# AnalogNAS: A Neural Network Design Framework for Accurate Inference with Analog In-Memory Computing

Hadjer Benmeziane*⬤, Corey Lammie†⬤, Irem Boybat†⬤, Malte Rasch‡⬤, Manuel Le Gallo†⬤, Hsinyu Tsai§⬤,
Ramachandran Muralidhar‡⬤, Smail Niar*⬤, Ouarnoughi Hamza*⬤, Vijay Narayanan‡,
Abu Sebastian†⬤ and Kaoutar El Maghraoui‡⬤,

*Univ. Polytechnique Hauts-de-France, CNRS, UMR 8201 - LAMIH, F-59313 Valenciennes, France
†IBM Research Europe, 8803 Rüschlikon, Switzerland
‡IBM T. J. Watson Research Center, Yorktown Heights, NY 10598, USA
§ IBM Research Almaden, 650 Harry Road, San Jose, CA USA

*Abstract*—The advancement of Deep Learning (DL) is driven by efficient Deep Neural Network (DNN) design and new hardware accelerators. Current DNN design is primarily tailored for general-purpose use and deployment on commercially viable platforms. Inference at the edge requires low latency, compact and power-efficient models, and must be cost-effective. Digital processors based on typical von Neumann architectures are not conducive to edge AI given the large amounts of required data movement in and out of memory. Conversely, analog/mixed-signal in-memory computing hardware accelerators can easily transcend the memory wall of von Neuman architectures when accelerating inference workloads. They offer increased area- and power efficiency, which are paramount in edge resource-constrained environments. In this paper, we propose *AnalogNAS*, a framework for automated DNN design targeting deployment on analog In-Memory Computing (IMC) inference accelerators. We conduct extensive hardware simulations to demonstrate the performance of AnalogNAS on State-Of-The-Art (SOTA) models in terms of accuracy and deployment efficiency on various Tiny Machine Learning (TinyML) tasks. We also present experimental results that show AnalogNAS models achieving higher accuracy than SOTA models when implemented on a 64-core IMC chip based on Phase Change Memory (PCM). The AnalogNAS search code is released[1].

*Index Terms*—Analog AI, Neural Architecture Search, Optimization, Edge AI, In-memory Computing

## I. INTRODUCTION

WITH the growing demands of real-time DL workloads, today's conventional cloud-based AI deployment approaches do not meet the ever-increasing bandwidth, real-time, and low-latency requirements. Edge computing brings storage and local computations closer to the data sources produced by the sheer amount of Internet of Things (IoT) objects, without overloading network and cloud resources. As DNNs are becoming more memory and compute intensive, edge AI deployments on resource-constrained devices pose significant challenges. These challenges have driven the need for specialized hardware accelerators for on-device Machine Learning (ML) and a plethora of tools and solutions targeting the development and deployment of power-efficient edge AI solutions. One such promising technology for edge hardware accelerators is analog-based IMC, which is herein referred to as *analog IMC*.

Analog IMC [1] can provide radical improvements in performance and power efficiency, by leveraging the physical properties of memory devices to perform computation and storage at the same physical location. Many types of memory devices, including Flash memory, PCM, and Resistive Random Access Memory (RRAM), can be used for IMC [2]. Most notably, analog IMC can be used to perform Matrix-Vector Mutliplication (MVM) operations in $O(1)$ time complexity [3], which is the most dominant operation used for DNN acceleration. In this novel approach, the weights of linear, convolutional, and recurrent DNN layers are mapped to crossbar arrays (tiles) of Non-Volatile Memory (NVM) elements. By exploiting basic Kirchhoff's circuit laws, MVMs can be performed by encoding inputs as Word-Line (WL) voltages and weights as device conductances. For most computations, this removes the need to pass data back and forth between Central Processing Units (CPUs) and memory. This back and forth data movement is inherent in traditional digital computing architectures, and is often referred to as the *von Neumann bottleneck*. Because there is greatly reduced movement of data, tasks can be performed in a fraction of the time, and with much less energy.

NVM crossbar arrays and analog circuits, however, have inherent non-idealities, such as noise, temporal conductance drift, and non-linear errors, which can lead to imprecision and noisy computation [4]. These effects need to be properly quantified and mitigated to ensure the high accuracy of DNN models. In addition to the hardware constraints that are prevalent in edge devices, there is the added complexity of designing DNN
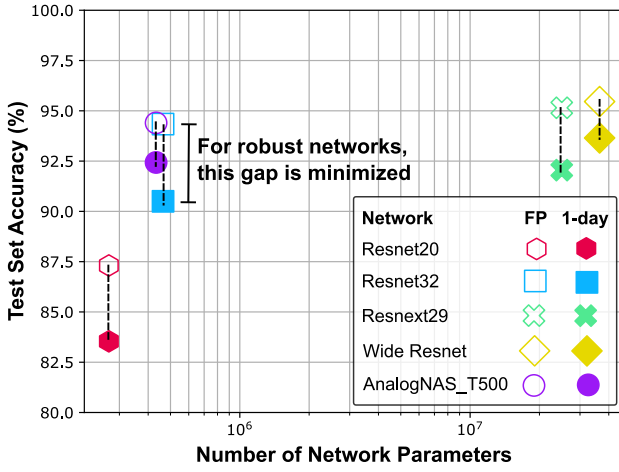
[1]https://github.com/IBM/analog-nas

Fig. 1. The effect of PCM conductance drift after one day on standard CNN architectures and one architecture (`AnalogNAS_T500`) obtained using HW-NAS, evaluated using CIFAR-10. *FP* refers to the original network accuracy, and *1-day* to the simulated analog network accuracy after 1-day device drift.

architectures which are optimized for the edge on a variety of hardware platforms. This requires hardware-software co-design approaches to tackle this complexity, as manually-designed architectures are often tailored for specific hardware platforms. For instance, MobileNet [5] uses a depth-wise separable convolution that enhances CPU performance but is inefficient for Graphics Processing Unit (GPU) parallelization [6]. These are bespoke solutions that are often hard to implement and generalize to other platforms.

HW-NAS [7] is a promising approach that seeks to automatically identify efficient DNN architectures for a target hardware platform. In contrast to traditional Neural Architecture Search (NAS) approaches that focus on searching for the most accurate architectures, HW-NAS searches for highly accurate models while optimizing hardware-related metrics. Existing HW-NAS strategies cannot be readily used with analog IMC processors without significant modification for three reasons: (i) their search space contains operations and blocks that are not suitable for analog IMC, (ii) lack of a benchmark of hardware-aware trained architectures, and (iii) their search strategy does not include noise injection and temporal drift on weights.

To address these challenges, we propose *AnalogNAS*, a novel HW-NAS strategy to design dedicated DNN architectures for efficient deployment on edge-based analog IMC inference accelerators. This approach considers the inherent characteristics of analog IMC hardware in the search space and search strategy. Fig. 1 depicts the necessity of our approach. As can be seen, when traditional DNN architectures are deployed on analog IMC hardware, non-idealities, such as conductance drift, drastically reduce network performance. Networks designed by *AnalogNAS* are extremely robust to these non-idealities and have much fewer parameters compared to equivalently-robust traditional networks. Consequently, they have reduced resource utilization.

Our specific contributions can be summarized as follows:

- We design and construct a search space for analog IMC, which contains ResNet-like architectures, including ResNext [8] and Wide-ResNet [9], with blocks of varying widths and depths;
- We train a collection of networks using Hardware-Aware (HWA) training for image classification, Visual Wake Words (VWW), and Keyword Spotting (KWS) tasks. Using these networks, we build a surrogate model to rank the architectures during the search and predict robustness to conductance drift;
- We propose a global search strategy that uses evolutionary search to explore the search space and efficiently finds the right architecture under different constraints, including the number of network parameters and analog tiles;
- We conduct comprehensive experiments to empirically demonstrate that AnalogNAS can be efficiently utilized to carry out architecture search for various edge tiny applications, and investigate what attributes of networks make them ideal for implementation using analog AI;
- We validate a subset of networks on hardware using a 64-core IMC chip based on PCM.

The rest of the paper is structured as follows. In Section II, we present related work. In Section III, relevant notations and terminology are introduced. In Section IV, the search space and surrogate model are presented. In Section V, the search strategy is presented. In Section VI, the methodology for all experiments is discussed. The simulation results are presented in Section VI-B, along with the hardware validation and performance estimation in Section VII. The results are discussed in Section VIII. Section IX concludes the paper.

## II. RELATED WORK

### A. NAS for TinyML

HW-NAS has been successfully applied to a variety of edge hardware platforms [7], [10] used to deploy networks for TinyMLPerf tasks [11] such as image classification, VWW, KWS, and anomaly detection. MicroNets [12] leverages NAS for DL model deployment on micro-controllers and other embedded systems. It utilizes a differentiable search space [13] to find efficient architectures for different TinyMLPerf tasks. For each task, the search space is an extension of current SOTA architectures. $\mu$-nas [14] includes memory peak usage and a number of other parameters as constraints. Its search strategy combines aging evolution and Bayesian optimization to estimate the objectives and explore a granular search space efficiently. It constructs its search space from a standard CNN and modifies the operators' hyper-parameters and a number of layers.

### B. NAS for Mixed-Signal IMC Accelerators

Many works [15]–[18] target IMC accelerators using HW-NAS. FLASH [15] uses a small search space inspired by DenseNet [19] and searches for the number of skip connections that efficiently satisfy the trade-off between accuracy, latency, energy consumption, and chip area. Its surrogate model uses linear regression and the number of

skip connections to predict model accuracy. NAS4RRAM [17] uses HW-NAS to find an efficient DNN for a specific RRAM-based accelerator. It uses an evolutionary algorithm, trains each sampled architecture without HWA training, and evaluates each network on a specific hardware instance. NACIM [16] uses co-exploration strategies to find the most efficient architecture and the associated hardware platform. For each sampled architecture, networks are trained considering noise variations. This approach is limited by using a small search space due to the high time complexity of training. UAE [18] uses a Monte-Carlo simulation-based experimental flow to measure the device uncertainty induced to a handful of DNNs. Similar to NACIM [16], evaluation is performed using HWA training with noise injection. AnalogNet [20] extends the work of Micronet by converting their final models to analog-friendly models, replacing depthwise convolutions with standard convolutions and tuning hyperparameters.

Compared to the above-mentioned SOTA HW-NAS strategies, our AnalogNAS is better tailored to analog IMC hardware for two reasons: (i) Our search space is much larger and more representative, featuring resnet-like connections. This enables us to answer the key question of what architectural characteristics are suitable for analog IMC which cannot be addressed with small search spaces. (ii) We consider the inherent characteristics of analog IMC hardware directly in the objectives and constraints of our search strategy in addition to noise injection during the HWA training as used by existing approaches.

## III. PRELIMINARIES

### A. Analog IMC Accelerator Mechanisms

Analog IMC accelerators are capable of performing MVM operations $\mathbf{Y}^T = \mathbf{X}^T \mathbf{W}$ using the laws of physics, where $\mathbf{W}$ is an $M \times N$ matrix, $\mathbf{X}$ is a $M \times 1$ vector, and $\mathbf{Y}$ is a $N \times 1$ vector. When arranged in a crossbar configuration, $M \times N$, NVM devices can be used to compute MVM operations. This is done by encoding elements of $\mathbf{X}$ as WL voltages, denoted using $\mathbf{V}$, and elements of $\mathbf{W}$ as conductances of the unit cells, denoted using $\mathbf{G}$. Negative conductance states cannot be directly encoded/represented using NVM devices. Consequently, differential weight mapping schemes are commonly employed, where either positive weights, i.e., $\mathbf{W}^+ = \max(\mathbf{W}, 0)$, and negative weights, i.e., $\mathbf{W}^- = -\min(\mathbf{W}, 0)$, are encoded within unit cells, using alternate columns, or on different tiles [3]. The analog computation, i.e., $\mathbf{I} = \mathbf{VG}$ is performed, where the current flow to the end of the $N$-th column is $I_N = \sum_{i=0}^{M} G_{i,N} V_i$. Typically, Digital-to-Analog Converters (DACs) are required to encode WL voltages and Analog-to-Digital Converters (ADCs) are required to read the output currents of each column. The employed analog IMC tile, its weight mapping scheme, and computation mechanism are depicted in Fig. 2.

### B. Temporal Drift of Non-Volatile Memory Devices

Many types of NVM devices, most prominently, PCM, exhibit temporal evolution of the conductance values referred to as the conductance drift. This poses challenges for maintaining synaptic weights reliably [2]. Conductance drift is most commonly modelled using Eq. (1), as follows:

$$G(t) = G(t_0)(t/t_0)^{-\nu}, \tag{1}$$

where $G(t_0)$ is the conductance at time $t_0$ and $\nu$ is the drift exponent. In practice, conductance drift is highly stochastic because $\nu$ depends on the programmed conductance state and varies across devices. Consequently, when reporting the network accuracy at a given time instance (after device programming), it is computed across multiple experiment instances (trials) to properly capture the amount of accuracy variations.

### C. HWA-training and analog hardware accuracy evaluation simulation

To simulate training and inference on analog IMC accelerators, the IBM Analog Hardware Acceleration Kit (AIHWKIT) [21] is used. The AIHWKIT is an open-source Python toolkit for exploring and using the capabilities of in-memory computing devices in the context of artificial intelligence and has been used for HWA training of standard DNNs with hardware-calibrated device noise and drift models [22].

### D. Hardware-aware Neural Architecture Search (HW-NAS)

HW-NAS refers to the task of automatically finding the most efficient DNN for a specific dataset and target hardware platform. HW-NAS approaches often employ black-box optimization methods such as evolutionary algorithms [23], reinforcement learning [24], [25], and Bayesian optimization [26], [27]. The optimization problem is either cast as a constrained or multi-objective optimization [7]. In AnalogNAS, we chose constrained optimization over multi-objective optimization for several reasons. First, constrained optimization is more computationally efficient than multi-objective optimization, which is important in the context of HW-NAS, to allow searching a large search space in a practical time frames. Multi-objective optimization is computationally expensive and can result in a
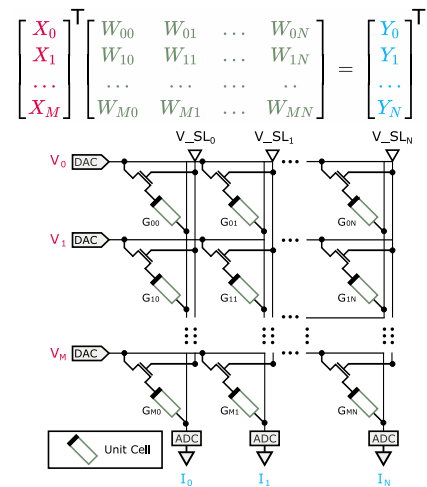


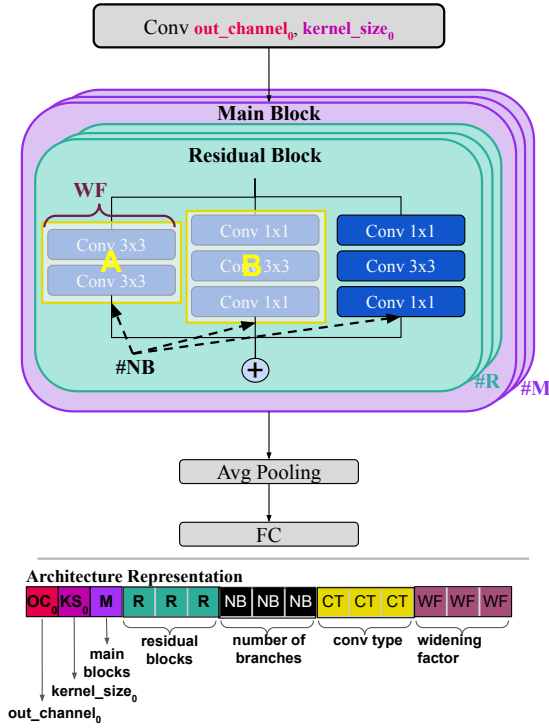Fig. 2. Employed analog IMC tile and weight mapping scheme.

Fig. 3. Resnet-like macro architecture.

a series of $M$ distinct main blocks. Each main block contains $R$ residual blocks. The residual blocks use skip connections with or without downsampling. Downsampling is performed using 1x1 convolution layers when required, i.e., when the input size does not match the output size. The residual block can have $B$ branches. Each branch uses a convolution block. We used different types of convolution blocks to allow the search space to contain all standard architectures such as Resnets [30], ResNext [8], and Wide Resnets [9]. The standard convolution blocks used in Resnets, commonly referred to as *BottleNeckBlock* and *BasicBlock*, are denoted as A and B respectively. We include variants of A and B in which we inverse the order of the ReLU and Batch normalization operations. The resulting blocks are denoted as C and D. Table I summarizes the searchable hyper-parameters and their respective ranges. The widening factor scales the width of the residual block. We sample architectures with different depths by changing the number of main and residual blocks. The total size of the search space is approximately 73B architectures. The larger architecture would contain 240 convolutions and start from an output channel of 128 multiplying that by 4 for every 16 blocks.

### B. Analog-accuracy Surrogate Model

*1) Evaluation Criteria:* To efficiently explore the search space, a search strategy requires evaluating the objectives of each sampled architecture. Training the sampled architectures is very time-consuming; especially when HWA retraining is performed, as noise injection and I/O quantization modeling greatly increases the computational complexity. Consequently, we build a surrogate model capable of estimating the objectives of each sampled architecture in IMC devices. To find architectures that maximize accuracy, stability, and resilience against IMC noise and drift characteristics, we have identified the following three objectives.

*a) The 1-day accuracy:* is the primary objective that most NAS algorithms aim to maximize. It measures the performance of an architecture on a given dataset. When weights are encoded using IMC devices, the accuracy of the architecture can drop over time due to conductance drift. Therefore, we have selected the 1-day accuracy as a metric to measure the architecture's performance.

*b) The Accuracy Variation over One Month (AVM):* is the difference between the 1-month and 1-sec accuracy. This objective is essential to measure the robustness over a fixed time duration. A 30-day period allows for a reasonable trade-off between capturing meaningful accuracy changes and avoiding short-term noise and fluctuations that may not reflect long-term trends.

*c) The 1-day accuracy standard deviation:* measures the variation of the architecture's performance across experiments, as discussed in Section III-B. A lower standard deviation indicates that the architecture produces consistent results on hardware deployments, which is essential for real-world applications.

large number of non-dominated solutions that can be difficult to interpret. Secondly, by using constrained optimization, we can explicitly incorporate the specific constraints of the analog hardware in our search strategy. This enables us to find DNN architectures that are optimized for the unique requirements and characteristics of analog IMC hardware, rather than simply optimizing for multiple objectives.

## IV. ANALOG-NAS

The objective of AnalogNAS is to find an efficient network architecture under different analog IMC hardware constraints. AnalogNAS comprises three main components: (i) a resnet-like search space, (ii) an analog-accuracy surrogate model, and (iii) an evolutionary-based search strategy. We detail each component in the following subsections.

### A. Resnet-like Search Space

Resnet-like architectures have inspired many manually designed SOTA DL architectures, including Wide ResNet [9] and EfficientNet [28]. Their block-wise architecture offers a flexible and searchable macro-architecture for NAS [29]. Resnet-like architectures can be implemented efficiently using IMC processors, as they are comprised of a large number of MVM and element-wise operations. Additionally, due to the highly parallel nature of IMC, Resnet architectures can get free processing of additional input/output channels. This makes Resnet-like architectures highly amenable to analog implementation.

Fig. 3 depicts the macro-architecture used to construct all architectures in our search space. The architecture consists of

TABLE I
SEARCHABLE HYPER-PARAMETERS AND THEIR RESPECTIVE RANGES.

| Hyper-parameter | Definition | Range |
|---|---|---|
| $OC_0$ | First layer's output channel | Discrete Uniform [8, 128] |
| $KS_0$ | First layer's kernel size | Discrete Uniform [3, 7] |
| M | Number of main blocks | Discrete Uniform [1, 5] |
| R* | Number of residual block per main block | Discrete Uniform [1, 16] |
| B* | Number of branches per main block | Discrete Uniform [1, 12] |
| CT* | Convolution block type per main block | Uniform Choice [A; B; C; D] |
| WF* | Widening factor per main block | Uniform [1, 4] |

*The hyper-parameter is repeated for each main block. ConvBlock refers to different Conv-Relu-BN blocks.

To build the surrogate model, we follow two steps: Dataset creation and Model training:

*2) Dataset Creation:* The surrogate model will predict the rank based on the 1-day accuracy and estimates the AVM and 1-day accuracy standard deviation using the Mean Squared Error (MSE). Since the search space is large, care has to be taken when sampling the dataset of architectures that will be used to train the surrogate model.

The architectures of the search space are sampled using two methods: (i) Latin Hypercube Sampling (LHS) [31] and (ii) NAS with full training. A more detailed description of the AnalogNAS algorithm is presented in Section V. We use LHS to sample architectures distributed evenly over the search space. This ensures good overall coverage of different architectures and their accuracies. NAS with full training is performed using an evolutionary algorithm to collect high-performance architectures. This ensures good exploitation when reaching well-performing regions. In Fig. 4, we present a visualization of the search space coverage, which does not show any clustering of similarly performing architectures at the edge of the main cloud of points. Thus, it is not evident that architectures with similar performance are located close to each other in the search space. This suggests that traditional search methods that
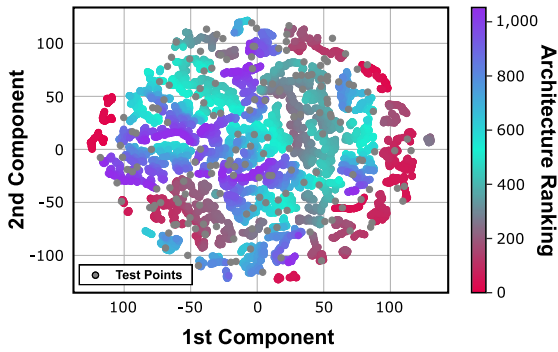


Fig. 4. t-Distributed Stochastic Neighbor Embedding (t-SNE) visualization of the sampled architectures for CIFAR-10.
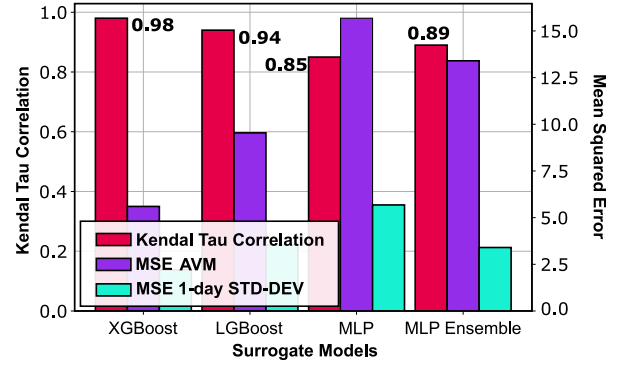


Fig. 5. Surrogate models comparison.

rely on local optimization may not be effective in finding the best-performing architectures. Instead, population-based search strategies, which explore a diverse set of architectures, could be more effective in finding better-performing architectures. Our search strategy extracted 400 test points, and we found that architectures were distributed throughout the main cloud, indicating that our dataset covers a diverse portion of the search space, despite the limited size of only 1,000.

Each sampled architecture is trained using different levels of weight noise and HWA training hyper-parameters using the AIHWKIT [21]. Specifically, we modify the standard deviation of the added weight noise between [0.1, 5.0] in increments of 0.1. The tile size was assumed to be symmetric and varied in [256, 512], representing 256-by-256 and 512-by-512 arrays respectively. Training with different configurations allowed us to generalize the use of the surrogate model across a range of IMC hardware configurations, and to increase the size of the constructed dataset.

*3) Model training:* To train the surrogate model, we used a hinge pair-wise ranking loss [32] with margin $m = 0.1$. The hinge loss, defined in Eq. (2), allows the model to learn the relative ranking order of architectures rather than the absolute accuracy values [32], [33].

$$L(\{a_j, y_j\}_{j=1,...,N}) = \sum_{j=1}^{N} \sum_{\{i,j|y_i>y_j\}} \max[0, m - P(a_i) - P(a_j)] \quad (2)$$

$a_j$ refers to architectures indexed $j$, and $y_j$ to its corresponding 1-day accuracy. $P(a)$ is the predicted score of architecture $a$. $P(a)$ during training, the output score is trained to be correlated with the actual ranks of the architectures. Several algorithms were tested. After an empirical comparison, we adopted Kendall's Tau ranking correlation [34] as the direct criterion for evaluating ranking surrogate model performance. Fig. 5 shows the comparison using different ML algorithms to predict the rankings and AVMs. Our dataset is tabular. It contains each architecture and its corresponding features. XGBoost outperforms the different surrogate models in predicting the architectures' ranking order, the AVM of each architecture, and the 1-day standard deviation.
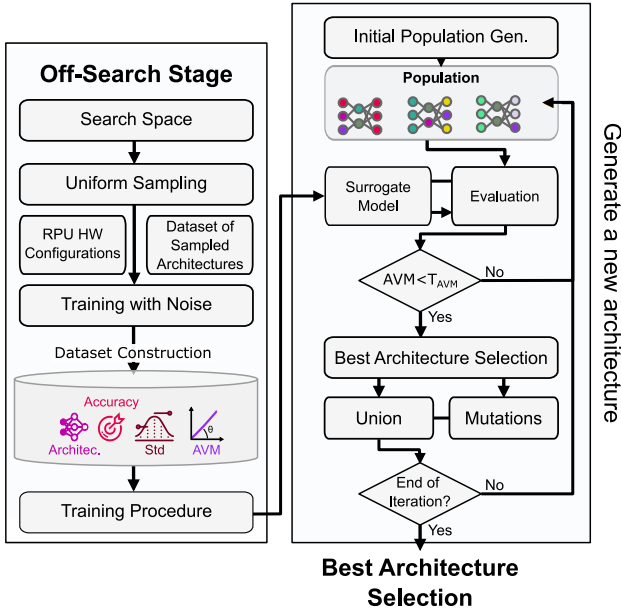
Fig. 6. Overview of the AnalogNAS framework.

**Algorithm 1** AnalogNAS algorithm.

**Input:** Search space: $S$, RPU Configuration: $rpu\_config$, target task: $task$, population size: $population\_size$, AVM threshold: $T_{AVM}$, parameter threshold: $T_p$, number of iterations: $N$, $time\_budget$
**Output:** Most efficient architecture for $rpu\_config$ in $S$
**Begin**
$D$ = sample($S$, dataset_size)
HW_Train($D$, $task$)
AVMs = compute_AVM($D$)
surrogate_model = XGBoost
train(surrogate_model, $D$, AVMs)
**repeat**
    population = LHS(population_size, $T_p$)
    AVMs, ranks = surrogate_model(population)
**until** AVMs > $T_{AVM}$
**while** $i < N$ or $time < time\_budget$ **do**
    top_50 = select(population, ranks)
    mutated = mutation(top_50, $T_p$)
    population = top_50 $\bigcup$ mutated
    AVMs, ranks = surrogate_model(population)
**end while**
**return** $top_1$(population, ranks)

## V. SEARCH STRATEGY

Fig. 6 depicts the overall search framework. Given a dataset and a hardware configuration readable by AIHWKIT, the framework starts by building the surrogate model presented in Section IV-B. Then, we use an optimized evolutionary search to efficiently explore the search space using the surrogate model. Similar to traditional evolutionary algorithms, we use real number encoding. Each architecture is encoded into a vector, and each element of the vector contains the value of the hyper-parameter, as listed in Table I.

### A. Problem Formulation

Given the search space $S$, our goal is to find an architecture $\alpha$, that maximizes the 1-day accuracy while minimizing the 1-day standard deviation, subject to constraints on the number of parameters and the AVM. The number of parameters is an important metric in IMC, because it directly impacts the amount of on-chip memory required to store the weights of a DNN. Eq. (3) formally describes the optimization problem as follows:

$$\max_{\alpha \in S} \quad \frac{\text{ACC}(\alpha)}{\sigma(\alpha)} \tag{3}$$
$$\text{s.t} \quad \psi(\alpha) < T_p, \quad \text{AVM}(\alpha) < T_{AVM}$$

ACC refers to the 1-day accuracy objective, $\sigma$ denotes the 1-day accuracy's standard deviation, and $\psi$ is the number of parameters. $T_p$ and $T_{AVM}$ are user-defined thresholds that correspond to the maximum number of parameters and AVM, respectively.

### B. Search Algorithm

Our evolutionary search algorithm, i.e., AnalogNAS, is formally defined using Algorithm 1. AnalogNAS is an algorithm to find the most accurate and robust neural network architecture

for a given analog IMC configuration and task. The algorithm begins by generating a dataset of neural network architectures, which are trained on the task and evaluated using AIHWKIT. A surrogate model is then created to predict the efficiency of new architectures. The algorithm then generates a population of architectures using an LHS technique and selects the top-performing architectures to be mutated and generate a new population. The process is repeated until a stopping criterion is met, such as a maximum number of iterations or a time budget. Finally, the most robust architecture is returned. In the following, we detail how the population initialization, fitness evaluation, and mutations are achieved.

*1) Population Initialization:* The search starts by generating an initial population. Using the LHS algorithm, we sample the population uniformly from the search space. LHS ensures that the initial population contains architectures with different architectural features. LHS is made faster with parallelization by dividing the sampling into multiple independent subsets, which can be generated in parallel using multiple threads.

*2) Fitness Evaluation:* We evaluate the population using the aforementioned analog-accuracy surrogate model. In addition to the rankings, the surrogate model predicts the AVM of each architecture. As previously described, the AVM is used to gauge the robustness of a given network. If the AVM is below a defined threshold, $T_{AVM}$, the architecture is replaced by a randomly sampled architecture. The new architecture is constrained to be sampled from the same hypercube dimension as the previous one. This ensures efficient exploration.

*3) Selection and Mutation:* We select the top 50% architectures from the population using the predicted rankings. These architectures are mutated. The mutation functions are classified

as follows:

*a) Depth-related mutations:* modify the depth of the architectures. Mutations include adding a main block, by increasing or decreasing $M$ or a residual block $R$, or modifying the type of convolution block, i.e., $\{A, B, C, D\}$, for each main block.

*b) Width-related mutations:* modify the width of the architectures. Mutations include modifying the widening factor $W$ of a main block or adding or removing a branch $B$, or modifying the initial output channel size of the first convolution, $OC$.

*c) Other mutations:* modify the kernel size of the first convolution, $KS$, and/or add skip connections, denoted using $ST$.

Depth- and width-related mutations are applied with the same probability of 80%. The other mutations are applied with a 50% probability. In each class, the same probability is given to each mutation. The top 50% architectures in addition to the mutated architectures constitute the new population. For the remaining iterations, we verify the ranking correlation of the surrogate model. If the surrogate model's ranking correlation is degraded, we fine-tune the surrogate model with the population's architectures. The degradation is computed every 100 iterations. The surrogate model is tested on the population architectures after training them. It is fine-tuned if Kendall's tau correlation drops below 0.9.

## VI. Experiments

This section describes the experiments used to evaluate AnalogNAS on three tasks: CIFAR-10 image classification, VWW, and KWS. The AIHWKIT was used to perform hardware simulations.

### A. Experimental Setup

*1) Training Details:* We detail the hyper-parameters used to train the surrogate model and different architectures on CIFAR-10, VWW, and KWS tasks.

*a) Surrogate model training:* We trained a surrogate model and dataset of HWA trained DNN architctures for each task. The sizes of the datasets were 1,200, 600, and 1,500, respectively. An additional 500 architectures were collected during the search trials for validation. All architectures were first trained without noise injection (i.e., using vanilla training routines), and then converted to AIHWKIT models for HWA retraining. The surrogate model architecture used was XGBoost. For VWW and KWS, the surrogate model was fine-tuned from the image classification XGBoost model.

*b) Image classification training:* We first trained the network architectures using the CIFAR-10 dataset [35], which contains 50,000 training and 10,000 test samples, evenly distributed across 10 classes. We augmented the training images with random crops and cutouts only. For training, we used Stochastic Gradient Descent (SGD) with a learning rate of 0.05 and a momentum of 0.9 with a weight decay of 5e-4. The learning rate was adjusted using a cosine annealing learning rate scheduler with a starting value of 0.05 and a maximum number of 400 iterations.

*c) Visual Wake Words (VWW) training:* We first trained the network architectures using the VWW dataset [36], which contains 82,783 train and 40,504 test images. Images are labeled 1 when a person is detected, and 0 when no person is present. The image pre-processing pipeline includeded horizontal and vertical flipping, scale augmentation [37], and random Red Green Blue (RGB) color shift. To train the architectures, we used the RMSProp optimizer [38] with a momentum of 0.9, a learning rate of 0.01, a batch normalization momentum of 0.99, and a $l_2$ weight decay of 1e-5.

*d) Keyword Spotting (KWS) training:* We first trained the network architectures using the KWS dataset [39], which contains 1-second long incoming audio clips. These are classified into one of twelve keyword classes, including "silence" and "unknown" keywords. The dataset contains 85,511 training, 10,102 validation, and 4,890 test samples. The input was transformed to $49 \times 10 \times 1$ features from the Mel-frequency cepstral coefficients [40]. The data pre-processing pipeline included applying background noise and random timing jitter. To train the architectures, we used the Adam optimizer [41] with a decay of 0.9, a learning rate of 3e-05, and a linear learning rate scheduler with a warm-up ratio of 0.1.

*2) Search Algorithm:* The search algorithm was run five times to compute the variance. The evolutionary search was executed with a population size of 200. If not explicitly mentioned, the AVM threshold was set to 10%. The width and depth mutation probability was set to 0.8. The other mutations' probability was set to 0.5. The total number of iterations was 200. After the search, the obtained architecture for each task was compared to SOTA baselines for comparison.

### B. Results

The final architecture compositions for the three tasks are listed in Table II. In addition, figure 10 highlights the architectural differences between AnalogNAS_T500 and resnet32. We modified $T_p$ to find smaller architectures. To determine the optimal architecture for different parameter thresholds, we use T$X$, where $X$ represents the threshold $T_p$ in K units (e.g., T100 refers to the architecture with a threshold of 100K parameters). When searching for T200 and T100, the probability of increasing the widening factor or depth to their highest values, was lessened to 0.2.

In Fig. 7, the simulated hardware comparison of the three tasks is depicted. Our models outperform SOTA architectures with respect to both accuracy and resilience to drift. On CIFAR-10, after training the surrogate model, the search took 17 minutes to run. We categorize the results based on the number of parameters threshold into two distinct groups. We consider edge models with a number of parameters below 1M and above 400k. Below 400K, architectures are suitable for TinyML deployment. The final architecture, T500, is smaller than Resnet32, and achieved +1.86% better accuracy and a drop of 1.8% after a month of inference, compared to 5.04%. This model is $\sim 86\times$ smaller than Wide Resnet [9], which has 36.5M parameters. Our smallest model, T100,
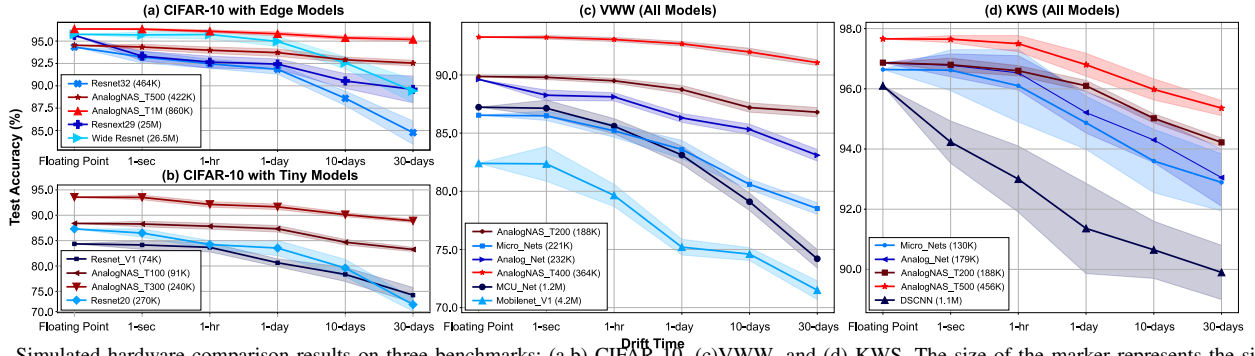
Fig. 7. Simulated hardware comparison results on three benchmarks: (a,b) CIFAR-10, (c)VWW, and (d) KWS. The size of the marker represents the size (i.e., the number of parameters) of each model. The shaded area corresponds to the standard deviation at that time.

TABLE II
FINAL ARCHITECTURES FOR CIFAR-10, VWW, AND KWS. OTHER NETWORKS FOR VWW AND KWS ARE NOT LISTED, AS THEY CANNOT EASILY BE REPRESENTED USING OUR MACRO-ARCHITECTURE.

| Network | $OC_0$ | $KS_0$ | Macro-Architecture Parameter | | | | |
| | | | M | R* | B* | CT* | WF* |
|---|---|---|---|---|---|---|---|
| **CIFAR-10** | | | | | | | |
| Resnet32 | 64 | 7 | 3 | (5, 5, 5) | (1, 1, 1) | (B, B, B) | (1, 1, 1) |
| AnalogNAS_T100 | 32 | 3 | 1 | (2, ) | (1,) | (C, ) | (2, ) |
| AnalogNAS_T300 | 32 | 3 | 1 | (3, 3) | (1, 1) | (A, B) | (2, 1) |
| AnalogNAS_T500 | 64 | 5 | 1 | (3, ) | (3, ) | (A, ) | (2, ) |
| AnalogNAS_T1M | 32 | 5 | 2 | (3, 3) | (2, 2) | (A, A) | (3, 3) |
| **VWW** | | | | | | | |
| AnalogNAS_T200 | 24 | 3 | 3 | (2, 2, 2) | (1, 2, 1) | (B, A, A) | (2, 2, 2) |
| AnalogNAS_T400 | 68 | 3 | 2 | (3, 5) | (2, 1) | (C, C) | (3, 2) |
| **KWS** | | | | | | | |
| AnalogNAS_T200 | 80 | 1 | 1 | (1, ) | (2, ) | (C, ) | (4, ) |
| AnalogNAS_T400 | 68 | 1 | 2 | (2, 1) | (1, 2) | (B, B) | (3, 3) |

*As depicted in Fig. 3, thse macro-architecture parameters comprise multiple instances.

was $1.23\times$ bigger than Resnet-V1, the SOTA model benchmarked by MLPerf [11]. Despite not containing any depth-wise convolutions, Resnet_V1 is extremely small, with only 70k parameters. Our model offers a +7.98% accuracy increase with a 5.14% drop after a month of drift compared to 10.1% drop for Resnet_V1. Besides, our largest model, *AnalogNAS_1M*, outperforms Wide Resnet with +0.86% in the 1-day accuracy with a drop of only 1.16% compared to 6.33%. In addition, the found models exhibit greater consistency across experiment trials, with an average standard deviation of 0.43 over multiple drift times as opposed to 0.97 for SOTA models.

Similar conclusions can be made about VWW and KWS. In VWW, current baselines use a depth-wise separable convolution that incurs a high accuracy drop on analog devices. Compared to AnalogNet-VWW and Micronets-VWW, the current SOTA networks for VWW in analog and edge devices, our T200 model has similar number of parameters (x1.23 smaller) with a +2.44% and +5.1% 1-day accuracy increase respectively. AnalogNAS was able to find more robust and consistent networks with an average AVM of 2.63% and a standard deviation of 0.24. MCUNet [42] and MobileNet-V1 present the highest AVM. This is due to the sole use of depth-wise separable convolutions.

On KWS, the baseline architectures, including DSCNN [43],

use hybrid networks containing recurrent cells and convolutions. The recurrent part of the model ensures high robustness to noise. While current models are already robust with an average accuracy drop of 4.72%, our model outperforms tiny SOTA models with 96.8% and an accuracy drop of 2.3% after a month of drift. Critically, our AnalogNAS models exhibit greater consistency across experiment trials, with an average standard deviation of 0.17 over multiple drift times as opposed to 0.36 for SOTA models.

### C. Comparison with HW-NAS

In accordance with commonly accepted NAS methodologies, we conducted a comparative analysis of our search approach with Random Search. Results, presented in Fig. 8, were obtained across five experiment instances. Our findings indicate that Random Search was unable to match the 1-day accuracy levels of our final models, even after conducting experiments for a duration of four hours and using the same surrogate model. We further conducted an ablation study to evaluate the effectiveness of our approach by analyzing the impact of the LHS algorithm and surrogate model. The use of a
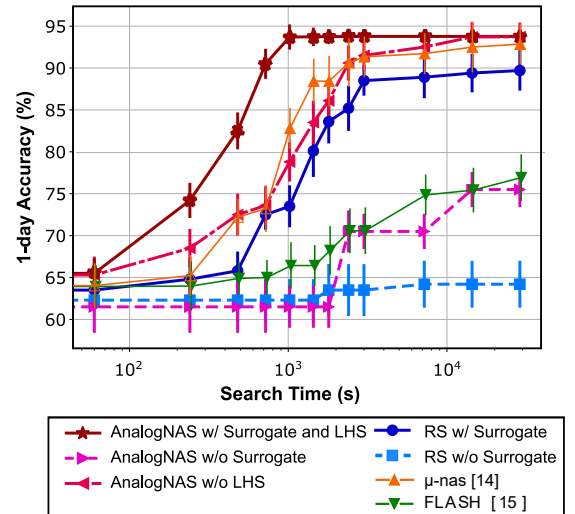


Fig. 8. Ablation study comparison against HW-NAS. Mean and standard deviation values are reported across five experiment instances (trials).

TABLE III
AVM THRESHOLD VARIATION RESULTS ON CIFAR-10.

| $T_{AVM}$ (%) | 1.0 | 3.0 | 5.0* |
|---|---|---|---|
| 1-day Accuracy | 88.7% | 93.71% | 93.71% |
| AVM | 0.85% | 1.8% | 1.8% |
| Search Time (min) | 34.65 | 28.12 | 17.65 |

*Overall results computed with $T_{AVM}$ (%) = 5.0.

random sampling strategy and exclusion of the surrogate model resulted in a significant increase in search time. The LHS algorithm helped in starting from a diverse initial population and improving exploration efficiency, while the surrogate model played a crucial role in ensuring practical search times.

Besides, AnalogNAS surpasses both FLASH [15] and $\mu$-nas [14] in performance and search time. FLASH search strategy is not adequate for large search spaces such as ours. As for $\mu$-nas, although it manages to achieve acceptable results, its complex optimization algorithm hinders the search process, resulting in decreased efficiency.

### D. Search Time and Accuracy Variation over One Month (AVM) Threshold Trade-Off

During the search, we culled architectures using their predicted AVM, i.e., any architecture with a higher AVM than the AVM threshold was disregarded. As listed in Table III, we varied this threshold to investigate the trade-off between $T_{AVM}$ and the search time. As can be seen, as $T_{AVM}$ is decreased, the delta between AVM and $T_{AVM}$ significantly decreases. The correlation between the search time and $T_{AVM}$ is observed to be non-linear.

### VII. EXPERIMENTAL HARDWARE VALIDATION AND ARCHITECTURE PERFORMANCE SIMULATIONS

### A. Experimental Hardware Validation

An experimental hardware accuracy validation study was performed using a 64-core IMC chip based on PCM [44]. Each core comprises a crossbar array of 256x256 PCM-based unit-cells along with a local digital processing unit [45]. This validation study was performed to verify whether the simulated network accuracy values and rankings are representative of those when the networks are deployed on real physical hardware. We deployed two networks for the CIFAR-10 image classification task on hardware: AnalogNAS_T500 and the baseline ResNet32 [30] networks from Fig. 7(a).

To implement the aforementioned models on hardware, after HWA training was performed, a number of steps were carried out. First, from the AIHWKIT, unit weights of linear (dense) and unrolled convolutional layers, were exported to a state dictionary file. This was used to map network parameters to corresponding network layers. Additionally, the computational inference graph of each network was exported. These files were used to generate proprietary data-flows to be executed in-memory. As only hardware accuracy validation was being performed, all other operations aside from MVMs were performed on a host machine connected to the chip

| Architecture | ResNet32 | AnalogNAS_T500 |
|---|---|---|
| **Hardware Experiments** | | |
| FP Accuracy | 94.34% | 94.54% |
| Hardware accuracy* | 89.55% | 92.07% |
| **Simulated Hardware Power Performance** | | |
| Total weights | 464,432 | 416,960 |
| Total tiles | 43 | 27 |
| Network Depth | 32 | 17 |
| Execution Time (msec) | 0.434 | 0.108 |
| Inferences/s/W | 43,956.7 | 54,502 |

*The mean accuracy is reported across five experiment repetitions.

through an Field-programmable Gate Array (FPGA). The measured hardware accuracy was 92.05% for T500 and 89.87% for Resnet32, as reported in Table IV. Hence, the T500 network performs significantly better than Resnet32 also when implemented on real hardware. This further validates that our proposed AnalogNAS approach is able to find networks with similar number of parameters that are more accurate and robust on analog IMC hardware.

### B. Simulated Hardware Energy and Latency

We conducted power performance simulations for Analog-NAS_T500 and ResNet32 models using a 2D-mesh based heterogeneous analog IMC system with the simulation tool presented in [46]. The simulated IMC system consists of one analog fabric with 48 analog tiles of 512x512 size, on-chip digital processing units, and digital memory for activation orchestration between CNN layers. Unlike the accuracy validation experiments on the 64-core IMC chip, the simulated power performance assumes all intermediate operations to be mapped and executed on-chip. Our results, provided in Table IV, show that AnalogNAS_T500 outperformed ResNet32 in terms of both execution time and energy efficiency.

We believe that this power performance benefit is realized because, in analog IMC hardware, wider layers can be computed in parallel, leveraging the $O(1)$ latency from analog tiles, and are therefore preferred over deeper layers. It is noted that both networks exhibit poor tile utilization and that the tile utilization and efficiency of these networks could be further improved by incorporating these metrics as explicit constraints. This is left to future work and is beyond the scope of AnalogNAS.

### VIII. DISCUSSION

During the search, we analyzed the architecture characteristics and studied which types of architectures perform the best on IMC inference processors. The favoured architectures combine robustness to noise and accuracy performance. Fig. 9 shows the evolution of the average depth, the average widening factor, the average number of branches, and the average first convolution's output channel size of the search population for every 20 iterations. The depth represents the number of convolutions.
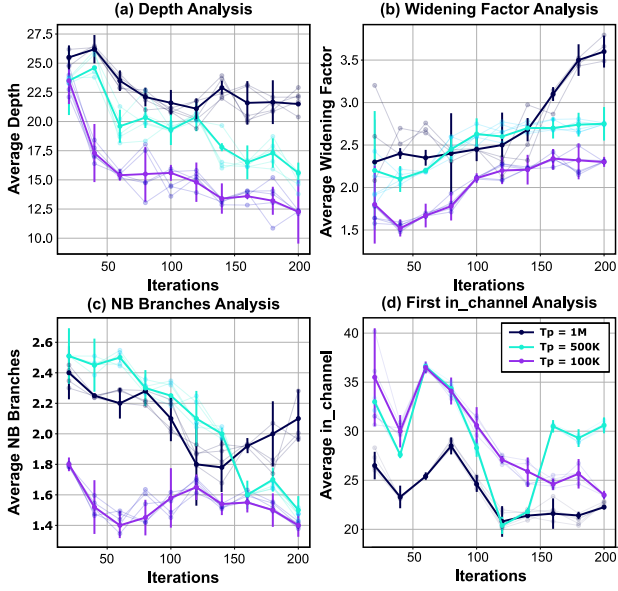
Fig. 9. Evolution of architecture characteristics in the population during the search for CIFAR-10. Random individual networks are shown.
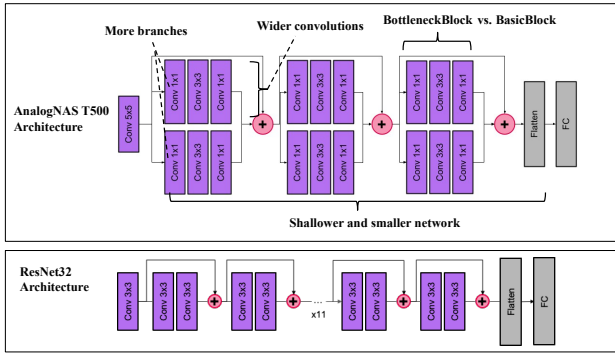


Fig. 10. Architectural differences between AnalogNAS_T500 and Resnet32.

A sampled architecture has a widening factor per block. To compute the average widening factor, we first computed the average widening factor per architecture by dividing the sum of the widening factors by the number of blocks contained in the architecture. Then, we calculated the average widening factor across all architectures. Similar computations were performed for the average number of branches.

For each plot, the search was run 5 times and the mean is represented in each point. The plotted error corresponds to one standard deviation from that mean. Starting from a random population obtained using LHS, the population evolves through different width and depth-related mutations. During this analysis, we want to answer the following questions: (i) does the search favor wide or deep networks? And subsequently, are wider architectures more noise resilient? (ii) what architectures are exploited by the search for different tasks when constraining the number of parameters?

## A. Are Wider or Deeper Networks More Robust to PCM Device Drift?

From Fig. 9, it can be observed that the depth of all networks decreases during the search. This trend is especially seen when we constrain the model's size to 100K and 500K parameters. During the search, the widening factor also increases, allowing the blocks to have wider convolutions. The number of branches is highly dependent on $T_p$. This number is, on average, between 1 and 2. The branches are the number of parallel convolutions in a block disregarding the skip connection. In the literature, architectures such as ResNext, that support a higher number of branches, have a number of parameters around 30M. It is still interesting to get blocks with two branches, which also reflects an increase in the width of the network by increasing the number of features extracted within the same block.The average output channel size of the first convolution decreases during the search. Its final value is around the same number of output channels as standard architectures, i.e., between 16 and 32. This follows the general trend of having wider convolutions in deeper network positions.

## B. Types Of Architectures

The architectures and parameter constraints differ for each task, but they all exhibit an increasing expansion ratio in the convolution block. This allows the convolutions to effectively utilize the tile and mitigate noise from unused cells in the crossbar.For CIFAR-10, architectures behave like Wide Resnet [9] while respecting the number of parameters constraint. For the VWW task, the architectures are deeper. The input resolution is $224 \times 224$, which requires more feature extraction blocks. However, they are still smaller than SOTA architectures, with a maximum depth of 22. As depth is essential to obtain high accuracy for the VWW task, no additional branches are added. For the KWS task, the architectures are the widest possible, maximizing the tile utilization for each convolutional layer.

## IX. CONCLUSION

In this paper, we propose an efficient NAS methodology dedicated to analog in-memory computing for TinyML tasks entitled AnalogNAS. The obtained models are accurate, noise and drift-resilient, and small enough to run on resource-constrained devices. Experimental results demonstrate that our method outperforms SOTA models on analog hardware for three tasks of the MLPerf benchmark: image classification on CIFAR-10, VWW, and KWS. Our AnalogNAS_T500 model implemented on physical hardware demonstrates $> 2\%$ higher accuracy experimentally on the CIFAR-10 benchmark than ResNet32. Calculated speed and energy efficiency estimates reveal a $> 4\times$ reduction in execution time, in addition to $> 1.2\times$ higher energy efficiency for AnalogNAS_T500 compared with ResNet32 when evaluated using a system-level simulator. While our paper has focused on a ResNet-like search space, it is important to note that our search strategy is adaptable and can be extended in future work to explore a broader range of architectures.

REFERENCES

[1] A. Sebastian, M. Le Gallo, R. Khaddam-Aljameh, and E. Eleftheriou, "Memory devices and applications for in-memory computing," *Nature Nanotechnology*, vol. 15, no. 7, pp. 529–544, Jul 2020.

[2] V. Joshi, M. L. Gallo, I. Boybat, S. Haefeli, C. Piveteau, M. Dazzi, B. Rajendran, A. Sebastian, and E. Eleftheriou, "Accurate deep neural network inference using computational phase-change memory," *CoRR*, vol. abs/1906.03138, 2019.

[3] C. Lammie, W. Xiang, B. Linares-Barranco, and M. Rahimi Azghadi, "Memtorch: An open-source simulation framework for memristive deep learning systems," *Neurocomputing*, vol. 485, pp. 124–133, 2022.

[4] I. Boybat, B. Kersting, S. G. Sarwat, X. Timoneda, R. L. Bruce, M. BrightSky, M. L. Gallo, and A. Sebastian, "Temperature sensitivity of analog in-memory computing using phase-change memory," in *2021 IEEE International Electron Devices Meeting (IEDM)*, 2021, pp. 28.3.1–28.3.4.

[5] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *CoRR*, vol. abs/1704.04861, 2017.

[6] G. Lu, W. Zhang, and Z. Wang, "Optimizing depthwise separable convolution operations on gpus," *IEEE Transactions on Parallel and Distributed Systems*, vol. 33, no. 1, pp. 70–87, 2022.

[7] H. Benmeziane, K. E. Maghraoui, H. Ouarnoughi, S. Niar, M. Wistuba, and N. Wang, "A comprehensive survey on hardware-aware neural architecture search," *CoRR*, vol. abs/2101.09336, 2021.

[8] S. Xie, R. B. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*. IEEE Computer Society, 2017, pp. 5987–5995.

[9] S. Zagoruyko and N. Komodakis, "Wide residual networks," in *Proceedings of the British Machine Vision Conference, BMVC*, R. C. Wilson, E. R. Hancock, and W. A. P. Smith, Eds. BMVA Press, 2016.

[10] L. Sekanina, "Neural architecture search and hardware accelerator co-search: A survey," *IEEE Access*, vol. 9, pp. 151 337–151 362, 2021.

[11] C. R. Banbury, V. J. Reddi, P. Torelli, N. Jeffries, C. Király, J. Holleman, P. Montino, D. Kanter, P. Warden, D. Pau, U. Thakker, A. Torrini, J. Cordaro, G. D. Guglielmo, J. M. Duarte, H. Tran, N. Tran, W. Niu, and X. Xu, "Mlperf tiny benchmark," in *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual*, J. Vanschoren and S. Yeung, Eds., 2021.

[12] C. R. Banbury, C. Zhou, I. Fedorov, R. M. Navarro, U. Thakker, D. Gope, V. J. Reddi, M. Mattina, and P. N. Whatmough, "Micronets: Neural network architectures for deploying tinyml applications on commodity microcontrollers," in *Proceedings of Machine Learning and Systems*, A. Smola, A. Dimakis, and I. Stoica, Eds. mlsys.org, 2021.

[13] H. Liu, K. Simonyan, and Y. Yang, "DARTS: differentiable architecture search," *CoRR*, vol. abs/1806.09055, 2018. [Online]. Available: http://arxiv.org/abs/1806.09055

[14] E. Liberis, L. Dudziak, and N. D. Lane, "µnas: Constrained neural architecture search for microcontrollers," in *EuroMLSys@EuroSys 2021, Proceedings of the 1st Workshop on Machine Learning and Systemsg Virtual Event, Edinburgh, Scotland, UK, 26 April, 2021*, E. Yoneki and P. Patras, Eds. ACM, 2021, pp. 70–79.

[15] G. Li, S. K. Mandal, Ü. Y. Ogras, and R. Marculescu, "FLASH: fast neural architecture search with hardware optimization," *ACM Trans. Embed. Comput. Syst.*, vol. 20, no. 5s, pp. 63:1–63:26, 2021.

[16] W. Jiang, Q. Lou, Z. Yan, L. Yang, J. Hu, X. S. Hu, and Y. Shi, "Device-circuit-architecture co-exploration for computing-in-memory neural accelerators," *IEEE Trans. Computers*, vol. 70, no. 4, pp. 595–605, 2021.

[17] Z. Yuan, J. Liu, X. Li, L. Yan, H. Chen, B. Wu, Y. Yang, and G. Sun, "NAS4RRAM: neural network architecture search for inference on rram-based accelerators," *Sci. China Inf. Sci.*, vol. 64, no. 6, 2021.

[18] Z. Yan, D.-C. Juan, X. S. Hu, and Y. Shi, "Uncertainty modeling of emerging device based computing-in-memory neural accelerators with application to neural architecture search," in *2021 26th Asia and South Pacific Design Automation Conference (ASP-DAC)*, 2021, pp. 859–864.

[19] G. Li, M. Zhang, J. Li, F. Lv, and G. Tong, "Efficient densely connected convolutional neural networks," *Pattern Recognit.*, vol. 109, p. 107610, 2021.

[20] C. Zhou, F. Redondo, J. Buchel, I. Boybat, X. Comas, S. R. Nandakumar, S. Das, A. Sebastian, M. L. Gallo, and P. N. Whatmough, "Ml-hw co-design of noise-robust tinyml models and always-on analog compute-in-memory edge accelerator," *IEEE Micro*, vol. 42, no. 06, pp. 76–87, 2022.

[21] M. J. Rasch, D. Moreda, T. Gokmen, M. L. Gallo, F. Carta, C. Goldberg, K. E. Maghraoui, A. Sebastian, and V. Narayanan, "A flexible and fast pytorch toolkit for simulating training and inference on analog crossbar arrays," in *3rd IEEE International Conference on Artificial Intelligence Circuits and Systems*. IEEE, 2021, pp. 1–4.

[22] M. J. Rasch, C. Mackin, M. L. Gallo, A. Chen, A. Fasoli, F. Odermatt, N. Li, S. R. Nandakumar, P. Narayanan, H. Tsai, G. W. Burr, A. Sebastian, and V. Narayanan, "Hardware-aware training for large-scale and diverse deep learning inference workloads using in-memory computing-based accelerators," 2023. [Online]. Available: https://arxiv.org/abs/2302.08469

[23] C. Termritthikun, Y. Jamtsho, J. Ieamsaard, P. Muneesawang, and I. Lee, "Eeea-net: An early exit evolutionary neural architecture search," *Eng. Appl. Artif. Intell.*, vol. 104, p. 104397, 2021.

[24] X. Dai, A. Wan, P. Zhang, B. Wu, Z. He, Z. Wei, K. Chen, Y. Tian, M. Yu, P. Vajda, and J. E. Gonzalez, "Fbnetv3: Joint architecture-recipe search using predictor pretraining," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*. Computer Vision Foundation / IEEE, 2021, pp. 16 276–16 285.

[25] W. Jiang, X. Zhang, E. H. Sha, L. Yang, Q. Zhuge, Y. Shi, and J. Hu, "Accuracy vs. efficiency: Achieving both through fpga-implementation aware neural architecture search," in *Proceedings of the 56th Annual Design Automation Conference 2019, DAC 2019, Las Vegas, NV, USA, June 02-06, 2019*. ACM, 2019, p. 5.

[26] A. Sarah, D. Cummings, S. N. Sridhar, S. Sundaresan, M. Szankin, T. Webb, and J. P. Munoz, "A hardware-aware system for accelerating deep neural network optimization," *CoRR*, vol. abs/2202.12954, 2022.

[27] Z. Dong, Y. Gao, Q. Huang, J. Wawrzynek, H. K. H. So, and K. Keutzer, "HAO: hardware-aware neural architecture optimization for efficient inference," in *29th IEEE Annual International Symposium on Field-Programmable Custom Computing Machines, FCCM 2021, Orlando, FL, USA, May 9-12, 2021*. IEEE, 2021, pp. 50–59.

[28] M. Tan and Q. V. Le, "Efficientnetv2: Smaller models and faster training," in *Proceedings of the 38th International Conference on Machine Learning, ICML*, ser. Proceedings of Machine Learning Research, M. Meila and T. Zhang, Eds., vol. 139. PMLR, 2021, pp. 10 096–10 106.

[29] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le, "Learning transferable architectures for scalable image recognition," in *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*. Computer Vision Foundation / IEEE Computer Society, 2018, pp. 8697–8710.

[30] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*. IEEE Computer Society, 2016, pp. 770–778.

[31] M. D. McKay, "Latin hypercube sampling as a tool in uncertainty analysis of computer models," in *Proceedings of the 24th Winter Simulation Conference, Arlington, VA, USA, December 13-16, 1992*, R. C. Crain, Ed. ACM Press, 1992, pp. 557–564.

[32] X. Ning, Y. Zheng, T. Zhao, Y. Wang, and H. Yang, "A generic graph-based neural architecture encoding scheme for predictor-based NAS," in *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part XIII*, ser. Lecture Notes in Computer Science, A. Vedaldi, H. Bischof, T. Brox, and J. Frahm, Eds., vol. 12358. Springer, 2020, pp. 189–204.

[33] H. Benmeziane, S. Niar, H. Ouarnoughi, and K. E. Maghraoui, "Pareto rank surrogate model for hardware-aware neural architecture search," in *International IEEE Symposium on Performance Analysis of Systems and Software, ISPASS*. IEEE, 2022, pp. 267–276.

[34] H. Abdi, "The kendall rank correlation coefficient," *Encyclopedia of Measurement and Statistics. Sage, Thousand Oaks, CA*, pp. 508–510, 2007.

[35] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.

[36] A. Chowdhery, P. Warden, J. Shlens, A. Howard, and R. Rhodes, "Visual wake words dataset," *CoRR*, vol. abs/1906.05721, 2019.

[37] H. Zhang, C. Wu, Z. Zhang, Y. Zhu, H. Lin, Z. Zhang, Y. Sun, T. He, J. Mueller, R. Manmatha, M. Li, and A. J. Smola, "Resnest: Split-attention networks," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, CVPR*. IEEE, 2022, pp. 2735–2745.

[38] M. C. Mukkamala and M. Hein, "Variants of rmsprop and adagrad with logarithmic regret bounds," in *Proceedings of the 34th International Conference on Machine Learning, ICML*, ser. Proceedings of Machine Learning Research, D. Precup and Y. W. Teh, Eds., vol. 70. PMLR, 2017, pp. 2545–2553.

[39] P. Warden, "Speech Commands: A Dataset for Limited-Vocabulary Speech Recognition," *ArXiv e-prints*, Apr. 2018. [Online]. Available: https://arxiv.org/abs/1804.03209

[40] J. Martínez, H. Pérez-Meana, E. E. Hernández, and M. M. Suzuki, "Speaker recognition using mel frequency cepstral coefficients (MFCC) and vector quantization (VQ) techniques," in *22nd International Conference on Electrical Communications and Computers, CONIELECOMP*, P. B. Sánchez, R. Rosas-Romero, and M. J. O. Galindo, Eds. IEEE, 2012, pp. 248–251.

[41] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *3rd International Conference on Learning Representations*, Y. Bengio and Y. LeCun, Eds., 2015.

[42] J. Lin, W. Chen, Y. Lin, J. Cohn, C. Gan, and S. Han, "Mcunet: Tiny deep learning on iot devices," in *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, Eds., 2020.

[43] P. M. Sørensen, B. Epp, and T. May, "A depthwise separable convolutional neural network for keyword spotting on an embedded system," *EURASIP J. Audio Speech Music. Process.*, vol. 2020, no. 1, p. 10, 2020.

[44] M. L. Gallo, R. Khaddam-Aljameh, M. Stanisavljevic, A. Vasilopoulos, B. Kersting, M. Dazzi, G. Karunaratne, M. Braendli, A. Singh, S. M. Mueller, J. Buechel, X. Timoneda, V. Joshi, U. Egger, A. Garofalo, A. Petropoulos, T. Antonakopoulos, K. Brew, S. Choi, I. Ok, T. Philip, V. Chan, C. Silvestre, I. Ahsan, N. Saulnier, V. Narayanan, P. A. Francese, E. Eleftheriou, and A. Sebastian, "A 64-core mixed-signal in-memory compute chip based on phase-change memory for deep neural network inference," 2022. [Online]. Available: https://arxiv.org/abs/2212.02872

[45] R. Khaddam-Aljameh, M. Stanisavljevic, J. F. Mas, G. Karunaratne, M. Braendli, F. Liu, A. Singh, S. M. Müller, U. Egger, A. Petropoulos, T. Antonakopoulos, K. Brew, S. Choi, I. Ok, F. L. Lie, N. Saulnier, V. Chan, I. Ahsan, V. Narayanan, S. R. Nandakumar, M. L. Gallo, P. A. Francese, A. Sebastian, and E. Eleftheriou, "Hermes core – a 14nm cmos and pcm-based in-memory compute core using an array of 300ps/lsb linearized cco-based adcs and local digital processing," in *2021 Symposium on VLSI Technology*, 2021, pp. 1–2.

[46] S. Jain, H. Tsai, C.-T. Chen, R. Muralidhar, I. Boybat, M. M. Frank, S. Woźniak, M. Stanisavljevic, P. Adusumilli, P. Narayanan, K. Hosokawa, M. Ishii, A. Kumar, V. Narayanan, and G. W. Burr, "A heterogeneous and programmable compute-in-memory accelerator architecture for analog-ai using dense 2-d mesh," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 31, no. 1, pp. 114–127, 2023.