

An Ultra-low Power TinyML System for Real-time Visual Processing at Edge

Kunran Xu[†], Huawei Zhang[†], Yishi Li, Yuhao Zhang, Rui Lai, *Member, IEEE*, and Yi Liu

Abstract—Tiny machine learning (TinyML), executing AI workloads on resource and power strictly restricted systems, is an important and challenging topic. This brief firstly presents an extremely tiny backbone to construct high efficiency CNN models for various visual tasks. Then, a specially designed neural co-processor (NCP) is interconnected with MCU to build an ultra-low power TinyML system, which stores all features and weights on chip and completely removes both of latency and power consumption in off-chip memory access. Moreover, an application specific instruction-set is further presented for realizing agile development and rapid deployment. Extensive experiments demonstrate that the proposed TinyML system based on our tiny model, NCP and instruction set yields considerable accuracy and achieves a record ultra-low power of 160mW while implementing object detection and recognition at 30FPS. The demo video is available on <https://www.youtube.com/watch?v=mZPxtJ-9EY>.

Index Terms—Convolutional neural network, tiny machine learning, internet of things, application specific instruction-set

I. INTRODUCTION

Running machine learning inference on the resource and power limited environments, also known as Tiny Machine Learning (TinyML), has grown rapidly in recent years. It is promising to drastically expand the application domain of healthcare, surveillance, and IoT, *etc* [1], [2]. However, TinyML presents severe challenges due to large computational load and memory demand of AI models, especially in vision applications. Popular solutions using CPU+GPU architecture has shown high flexibility in MobileML applications [3], but it is no longer feasible in TinyML for the much stricter constraints on hardware resources and power consumption. A typical TinyML system based on microcontroller unit (MCU) usually has only < 512KB on-chip SRAM, <2MB Flash, <1GOP/s computing capability, and <1W power limitation [2], [4]. Meanwhile, it is difficult to use off-chip memory (*e.g.*, DRAM) in TinyML system for the very limited energy budget, showing a huge gap between the desired and available storage capacity for running visual AI models.

[†] Authors contributed equally to this work.

This work was supported in part by the National Key R&D Program of China under Grant 2018YF70202800, Natural Science Foundation of China (NSFC) under Grant 61674120. (Corresponding author: Rui Lai).

Kunran Xu, Huawei Zhang, Yishi Li, Yuhao Zhang and Rui Lai are with the School of Microelectronics, Xidian University, Xi'an 710071, and also with the Chongqing Innovation Research Institute of Integrated Circuits, Xidian University, Chongqing 400031, China. (e-mail: aazzttcc@gmail.com; myyzhww@gmail.com; yshlee1994@outlook.com; stuyuh@163.com; rlai@mail.xidian.edu.cn).

Yi Liu is with the School of Microelectronics, Xidian University, Xi'an 710071, and also with the Guangzhou Institute of Technology, Xidian University, Guangzhou 510555, China. (yliu@mail.xidian.edu.cn).

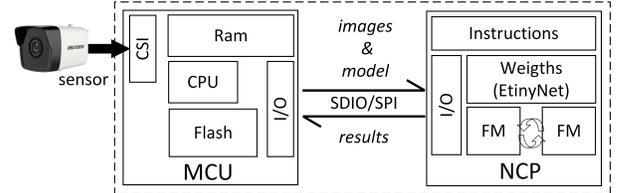


Fig. 1. The overview of the proposed TinyML system for visual processing.

Recently, the continuously emerging studies on TinyML achieve to deploy CNNs on MCUs by introducing memory-efficient inference engines [1], [4] and more compact CNN models [5], [6]. However, the existing TinyML systems still struggle to implement high-accuracy and real-time inference with ultra-low power consumption. Such as the state-of-the-art MCUNet [1] obtains 5FPS on STM32F746 but only achieves 49.9% top-1 accuracy on ImageNet. When the frame rate is increased to 10FPS, the accuracy of MCUNet further drops to 40.5%. What's more, running CNNs on MCUs is still not an extremely power-efficient solution due to the low efficiency of general purpose CPU in intensive convolution computing and massive weight data transmission. Considering this, we propose to greatly promote TinyML system by jointly designing more efficient CNN models and specific CNN co-processor. Specifically, we firstly design an extremely tiny CNN backbone EtinyNet aiming at TinyML applications, which has only 477KB model weights and maximum feature map size of 128KB and still yields remarkable 66.5% ImageNet Top-1 accuracy. Then, an ASIC-based neural co-processor (NCP) is specially designed for accelerating the inference. Since implementing CNN inference in a fully on-chip memory access manner, the proposed NCP achieves up to 180FPS throughput with 73.6mW ultra-low power consumption. On this basis, we propose a state-of-the-art TinyML system shown in Fig.2 for visual processing, which yields a record low power of 160mW in object detecting and recognizing at 30FPS.

In summary, we make the following contributions:

- 1) An extremely tiny CNN backbone named EtinyNet is specially designed for TinyML. It is far more efficient than existing lightweight CNN models.
- 2) An efficient neural co-processor (NCP) with specific designs for tiny CNNs is proposed. While running EtinyNet, NCP provides remarkable processing efficiency and convenient interface with extensive MCUs via SDIO/SPI.
- 3) Building upon the proposed EtinyNet and NCP, we promote the visual processing TinyML system to achieve a record ultra-low power and real-time processing efficiency, greatly advancing the TinyML community.

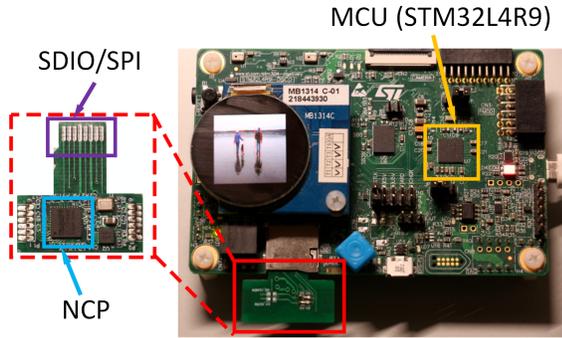


Fig. 2. The TinyML system for verification.

II. SOLUTION OF OUR TINYML SYSTEM

Fig.1 shows the overview of the proposed TinyML system. It integrates MCU with the specially designed energy-efficient NCP on a compact board to achieve superior efficiency in a collaborative work manner. To the best of our knowledge, we are the first to propose such a collaborative architecture in TinyML field, which successfully balances the efficiency and flexibility in the inference.

Initially, MCU sends the model weights and instructions to NCP who has sufficient on-chip SRAM to cache all these data. During inference, NCP computes the intensive CNN backbone efficiently while MCU only performs the light-load pre-processing (color normalization) and post-processing (fully-connected layer, non-maximum suppression, *etc*), which improves the overall energy efficiency to the greatest extent. Besides, the inference process of NCP only involves two kinds of data transfer, which are the input image and the output results. This working mode greatly reduces the off-chip data transfer power consumption and overall processing latency, and helps the system to achieve high energy efficiency in computing, which will be demonstrated in Section VI.

Considering real-time application, we interconnects NCP and MCU with SDIO/SPI interface. SDIO could provide up to 500Mbps bandwidth, which can transmit about 300FPS for 256×256 RGB image and 1200FPS for 128×128 one. As for SPI, it still reaches 100Mbps, or an equivalent throughput of 60FPS for 256×256 RGB image. These two buses are widely supported by MCUs available in the market, which makes NCP can be applied in a wide range of TinyML systems.

Fig.2 shows the prototype verification system only consisting of STM32L4R9 MCU and our proposed NCP. Thanks to the innovative model (EtinyNet), co-processor (NCP) and application specific instruction-set, the entire system yields both of efficiency and flexibility.

III. PARAMETER-EFFICIENT ETINYNET MODEL

Since NCP handles CNN workloads entirely on-chip for pursuing extreme efficiency, we focus on reducing the model size for satisfying the memory constrains of IoT devices in TinyML, which is totally different from MobileML targeting at the reduction of MAdds. By presenting Linear Depthwise Block (LB) and Dense Linear Depthwise Block (DLB), we derive an extremely tiny CNN backbone EtinyNet, shown in Fig.3.

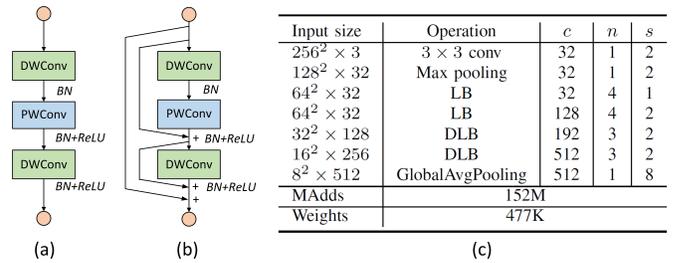


Fig. 3. The proposed building blocks that make up the EtinyNet. (a) is the Linear Depthwise Block (LB), (b) is the Dense Linear Depthwise Block (DLB), and (c) is the configuration of the backbone.

A. Design of Proposed Blocks

We present the linear depthwise convolution by removing the ReLU behind $DWConv$ of ϕ_{d1} under the observation that this non-linearity harms accuracy in the design of extremely parameter-efficient architectures, forming a specific case of sparse coding. Then, we introduce additional $DWConv$ of ϕ_{d2} behind $PWConv$ of ϕ_p to build a novel linear depthwise block (LB) by utilizing $DWConv$'s parameter efficiency [7]. The LB is defined as

$$\mathbf{O} = \sigma(\phi_{d2}(\sigma(\phi_p(\phi_{d1}(\mathbf{I})))))) \quad (1)$$

As shown in Fig 3(a), the structure of proposed LB can be represented as $DWConv-PWConv-DWConv$, which is apparently different from the commonly used bottleneck block of $PWConv-DWConv-PWConv$ in mobile models, explained by the fact that increasing the proportion of $DWConv$ is beneficial to the accuracy of tiny models.

Additionally, we introduce the dense connection into LB for increasing its equivalent width, which is important and necessary for a higher accuracy [8], as well as the very limited size of features and weights. We refer the resulting block to Dense Linear Depthwise Block (DLB) depicted in Fig 3(b). Note that we take the ϕ_{d1} and ϕ_p as a whole due to the removal of ReLU, and add the shortcut connection at the ends of these two layers.

B. Architecture of EtinyNet Backbone

By stacking LBs and DLBs, we configure the EtinyNet backbone as indicated in Fig 3(c), where n , c and s represent block repeated times, the number of output channels, and the first layer's stride in each block (other layers' stride equaling one) respectively. Since dense connection consumes more memory space, we only utilize DLB at high level stages with much smaller feature maps. It's encouraging that EtinyNet backbone has only 477KB parameters and still achieves 66.5% ImageNet Top-1 accuracy. The extreme compactness of EtinyNet makes it possible to design small footprint NCP that could run without off-chip DRAM.

IV. APPLICATION SPECIFIC INSTRUCTION-SET FOR NCP

For easily deploying tiny CNN models on NCP, we define an application specific instruction-set. As shown in Table I, the set contains 13 instructions, belonging to neural operation

type and control type respectively. It includes basic operations for tiny CNN models, and each instruction consists of 128 bits: 5 bits for operation code, and the rest for attributes of operations and operands. With each neural type instruction encoding an entire layer, the proposed instruction-set has a relatively coarser granularity, which simplifies the control complexity of hardware. Moreover, the basic operations included in the instruction-set provide sufficient ability to execute commonly-used CNN architectures (e.g., MobileNetV2 [9], MobileNeXt [10], *etc.*).

TABLE I
INSTRUCTION SET FOR PROPOSED NCP

Instruction format	Description	Type
bn	batch normalization	N
relu	non-linear activation operation	N
conv	1x1 and 3x3 convolution & bn, relu	N
dwconv	3x3 depthwise conv & bn, relu	N
add	elementwise addition	N
move	move tensor to target address	N
dsam	down-sampling by factor of 2	N
usam	up-sampling by factor of 2	N
maxp	max pooling by factor of 2	N
gap	global average pooling	N
jump	set program counter (PC) to target	C
sup	suspend processor	C
end	suspend processor and reset PC	C

V. DESIGN OF NEURAL CO-PROCESSOR

As shown in Fig.4, the proposed NCP consists of five main components: Neural Operation Unit (NOU), Tensor Memory (TM), Instruction Memory (IM), I/O and System Controller (SC). When NCP works, SC decodes one instruction fetched from IM and informs the NOU to start computing with decoded signals. The computing process takes multiple cycles, during which NOU reads operands from TM and writes results back automatically. Once completing the writing back process, SC continues to process the next instruction until an **end** or **suspend** instruction is encountered. When NOU is idle, TM is accessed through I/O. We will fully describe each component in the following parts.

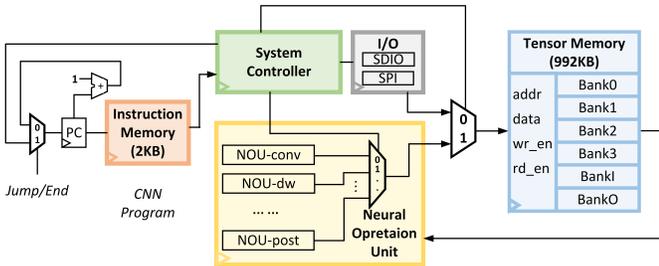


Fig. 4. The overall block diagram of the proposed NCP.

A. Neural Operation Unit

CNN workloads mainly come from operations of int8 **conv**, **dwconv** and float32 **bn**. To achieve a high energy efficiency, we respectively design special hardware units, termed NOU-conv, NOU-dw and NOU-post, focusing on optimizing the

implementation of each operation. Furthermore, we deal with the design details in the following three aspects.

1) Different from other designs [11], [12] with fine grained instructions, we implement NOU-conv with a hardwired matrix multiply-accumulate (MAC) [13] array, which helps to improve efficiency with a simpler control logic. The MAC array is designed to perform matrix outer product with parallelism in spatial and output channel dimension for handling the most computational costly 3×3 Conv and *PWConv* with *im2col* operation. In this way, the number of effective multiplications in each cycle is fixed to $T_{oc} \times T_{hw}$. Note that the number of channels varies across different convolution layers, which may lead to inefficient computation for other ways of implementation (e.g., dot product). Conversely, our implementation manner can avoid the above-mentioned problem and improve the overall efficiency in running *PWConv* of entire network. Moreover, the addition is realized by simple accumulation process instead of commonly-used adder tree with extra hardware overhead.

2) As for the implementation of *DWConv*, the above designed MAC array proves its efficiency only in diagonal units. Given this, we turn to the classical convolution processing pipeline [14], where nine multipliers and eight adders are arranged to compute *DWConv* in each channel. The independence between channels allows us to extend pipelines easily, implementing a parallelism of T_{oc} to build NOU-dw. Since the feature length in spatial dimension is usually much larger than the pipeline depth, the *DWConv* can be performed in a fully pipelined manner, which yields NOU-dw an ultra high efficiency up to nearly 100% .

3) In NOU-post unit, modules of int2float, float32 multiply-add, float2int and ReLU are designed and interconnected to perform post-operations of float32 BN, ReLU and elementwise addition. To reduce memory access as much as possible, multiplexers are further utilized to select data from the output of NOU-conv, NOU-dw or TM, and connect modules as needed, allowing flexible fusion of post-operations with the previous convolution layer. By implementing T_{oc} pipelines to match the throughput of convolution, we effectively maximize the efficiency of fusion operations.

B. Tensor Memory and Tensor Layout

1) TM is a single-port SRAM consisting of 6 banks, whose width is $T_{tm} \times 8$ bits, as shown in Fig 4. Thanks to the compactness of EtinyNet, NCP only requires totally 992KB on-chip SRAM. The Bank1 (192KB) is responsible for caching 256×256 input RGB images. The 128KB sized Bank0 and Bank1 are arranged for caching feature maps, while Bank2 and Bank3 with a larger size of 256KB are used for storing weights. The Bank0 (32KB) is used to store final results, such as feature vectors and bonding boxes, *etc.* TM's small capacity and simple structure yield our NCP a small footprint.

2) The highly efficient NOU brings 2 types of tensor layouts, named pixel-major layout and interleaved layout respectively shown in Fig.5. For the former, all pixels of the first channel are sequentially mapped to TM in a row-major order. Then, the next channel's counterpart is arranged in the same pattern

until the last channel’s pixels of a tensor are stored. For the latter, the whole tensor is divided into $N_c // T_{tm}$ tiles and are placed in TM sequentially, while each tile is arranged in a channel-major order. Different layouts are required for NOUs to achieve the maximum efficiency. For example, the input of NOU-conv prefers pixel-major layout because spatially continuous T_{hw} pixels of a channel need to be multiplied and added at a time by MAC array, while the reverse is the case for NOU-dw.

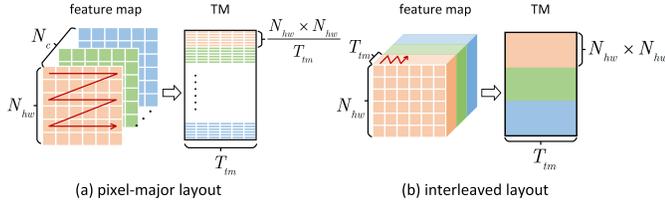


Fig. 5. Illustration of different tensor layouts. (a) Pixel-major layout. (b) Interleaved layout.

3) Running the proposed LB, DLB, and other blocks with NOU, the layout between adjacent *DWConv* and *PWConv* is constantly varying, which seriously decreases the computing efficiency because of the discontinuous memory access. It takes NOU-conv T_{oc} times to read the output of NOU-dw stored in an interleaved layout for performing a single matrix outer product operation. Hence, an efficient layout conversion circuit is designed to tackle this problem. As shown in Fig.6, the circuit is composed of two $T_{oc} \times T_{hw}$ register arrays A and B, working in a ping-pong mechanism. At the beginning, array A receives T_{oc} inputs at a time, after T_{hw} cycles, A will be filled and start to output T_{hw} results at a time in the transposed dimension. Since reading A empty requires T_{oc} cycles, the new coming data to be converted will be sent to array B in order to maintain the pipeline. When B is full and A completes the readout, the role of them are exchanged. This strategy obviously boosts the efficiency of valid memory access for computing.

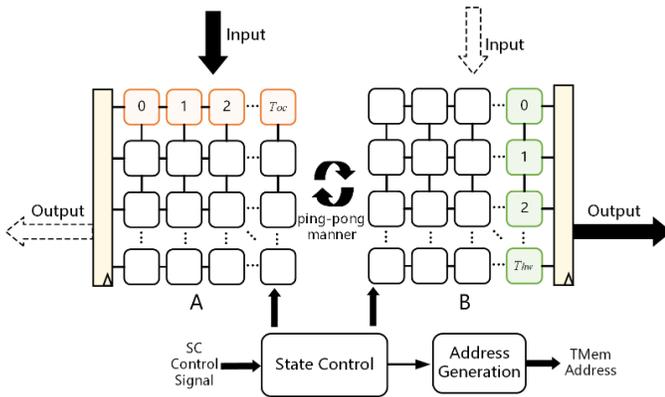


Fig. 6. The proposed efficient layout conversion circuit.

C. Characteristics

We implement our NCP using TSMC 65nm low-power technology. While $T_{tm} = 32$, $T_{oc} = 16$ and $T_{hw} = 32$, NCP contains 512 of 8-bit MACs in NOU-conv, 144 of 8-bit multipliers and 16 of adder trees in NOU-dw, and 16 of float32 MACs in NOU-post. When running at the maximum frequency of 250MHz, NOU-conv and NOU-post are active every cycle, achieving a peak performance of 264 GOP/s.

VI. EXPERIMENTAL RESULTS

A. EtinyNet Evaluation

Table II lists the ImageNet-1000 classification results of well-known lightweight CNNs, including MobileNetV2 [9], MobileNeXt [10], ShuffleNetV2 [15], and MCUNet series [16]. We pay more attention to the backbone because the fully-connected layer is generally not involved in most of visual models. Among these competitive models, MCUNet gets the highest accuracy at the cost of model size up to 2048K. Compared with tiny models in similar size, our EtinyNet reaches 66.5% top-1 and 86.8% top-5 accuracy, outperforming the most competitive MCUNetV2-M4 by significant 1.6% top-1 accuracy. Moreover, EtinyNet-0.75, the width of each layer is shrunk by 0.75, outperforms MCUNet-320kB by significant 2.6% top-1 accuracy with 60K fewer parameters. Obviously, EtinyNet yields much higher accuracy at the same level of storage consumption, and is more suitable for TinyML systems.

TABLE II
COMPARISON OF STATE-OF-THE-ART TINY MODELS OVER ACCURACY ON IMAGE NET. "B" DENOTES BACKBONE. "-" DENOTES NOT REPORTED.

Model	#Params. (K)	Top-1 Acc.	Top-5 Acc.
MobileNeXt-0.35	812(B) / 1836	64.7	85.7
MobileNetV2-0.35	740(B) / 1764	60.3	82.9
ShuffleNetV2-0.5	566(B) / 1590	61.1	82.6
MCUNet	-(B) / 2048	70.7	-
MCUNet-320kB	-(B) / 740	61.8	84.2
MCUNetV2-M4	-(B) / 1034	64.9	86.2
EtinyNet	477(B) / 989	66.5	86.8
EtinyNet-0.75	296(B) / 680	64.4	85.2
EtinyNet-0.5	126(B) / 446	59.3	81.2

B. NCP Evaluation

As shown in Table III, running general CNN models usually needs DRAMs to store their enormous weights and features [11], [18], resulting in considerable power consumption and processing latency. As for no DRAM access methods, YodaNN [17] yields the highest peak performance and energy efficiency, but it is a dedicated accelerator only for binarized networks with very limited accuracy. Except that, Vega [12] gets the lowest power and the maximum latency, which leads to the lowest peak performance. To comprehensively assess the throughput, energy consumption and speed of various neural processors in TinyML application, we prefer to use the metric of processing efficiency, which is the number of frames processed per unit time and per unit power consumption. Our proposed NCP reaches an extremely high processing

TABLE III
COMPARISON WITH STATE-OF-THE-ART NEURAL PROCESSORS. “-”
DENOTES NOT REPORTED.

Component	NullHop	ConvAix	YodaNN	Vega	NCP
Technology	28nm	28nm	65nm	22nm	65nm
Area (mm^2)	6.3	3.53	1.9	12	10.88
DRAM Used	yes	yes	none	none	none
FC Support	none	yes	none	yes	none
CNN model	VGG16	MbV1	VGG19	RVGGA0	EtinyNet
ImageNet Acc.	68.3%	70.6%	-	72.4%	66.5%
Latency	72.9ms	14.2ms	75.2ms	118ms	5.5ms
Typ. Power (mW)	155.0	313.1	153	37.3	73.6
Peak Perf. (GOP/s)	128	262.6	1500	32.2	264
Energy Eff. (GOP/s/W)	2714.8	256.3	8500	631.4	751.0
Processing Eff. (Frames/s/mJ)	1.21	15.18	2.95	1.93	449.1

efficiency up to 449.1 Frames/s/mJ, at least $29\times$ higher than other solutions, suggesting the unique superiority of NCP in this particular field. As for the reason, the specially designed NOU, tensor layout and coarse-grained instruction-set jointly decrease the delay and the power of inference.

C. TinyML System Verification

We compare our proposed system with existing prominent MCU-based TinyML systems. As shown in Table IV, CMSIS-NN obtains 59.5% ImageNet accuracy at 2FPS, promoted by MCUNet to 5FPS at the expense of accuracy dropping to 49.9%. In comparison, our solution reaches up to 66.5% accuracy and 30FPS, achieving the goal of real-time visual processing in TinyML. Furthermore, since existing methods burden MCUs with entire CNN models, high-performance MCUs (STM32H743/STM32F746) running at the upper-limit frequency (480MHz/216MHz) are necessary. Although flexible, general-purpose MCU is of low energy efficiency in computing massive tensors, which results in considerable power consumption up to about 600mW. In contrast, the proposed solution allows us to perform the same flexible task only with a low-end MCU (STM32L4R9, 120MHz) and proposed NCP, which boosts the energy efficiency of the entire system and achieves an ultra-low power of 160mW.

TABLE IV
COMPARISON WITH MCU-BASED DESIGNS ON IMAGE CLASSIFICATION (CLS) AND OBJECT DETECTION (DET). * DENOTES REPRODUCED RESULTS.

	Method	Hardware	Acc/mAP	FPS	Power
Cls	CMSIS-NN	H743	59.5%	2	*675 mW
	MCUNet	F746	49.9%	5	*525 mW
	Ours	L4R9+NCP	66.5%	30	160 mW
Det	CMSIS-NN	H743	31.6%	10	*640 mW
	MCUNet	H743	51.4%	3	*650 mW
	Ours	L4R9+NCP	56.4%	30	160 mW

In addition, we benchmark the object detection performance on Pascal VOC dataset. The results indicate that our system also greatly improves its performance, which makes AIoT more promising in extensive applications.

VII. CONCLUSION

In this brief, we propose an ultra-low power TinyML system for real-time visual processing by designing 1) an extremely tiny CNN backbone EtinyNet, 2) an ASIC-based neural co-processor and 3) an application specific instruction-set. Our study greatly advances the TinyML community and promises to drastically expand the application scope of AIoT.

REFERENCES

- [1] J. Lin, W.-M. Chen, Y. Lin, C. Gan, S. Han *et al.*, “McuNet: Tiny deep learning on iot devices,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 11 711–11 722, 2020.
- [2] M. Shafique, T. Theodorides, V. J. Reddy, and B. Murmann, “Tinyml: Current progress, research challenges, and future roadmap,” in *2021 58th ACM/IEEE Design Automation Conference (DAC)*. IEEE, 2021, pp. 1303–1306.
- [3] S. Mittal, “A survey on optimized implementation of deep learning models on the nvidia jetson platform,” *Journal of Systems Architecture*, vol. 97, pp. 428–442, 2019.
- [4] L. Lai, N. Suda, and V. Chandra, “Cmsis-nn: Efficient neural network kernels for arm cortex-m cpus,” *arXiv preprint arXiv:1801.06601*, 2018.
- [5] C. Banbury, C. Zhou, I. Fedorov, R. Matas, U. Thakker, D. Gope, V. Janapa Reddi, M. Mattina, and P. Whatmough, “Micronets: Neural network architectures for deploying tinyml applications on commodity microcontrollers,” *Proceedings of Machine Learning and Systems*, vol. 3, pp. 517–532, 2021.
- [6] R. T. N. Chappa and M. El-Sharkawy, “Deployment of se-squeezenext on nxp bluebox 2.0 and nxp i. mx rt1060 mcu,” in *2020 IEEE Midwest Industry Conference (MIC)*, vol. 1. IEEE, 2020, pp. 1–4.
- [7] K. Xu, Y. Li, H. Zhang, R. Lai, and L. Gu, “EtinyNet: Extremely tiny network for tinyml,” in *AAAI Conference on Artificial Intelligence*. AAAI, 2022, pp. 4628–4636.
- [8] S. Zagoruyko and N. Komodakis, “Wide residual networks,” *arXiv preprint arXiv:1605.07146*, 2016.
- [9] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, “Mobilenetv2: Inverted residuals and linear bottlenecks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4510–4520.
- [10] D. Zhou, Q. Hou, Y. Chen, J. Feng, and S. Yan, “Rethinking bottleneck structure for efficient mobile network design,” in *European Conference on Computer Vision*. Springer, 2020, pp. 680–697.
- [11] A. Bytyn, R. Leupers, and G. Ascheid, “Convaix: An application-specific instruction-set processor for the efficient acceleration of cnns,” *IEEE Open Journal of Circuits and Systems*, vol. 2, pp. 3–15, 2020.
- [12] D. Rossi, F. Conti, M. Eggiman, A. Di Mauro, G. Tagliavini, S. Mach, M. Guermandi, A. Pullini, I. Loi, J. Chen *et al.*, “Vega: A ten-core soc for iot endnodes with dnn acceleration and cognitive wake-up from mram-based state-retentive sleep mode,” *IEEE Journal of Solid-State Circuits*, vol. 57, no. 1, pp. 127–139, 2021.
- [13] M. E. Nojehdeh, S. Parvin, and M. Altun, “Efficient hardware implementation of convolution layers using multiply-accumulate blocks,” in *2021 IEEE Computer Society Annual Symposium on VLSI (ISVLSI)*. IEEE, 2021, pp. 402–405.
- [14] T. Chen, Z. Du, N. Sun, J. Wang, C. Wu, Y. Chen, and O. Temam, “Diannao: A small-footprint high-throughput accelerator for ubiquitous machine-learning,” *ACM SIGARCH Computer Architecture News*, vol. 42, no. 1, pp. 269–284, 2014.
- [15] N. Ma, X. Zhang, H.-T. Zheng, and J. Sun, “Shufflenet v2: Practical guidelines for efficient cnn architecture design,” in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 116–131.
- [16] J. Lin, W. Chen, H. Cai, C. Gan, and S. Han, “McuNetv2: Memory-efficient patch-based inference for tiny deep learning. arxiv 2021,” *arXiv preprint arXiv:2110.15352*.
- [17] R. Andri, L. Cavigelli, D. Rossi, and L. Benini, “Yodann: An architecture for ultralow power binary-weight cnn acceleration,” *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 37, no. 1, pp. 48–60, 2017.
- [18] A. Aimar, H. Mostafa, E. Calabrese, A. Rios-Navarro, R. Tapiador-Morales, I.-A. Lungu, M. B. Milde, F. Corradi, A. Linares-Barranco, S.-C. Liu *et al.*, “Nullhop: A flexible convolutional neural network accelerator based on sparse representations of feature maps,” *IEEE transactions on neural networks and learning systems*, vol. 30, no. 3, pp. 644–656, 2018.