU\ proc}} \cdot 100\% \tag{8}$$

**Network Evaluation:** After mapping the selected SNN on the processor, we evaluate its performance and efficiency to observe if the selected SNN meets the given requirements. If the systems need to update their knowledge, then this SNN model needs to be updated through an on-chip learning.

# C. On-Chip Learning Strategy

It aims at *learning new classes on-chip, thereby enabling* an efficient fine-tuning for the existing SNN model. Here, possible on-chip learning strategies for neuromorphic processors mainly also depend on the availability of related API from their development framework.

**On-chip Learning:** For a case study, we use the available on-chip learning strategy from the Akida's MetaTF framework [17]. To enable the on-chip learning, we need to fulfill the following learning constraints for the last network layer, since this last layer is the only part that will be trained in on-chip learning: (1) it must be a fully connected type, (2) it must have binary weights, and (3) it must receive binary inputs. We fulfill these requirements through the following steps.

• We replace the last layer with a new layer that meets the learning constraints/characteristics.

- The new layer should accommodate classifying both the old and new classes. Hence, multiple neurons for each old class are used to provide spaces for learning new classes on-chip.
- We perform a few-shot learning on-chip with a few samples for each new class using the Akida's built-in algorithm. The new classes are associated with specific neurons in last layer.

**Network Evaluation:** After mapping the updated SNN on the processor, we evaluate its performance and efficiency to observe if the updated SNN meets the given requirements.

#### IV. EVALUATION METHODOLOGY

To evaluate our proposed design methodology, we employ the experimental setup and tools flow presented in Fig. 7.

**Software Development Part:** We employ MetaTF framework [17], which is based on TensorFlow and Keras libraries, to convert the pre-trained NN model into an SNN model. Afterward, this SNN model is mapped on the neuromorphic processor and executed accordingly. For the on-chip learning, the SNN model is modified to facilitate learning new classes; as discussed in Section III-C. In experiments, we record results like accuracy, latency, and power/energy consumption.

Hardware Development Part: We develop a real-world edge AI system, comprising neuromorphic processor and host CPU. For the neuromorphic processor, we employ a single Akida NSoC (Akida v1.0 AKD1000) [7]. Meanwhile, for the host CPU, we employ an ARM Cortex-A72 through the Raspberry Pi Compute Module 4 (CM4), which runs the Ubuntu 22.04 OS. These host and neuromorphic processors are connected through the PCIe interface, thereby providing a high-speed serial computer expansion bus.

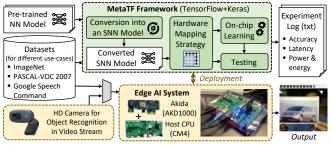


Fig. 7. Experimental setup and tools flow.

**Application Use-Cases:** For showing the generality of our design methodology, we consider three different applications.

- 1) Classification of Static Images: It considers the ImageNet dataset [25]. Its application requirements include the maximum 50ms latency and 250mW power consumption.
- 2) Real-time Object Recognition in Video Streaming: It considers the PASCAL-VOC 2007 dataset [26]. We perform real-time object detection in video streaming using a complete edge AI system utilizing a Logitech C270 HD WebCam; see Fig. 7. Its application requirements include the maximum 200ms latency and 250mW power consumption.
- 3) *Keyword Recognition:* It uses the Google Speech Command dataset [27]. Its application requirements include the maximum 5ms latency and 250mW power consumption.

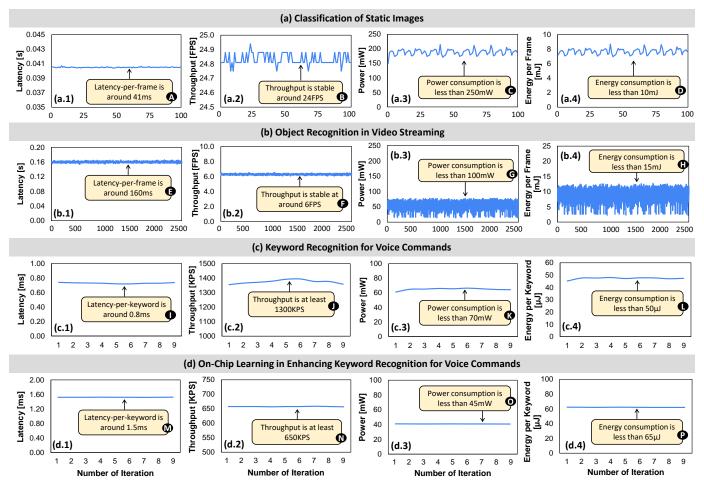


Fig. 8. Experimental results of running SNN models on Akida for (a) classification of static images, (b) real-time object detection in video stream, (c) keyword recognition for voice commands, and (d) on-chip learning in enhancing keyword recognition for voice commands; encompassing the latency, throughput, as well as power and energy consumption. Note, FPS denotes the frame-per-second, while KPS denotes the keyword-per-second.

# V. RESULTS AND DISCUSSION

# A. Classification of Static Images

**Network Selection:** Our network compatibility analysis (discussed in Section III-A) leads to the selection of AkidaNet\_0.5\_224 from many possible network models due to the following reasons.

- The network size meets the memory budget of an Akida chip (i.e., 8MB), as shown by **1** in Fig. 2(a).
- The number of NPUs required for a complete computation of the network meets the NPU budget of an Akida chip (i.e., 80 NPUs), as shown by ② in Fig. 2(a).
- The network achieves higher accuracy as compared to other network models, but slightly lower than AkidaNet\_1\_160 and Spiking-MobileNetV1\_1\_160, as shown by 3 in Fig. 2(a). We select AkidaNet\_0.5\_224 as it can handle higher input resolution 224x224, which is beneficial for systems with high resolution sensors.

**Accuracy:** Here, we perform inference by presenting 10 images from the ImageNet to the system over 100 iterations. The experimental results show that, running AkidaNet\_0.5\_224 on the Akida processor achieves 80% accuracy. This high accuracy comes from an effective training process in ANN domain

that employs the accurate backpropagation technique, and an effective conversion technique that accurately translates the trained ANN components and parameters into representative SNN components and parameters.

Latency, Throughput, Power and Energy Consumption: The results for latency and throughput are shown in Fig. 8(a.1) and Fig. 8(a.2), respectively. Latency is stable around 41ms across 100 iterations of experiments (see A), which leads to 24FPS throughput (see **B**), thereby meeting the design requirement of maximum 50ms latency. Such low latency and high throughput mainly come from the selected mapping strategy, which considers maximizing hardware resources to fully map the entire SNN model on the Akida's NPU fabrics, hence avoiding the time-consuming execution of network partitions. Meanwhile, the results for power and energy consumption are presented in Fig. 8(a.3) and Fig. 8(a.4), respectively. Overall, power consumption is about 215mW (see **©**), and energy consumption is about 9mJ (see **D**), thereby meeting the design requirement of maximum 250mW power. Such low power and low energy consumption come from the sparse spike-driven computation, and the selected mapping strategy that optimizes data movements through efficient NPU allocation, and hence minimizing power consumption for the respective operations.

#### B. Real-Time Object Detection in Video Streaming

**Network Selection:** Our compatibility analysis (from Section III-A) leads to the selection of Spiking-YOLOv2 [17] due to the following reasons.

- The network size (i.e.,  $\sim$ 3MB) meets the memory budget of Akida chip (i.e., 8MB).
- The number of NPUs required for a complete computation of the network (i.e., 71 NPUs) meets the NPU budget of an Akida chip (i.e., 80 NPUs).

Accuracy: We perform inference using with 2500 iterations of object presentation (i.e., person and car). Screenshots of the real-time object detection in video streaming are presented in Fig. 9. The experimental results show that, running Spiking-YOLOv2 on the Akida processor can achieve 94.44% accuracy for detecting the presented objects. This high accuracy is due to the training process that utilizes accurate backpropagation in ANN domain, and the conversion process that accurately translates the ANN model into a representative SNN model.

Latency, Throughput, Power and Energy Consumption: The experimental results for latency and throughput are shown in Fig. 8(b.1) and Fig. 8(b.2), respectively. Processing latency is around 160ms across 2500 iterations of experiments (see **(B)**, which leads to 6FPS throughput (see **(B)**), and thereby meeting the design requirement of maximum 200ms latency. These results show that, the system achieves relatively low latency and high throughput for real-time object detection in video streaming. It is due to the selected mapping strategy, which fully maps the entire SNN model on the NPU fabrics, which avoids the time-consuming execution of network partitions. Meanwhile, the experimental results for power and energy consumption are shown in Fig. 8(b.3) and Fig. 8(b.4), respectively. Power consumption is  $\sim$ 78mW (see **G**), and energy consumption is  $\sim$ 13mJ (see **(1)**), thereby meeting the design requirement of maximum 250mW power. Such low power and low energy consumption come from the sparse spike-driven computation, and the selected mapping strategy that optimizes data movements through judicious NPU allocation, hence minimizing the power consumption.

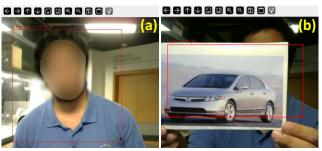


Fig. 9. Experimental results of the edge AI system prototyping based on the Akida neuromorphic processor for object detection in video streaming: (a) person detection, and (b) car detection.

# C. Keyword Recognition for Voice Commands

**Network Selection:** Our compatibility analysis (from Section III-A) leads to the selection of Spiking Depthwise Separable Convolutional Neural Network (Spiking-DSCNN) [17] due to the following reasons.

- The network size (i.e., 23KB) meets the memory budget of Akida chip (i.e., 8MB).
- The number of NPUs required for a complete computation of the network (i.e., 5 NPUs) meets the NPU budget of an Akida chip (i.e., 80 NPUs).

**Accuracy:** Here, we perform inference by presenting 1000 keywords from the Google Speech Command dataset to the system over 9 iterations. The experimental results show that, running Spiking-DSCNN on the Akida processor can achieve 91.73% accuracy. This high accuracy is due to an effective ANN training and its accurate ANN-to-SNN conversion.

Latency, Throughput, Power and Energy Consumption: The experimental results for processing latency and throughput are presented in Fig. 8(c.1) and Fig. 8(c.2), respectively. Here, processing latency is stable around 0.72ms across 9 iterations of experiments (see **1**), which leads to more than 1300KPS (keyword-per-second) throughput (see **①**), thereby meeting the design requirement of maximum 5ms latency. Such low latency and high throughput mainly come from the small size of Spiking-DCNN with 23KB, which makes the entire network easy to map on the NPU fabrics. Consequently, this avoids the time-consuming execution of network partitions, while incurring small data movements and operations. Meanwhile, the experimental results for power and energy consumption are presented in Fig. 8(c.3) and Fig. 8(c.4), respectively. Power consumption is about 68mW (see **(S)**), and energy consumption is about  $49\mu J$  (see **1**), thereby meeting the design requirement of maximum 250mW power. These low power and low energy consumption mainly come from the small size of Spiking-DCNN which makes the entire network model can be efficiently executed on the NPU fabrics.

#### D. On-Chip Learning for Knowledge Updates

We select the keyword recognition application and employ the pre-trained Spiking-DSCNN. From the Google Speech Command dataset, we use 32 keywords for the offline training and 3 new keywords (i.e., 'backward', 'follow', and 'forward') as new classes for on-chip learning.

**Accuracy:** We first perform on-chip learning on the Akida for 3 new keywords, and each one is trained using 160 samples. Then, we perform inference using 6 samples for 'backward', 7 samples for 'follow', and 6 samples for 'forward' over 9 iterations. The experimental results show that, running Spiking-DSCNN on the Akida can achieve 94.74% accuracy. This high accuracy is due to the effective few-shot learning algorithm provided by the Akida's MetaTF framework.

Latency, Throughput, Power and Energy Consumption: The experimental results for latency and throughput for onchip learning with the Akida are presented in Fig. 8(d.1) and Fig. 8(d.2), respectively. Latency is stable around 1.5ms across 9 iterations of experiments (see ), which leads to more than 650KPS throughput (see ). Such low latency and high throughput come from the efficient few-shot learning that utilizes relatively small number of samples for new classes. The same reason also leads the on-chip learning to incur low power and energy consumption. Power consumption is about

Summary of comparison between our neuromorphic platform (Akida) with existing conventional AI solutions for object detection using YOLOv2; based on our results and data from state-of-the-art [28]–[30].

	Desktop CPU	Desktop GPU	Embedded CPU	Embedded GPU	FPGA			Our Akida
	Intel	Nvidia	ARM	Nvidia	ZedBoard	ZCU102	Virtex-7	Neuromorphic
	i7-6700HQ	GTX 960M	Cortex-A57	Jetson TX2	Zeuboaru	ZCU102	XC7V690t	Platform
Performance [FPS]	78.2	219.7	0.23	7.8	1.02	40.81	302.3	6
Power [W]	29.88	46.67	4	5.8	1.2	4.5	11.35	0.078
Efficiency [FPS/W]	2.62	4.71	0.06	1.34	0.85	9.06	26.63	76.92

41mW (see  $\bigcirc$ ), and energy consumption is about  $62\mu J$  (see  $\bigcirc$ ), as shown in Fig. 8(d.3) and Fig. 8(d.4), respectively.

#### E. Further Discussion

It is important to compare neuromorphic-based solutions against the state-of-the-art ANN-based solutions, which typically employ conventional hardware platforms, such as CPUs, GPUs, and specialized accelerators (e.g., FPGA or ASIC). To ensure a fair comparison, we select object recognition as the application and YOLOv2 as the network, while considering performance efficiency (FPS/W) as the comparison metric. Summary of the comparison is provided in Table I, and it clearly shows that our Akida-based neuromorphic solution achieves the highest performance efficiency. This is due to the sparse spike-driven computation that is fully exploited by neuromorphic processor, thus delivering highly power/energyefficient SNN processing. Moreover, our Akida-based neuromorphic solution also offers an on-chip learning capability, which gives it further advantages over the other solutions. This comparison highlights the immense potentials of neuromorphic computing for enabling efficient edge AI systems.

#### VI. CONCLUSION

We propose a novel design methodology to enable efficient SNN processing on commodity neuromorphic processors. It is evaluated using a real-world edge AI system implementation with the Akida processor. The experimental results demonstrate that, our methodology leads the system to achieve high performance and high energy efficiency across different applications. It achieves low latency of inference (i.e., less than 50ms for image classification, less than 200ms for real-time object detection in video streaming, and less than 1ms for keyword recognition) and low latency of on-chip learning (i.e., less than 2ms for keyword recognition), while consuming less than 250mW of power. In this manner, our design methodology potentially enables ultra-low power/energy design of edge AI systems for diverse application use-cases.

#### ACKNOWLEDGMENT

This work was partially supported by the NYUAD Center for Artificial Intelligence and Robotics (CAIR), funded by Tamkeen under the NYUAD Research Institute Award CG010.

#### REFERENCES

- D. McNulty et al., "A review of li-ion batteries for autonomous mobile robots: Perspectives and outlook for the future," *Journal of Power Sources (JPS)*, vol. 545, p. 231943, 2022.
- [2] G. Li et al., "Brain-inspired computing: A systematic survey and future trends," Proceedings of the IEEE, vol. 112, no. 6, pp. 544–584, 2024.
- [3] R. V. W. Putra and M. Shafique, "Fspinn: An optimization framework for memory-efficient and energy-efficient spiking neural networks," *IEEE Trans. on Computer-Aided Design of Integrated Circuits and Systems* (TCAD), vol. 39, no. 11, pp. 3601–3613, 2020.

- [4] A. Basu *et al.*, "Spiking neural network integrated circuits: A review of trends and future directions," in *CICC*, 2022, pp. 1–8.
- [5] R. V. W. Putra et al., "Embodied neuromorphic artificial intelligence for robotics: Perspectives, challenges, and research development stack," in ICARCV, 2024, pp. 612–619.
- [6] B. Vogginger et al., "Neuromorphic hardware for sustainable ai data centers," arXiv preprint arXiv:2402.02521, 2024.
- [7] BrainChip. Akida neural processor soc. [Online]. Available: https://brainchip.com/akida-neural-processor-soc/
- [8] SynSense. Dynap-cnn: The world's first fully scalable, eventdriven neuromorphic processor with up to 1m configurable spiking neurons and direct interface with external dvs. [Online]. Available: https://www.synsense.ai/products/dynap-cnn/
- [9] J. Dupeyroux et al., "Neuromorphic control for optic-flow-based landing of mavs using the loihi processor," in ICRA. IEEE, 2021, pp. 96–102.
- [10] S. Stroobants, J. Dupeyroux, and G. De Croon, "Design and implementation of a parsimonious neuromorphic pid for onboard altitude control for mavs using neuromorphic processors," in *ICONS*, 2022, pp. 1–7.
- [11] H. Patel *et al.*, "Bringing touch to the edge: A neuromorphic processing approach for event-based tactile systems," in *AICAS*, 2023, pp. 1–5.
- [12] S. Venkatachalam et al., "Realtime person identification via gait analysis using imu sensors on edge devices," in ICONS, 2024, pp. 371–375.
- [13] C. Kadway et al., "Low power & low latency cloud cover detection in small satellites using on-board neuromorphic processors," in *IJCNN*, 2023, pp. 1–8.
- [14] G. Lenz et al., "Low-power ship detection in satellite images using neuromorphic hardware," arXiv preprint arXiv:2406.11319, 2024.
- [15] D. Silva et al., "End-to-end edge neuromorphic object detection system," in AICAS, 2024, pp. 194–198.
- [16] R. V. W. Putra and M. Shafique, "Spikedyn: A framework for energy-efficient spiking neural networks with continual and unsupervised learning capabilities in dynamic environments," in DAC, 2021, p. 1057.
- [17] BrainChip. Metatf: The akida neuromorphic ml framework. [Online]. Available: https://doc.brainchipinc.com/index.html
- [18] R. V. W. Putra and M. Shafique, "Ipspikecon: Enabling low-precision spiking neural network processing for efficient unsupervised continual learning on autonomous agents," in *IJCNN*, 2022, pp. 1–8.
- [19] M. F. Minhas et al., "Continual learning with neuromorphic computing: Theories, methods, and applications," arXiv preprint:2410.09218, 2024.
- [20] R. V. W. Putra and M. Shafique, "Q-spinn: A framework for quantizing spiking neural networks," in *IJCNN*, 2021, pp. 1–8.
- [21] R. V. W. Putra, M. A. Hanif, and M. Shafique, "Respawn: Energy-efficient fault-tolerance for spiking neural networks considering unreliable memories," in *ICCAD*, 2021, pp. 1–9.
- [22] ——, "Softsnn: Low-cost fault tolerance for spiking neural network accelerators under soft errors," in DAC, 2022, pp. 151–156.
- [23] M. Demler, "Brainchip akida is a fast learner, spiking-neural-network processor identifies patterns in unlabeled data," *The Linley Group: Microprocessor Report*, vol. 28, 2019.
- [24] R. V. W. Putra, M. A. Hanif, and M. Shafique, "Sparkxd: A framework for resilient and energy-efficient spiking neural network inference using approximate dram," in *DAC*, 2021, pp. 379–384.
- [25] J. Deng et al., "Imagenet: A large-scale hierarchical image database," in CVPR, 2009, pp. 248–255.
- [26] M. Everingham et al., "The PASCAL Visual Object Classes Challenge 2007 (VOC2007) Results," http://www.pascalnetwork.org/challenges/VOC/voc2007/workshop/index.html.
- [27] P. Warden, "Speech commands: A dataset for limited-vocabulary speech recognition," arXiv preprint arXiv:1804.03209, 2018.
- [28] C. Liu, "Yolov2 acceleration using embedded gpu and fpgas: pros, cons, and a hybrid method," *Evolutionary Intelligence*, vol. 15, no. 4, 2022.
  [29] H. Nakahara *et al.*, "A lightweight yolov2: A binarized cnn with a
- [29] H. Nakahara *et al.*, "A lightweight yolov2: A binarized cnn with a parallel support vector regression for an fpga," in *FPGA*, 2018, p. 31.
- [30] S. Yan *et al.*, "An fpga-based mobilenet accelerator considering network structure characteristics," in *FPL*, 2021, pp. 17–23.