

Standing on the Shoulders of Giants: Hardware and Neural Architecture Co-Search with Hot Start

Weiwen Jiang, Lei Yang, Sakyasingha Dasgupta, Jingtong Hu, and Yiyu Shi

Abstract—Hardware and neural architecture co-search that automatically generates Artificial Intelligence (AI) solutions from a given dataset is promising to promote AI democratization; however, the amount of time that is required by current co-search frameworks is in the order of hundreds of GPU hours for one target hardware. This inhibits the use of such frameworks on commodity hardware. The root cause of the low efficiency in existing co-search frameworks is the fact that they start from a “cold” state (i.e., search from scratch). In this paper, we propose a novel framework, namely HotNAS, that starts from a “hot” state based on a set of existing pre-trained models (a.k.a. model zoo) to avoid lengthy training time. As such, the search time can be reduced from 200 GPU hours to less than 3 GPU hours. In HotNAS, in addition to hardware design space and neural architecture search space, we further integrate a compression space to conduct model compressing during the co-search, which creates new opportunities to reduce latency, but also brings challenges. One of the key challenges is that all of the above search spaces are coupled with each other, e.g., compression may not work without hardware design support. To tackle this issue, HotNAS builds a chain of tools to design hardware to support compression, based on which a global optimizer is developed to automatically co-search all the involved search spaces. Experiments on ImageNet dataset and Xilinx FPGA show that, within the timing constraint of 5ms, neural architectures generated by HotNAS can achieve up to 5.79% Top-1 and 3.97% Top-5 accuracy gain, compared with the existing ones.

I. INTRODUCTION

The success of Deep Neural Networks (DNN), has propelled Artificial Intelligence (AI) in entering every aspect of our lives and is being widely employed for diverse applications on different types of hardware. Neural Architecture Search (NAS), a successful product of Automatic Machine Learning (AutoML), has paved the way from a given dataset to a neural architecture with state-of-the-art accuracy. Moving forward, to be able to use AI for enabling and accelerating different applications, we need to be able to design the neural network in a way that the design specifications are met on our target hardware; for instance, real-time constraints for edge devices, low power budgets for IoT devices, etc.

W. Jiang and Y. Shi are with the Department of Computer Science and Engineering, University of Notre Dame, Notre Dame, IN 46556.

L. Yang is with the Department of Electrical and Computer Engineering, University of New Mexico, Albuquerque, NM 87131.

S. Dasgupta is with Edgecortex Inc., Tokyo, Japan, 1410031.

J. Hu is with the Department of Electrical and Computer Engineering, University of Pittsburgh, Pittsburgh, PA 15261.

W. Jiang is the corresponding author (e-mail: wjiang2@nd.edu).

This work is partially supported by the National Science Foundation under Grants CCF-1919167, CCF-1820537 and CNS-1822099.

This article was presented in the International Conference on Hardware/Software Codesign and System Synthesis 2020 and appears as part of the ESWEK-TCAD special issue.

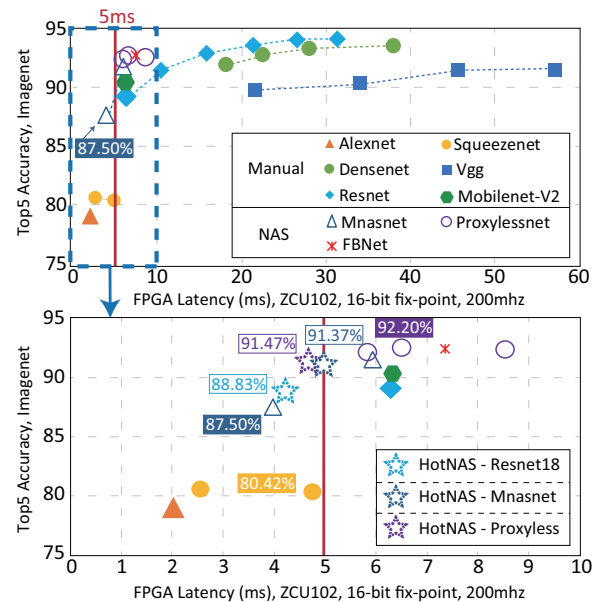


Figure 1. Architecture-hardware co-search by HotNAS: (top) latency and accuracy of models in the model zoo; (bottom) architectures in model zoo and that identified by HotNAS with 5ms timing constraint (Best viewed in color)

Recently, neural architecture and hardware design (abbr. architecture-hardware) co-search frameworks [1]–[13] have been proposed to bridge the gap between neural architecture and hardware design. These frameworks have demonstrated promising results in generating high-accuracy and low-cost systems. However, their search efficiency is low: existing co-search frameworks commonly take hundreds of GPU hours per target hardware. This may become the bottleneck in many emerging applications where fast turn-around or short time-to-market is desired. On the other hand, it has already been shown that the carbon footprint (pounds of CO_2) of NAS for one model is nearly equivalent to five times the lifetime emissions of a car [14]. In this work, we are revisiting the default setting used by existing co-search frameworks, where: the exploration always starts from scratch (i.e., cold start), which results in large search time and low efficiency. However, is a cold start really necessary?

We claim that the architecture-hardware co-search could stand on the shoulders of giants and start the search from a hot state, i.e., using an existing pre-trained model in a model zoo. The model zoo can be efficiently created, consisting of the existing neural architectures manually designed by domain experts, identified by NAS, or transferred from models from different datasets. To make full use of the candidates in

the model zoo, in this paper, we propose a novel co-search framework, namely ‘‘HotNAS’’, to start searching from a hot state. In this way, compared with the cold-start co-search, HotNAS can reduce the search time from hundreds of GPU hours to **less than 3 GPU hours for ImageNet and 20 GPU minutes for CIFAR-10 without proxy**; while achieving accuracy comparable with the state-of-the-art models.

Figure 1 shows the results of co-search using a model zoo with 24 models on ImageNet dataset, targeting a system with 5ms on Xilinx ZCU 102 FPGA. From the top figure, there are only 4 models that can satisfy the timing constraint and the highest accuracy is 87.50%; however, within the range from 5ms to 10ms, there are a lot of good candidates with accuracy higher than 90%. The existing co-search frameworks ignore these candidates and search from scratch, leading to hundreds of GPU hours. Viewing from the opposite angle, HotNAS takes full use of these pre-trained models and customize the models that violate time constraints but have high accuracy to the target hardware. As such, HotNAS can avoid lengthy training procedure to generate the solution in a couple of hours, which is guaranteed to meet timing constraints while greatly improving accuracy to 91.47%, as shown in the bottom figure.

Seemingly straightforward, the architecture-hardware co-search from a hot start is not a simple matter: a fundamental challenge is the discovery of the best search space. Some of the prior co-search works [4]–[6], [15] consider hardware design space of loop tiling and loop order, and neural architecture search space with flexibility across the number of channels, filter size, and model quantization. However, we observe that one of the most efficient techniques, model pruning [16]–[20], has hitherto not been combined in the co-search. Integrating model pruning faces a lot of challenges: First, without the full consideration of hardware design, the model pruning can easily become useless since it introduces overheads. Second, one compression technique does not work for all performance bottlenecks. Finally, the model compression techniques are tightly coupled with hardware design and neural architecture search: As such, a difficult challenge is to simultaneously optimize all these spaces.

In HotNAS, we address the above challenges by collaboration among four sub-components: iSearch, iSpace, iDesign, and iDetect. First, iDesign provides hardware design support for different compression techniques. Second, following the observation that different pruning techniques work for different types of bottlenecks; iDetect is developed to identify performance bottleneck for each layer so that we can select the most suitable compression techniques to alleviate performance bottlenecks. According to the detected bottlenecks, iSpace creates a dedicated search space for each layer. Finally, iSearch is devised to jointly search the hardware, neural architecture, and model compression using specification from iSpace.

The main contributions of this paper are threefold:

- We propose a novel neural architecture search mechanism to search from a hot state (i.e., a pre-trained model), which allows us to reduce search time from 200 GPU Hours to 3 GPU Hours; meanwhile, the solution can improve the Top-1 and Top-5 accuracy on the ImageNet

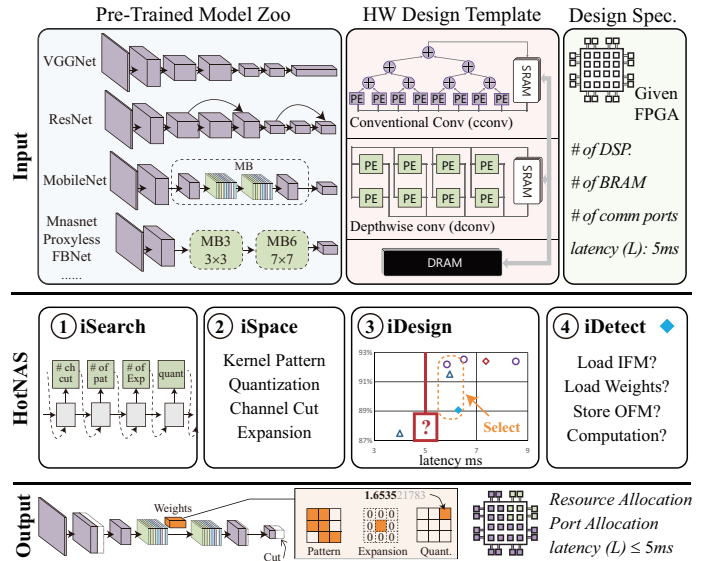


Figure 2. Implementation from the model zoo to hardware: (top) the given pre-trained model and design spec.; (middle) four components in the proposed HotNAS framework; (bottom) outputting the best neural architecture and hardware design.

dataset by 5.79% and 3.97%, respectively.

- An automated HotNAS framework is proposed to link the hardware design, neural architecture search, and model compression to automatically generate the architecture and hardware pair, such that the timing constraint can be met with the maximum model accuracy.
- In HotNAS, dedicated hardware designs to support the existing model compression are proposed, without which the model compression techniques may not achieve any performance gain at all.

The remainder of the paper is organized as follows. Section 2 presents design challenges and motivation. Section 3 presents the proposed HotNAS. Experimental results are shown in Section 4. Finally, concluding remarks are given in Section 5.

II. CHALLENGES AND MOTIVATION

This section demonstrates the challenges in the architecture-hardware co-search, and gives the motivation of this work.

Figure 2 demonstrates the architecture-hardware co-search problem, where we have a set of pre-trained models (called model zoo), the hardware design templates and the design specifications (e.g., constraints on the resource, area, and latency) as inputs. The co-search is to optimize neural architectures in the model zoo and hardware design to guarantee all design constraints to be met while maximizing accuracy, as shown in the bottom part in Figure 2

HotNAS framework is proposed in this paper to solve the above problem. As illustrated in the middle part of Figure 2, it is composed of four sub-components: ① iSearch, starting the co-search from hot instead of cold; ② iSpace, building an integrated search space which is in accordance with the performance bottleneck in the implementation; ③ iDesign, providing the design to support compression techniques on FPGAs; ④ iDetect, detecting the performance bottleneck to

Table I

SEARCH COST (GPU HOURS) OF NEURAL ARCHITECTURE SEARCH IS TOO HIGH. NOTE THAT THE HARDWARE-AWARE APPROACH NEEDS AN ENTIRE SEARCH FOR EACH SPECIFIC HARDWARE

NAS	NASNet	DARTS	MnasNet	FNAS	FBNet	ProxylessNAS
GPU Hours	48,000	90	40,000	267	216	200

guide the creation of iSpace. In the following text, we will show that there exists a couple of challenges in architecture-hardware co-search and all components work collaboratively to address these challenges.

Challenge 1: *How to efficiently explore neural architectures.*

The order of hundreds of GPU hours in architecture-hardware co-search cannot satisfy the short time-to-market requirements in many applications; as reported in Table I, the state-of-the-art hardware agnostic neural architecture search techniques (DARTS [21]) requires 90 GPU Hours, while the NAS for a specific type of hardware (MnasNet [3], FNAS [5], FBNet [1], and ProxylessNAS [2]) requires over 200 GPU Hours. Considering that the current computing system is composed of a large variety of hardware, the search process is simply unacceptable. In addition, the long search time leads to excessive CO_2 emission, which has already been known as a serious problem of existing NAS techniques [14].

We observe that the long search time in NAS framework is caused by the cold start. This leads to more than 40,000 GPU hours for MnasNet and NASNet to train a large number of potential architectures from scratch, and over 200 GPU hours when the hardware is considered. However, there exists a large set of pre-trained neural networks. We revisit the default configuration in the co-search framework: i.e. whether it is necessary to start the exploration from scratch, which results in low efficiency. In ① iSearch, we propose to make full use of the existing models and start the exploration from a hot state (e.g., pre-trained models).

Challenge 2: *Meet real-time constraint on specific hardware.*

Arbitrarily picking neural networks from the model zoo and plugging into the given hardware will lead to violations of the design specification, e.g., missing deadline in real-time systems. On the other hand, due to the large variety of hardware (different types of CPU, GPU, FPGA), it is infeasible to conduct co-search for all off-the-shelf hardware in advance.

Table II reports the hardware optimization results for the existing model zoo (Figure 1(a)) for Xilinx ZCU102 FPGA board with latency requirements of $\leq 5ms$. We can see that only four models can satisfy the timing constraints with the highest accuracy of 87.50%; while there are a group of networks whose accuracy is much higher, up to 92.54%, with the latency slightly exceeding the timing constraint.

The challenge here is, how can we compress the models to satisfy the timing constraints using its pre-trained weights, while a competitive model accuracy can be achieved. ② iSpace is developed to involve model compression in the search space, together with the hardware design space and neural architectures search space.

Table II

COMPARISON OF THE DESIGN, BOTTLENECKS, LATENCY, AND ACCURACY OF DIFFERENT NETWORKS (LATENCY $\leq 10ms$) ON XILINX ZCU102 BOARD USING 16-BIT FIX-POINT AT 200MHZ OF FREQUENCY. TARGETING 5ms FOR IMAGENET DATASET. (ORDERED BY LATENCY)

Network	Comp. Des.			Comm. Des.			Bottleneck				Lat. (ms)	Top-5 (%)
	Tm	Tn	Tm(d)	I_p	O_p	W_p	C	I	W	O		
alexnet	70	36	-	18	6	2	4	0	1	0	2.02	79.066
squeezenet 1.1	70	36	-	6	14	10	5	9	4	8	2.47	80.624
mnasnet 0.5	100	16	832	6	10	14	8	16	16	12	3.99	87.498
squeezenet 1.0	130	19	-	6	10	14	8	7	4	7	4.61	80.422
proxyless mobile	100	16	832	10	10	10	17	15	11	18	5.83	92.202
mnasnet1.0	100	18	704	10	10	10	7	21	10	14	5.94	91.51
resnet18	70	36	-	18	6	6	8	2	10	0	6.27	89.078
mobilenet v2	160	12	576	10	10	10	2	20	8	22	6.29	90.286
proxyless gpu	130	12	832	10	10	10	8	15	11	9	6.49	92.538
fbnet	100	16	832	10	10	10	11	29	10	15	7.37	92.386
proxyless cpu	220	8	704	10	10	10	18	5	10	28	8.53	92.394

Challenge 3: *One technique is not for all.*

One compression technique cannot solve all kinds of performance bottleneck in the accelerator. To effectively reduce latency, we need to specify where the bottleneck is and apply the most suitable compression technique to alleviate it.

Table II reports the performance bottleneck analysis. In the table, column ‘‘Comm. Des.’’ indicates the communication subsystem design, where I_p , O_p , and W_p are communication ports allocated for input feature maps (IFM), output feature maps (OFM) and Weights, respectively. The total communication bandwidth is limited, as a result, the performance of a layer may be bounded by the specific type of data movement. Using this information, the performance bottleneck can be computed. Under the column ‘‘Bottleneck’’, the above 4 types of bottlenecks are denoted by I, W, O, C. From the table, we can clearly observe that layers in one network lead to different types of bottlenecks.

In HotNAS, ③ iDesign will devise hardware design to support compression techniques, and provide the performance model; while ④ iDetect will match the compression technique with the corresponding type of performance bottleneck.

III. PROPOSED FRAMEWORK: HOTNAS

In response to all the challenges described in the previous section, we propose HotNAS framework in this section. As shown in Figure 2, HotNAS is composed of four sub-components, ① iSearch, ② iSpace, ③ iDesign, and ④ iDetect. This section will introduce these sub-components in detail.

① iSearch: Search from Hot Start

Figure 3 illustrates the overview of iSearch, which conducts the neural architecture search in two steps: (1) (top part of figure), it selects backbone architectures to be optimized; then, (2) (bottom part of figure), an optimizer tunes hyperparameters of neural architecture and hardware design simultaneously. The goal of iSearch is to find the architecture with the highest accuracy while meeting hardware design specifications. In the following texts, we will formally define the problem, and introduce the optimizer at the end of this section.

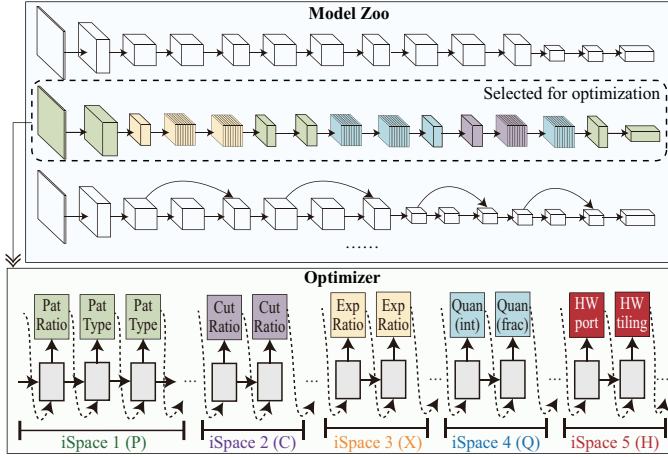


Figure 3. iSearch, iSpace, iDesign, and iDetect: (top) iSearch selects an architecture from the model zoo for optimization based on the results of iDesign and iDetect; (bottom) iSpace creates the search space in terms of performance bottlenecks, and iSearch determines hyperparameters for each search space.

i) Neural architectures and model zoo:

A neural architecture is defined as $A = \langle V, E, r, c, ch, o, f, para, acc \rangle$, composed of a set of nodes V representing the intermediate data (i.e., input and output feature maps), a set of edges $E \subseteq V \times V$ representing the dependency between a pair of nodes. For a node v_i in V , it has three hyperparameters r_i , c_i and ch_i representing the number of row, column and channel of v_i . For an edge $e_j \in E$, an operator o_j (e.g., convolution, depthwise convolution, or pooling, etc.) is associated to it. f_j represents the filter (i.e., weights) used in operator o_j , which is composed of a set of kernels. Each filter is associated with two hyperparameters: $s(f_i)$ indicates the size of the filter (e.g., 1×1 , 3×3 etc.), and $p(f_i)$ is a pattern applied to prune f_i . Both the size and the pattern of the filter are tunable, which will be introduced in ② iSpace. After all the above hyperparameters are determined and the neural architecture A is identified, it can be trained/fine-tuned on the training datasets (e.g., ImageNet) to obtain the parameters/weights $para(A)$, and finally we can obtain its test accuracy $acc(A)$ on the test dataset.

A pre-trained neural architecture is also called a model, and a model zoo $M = \{A_0, A_1, \dots, A_{N-1}\}$ is composed of N models. These models can be manually designed by experts, like AlexNet, VGGNet, ResNet, automatically searched via neural architecture search, like MnasNet, ProxylessNas, FBNet, or transferred from models for other datasets, like BiT [22]. In this work, we use the existing model zoo from torchvision, and collect the state-of-the-art pre-trained models from github; hence, the cost of building the model zoo can be neglected. Kindly note that, how to create the model zoo is out of the scope of this paper; related works can be found in [23], [24].

ii) FPGA and accelerator design:

The hardware efficiency is not only related to the neural architecture but also the resource on the FPGA and the accelerator design on it. An FPGA fp has 3 attributes: mem_{fp} ,

$comp_{fp}$, and BW_{fp} , referring to the on-chip memory size, the number of computing resources (e.g., DSPs), and the bandwidth between off-chip and on-chip memories, respectively.

The accelerator design should meet all resource constraints of a given FPGA. It is composed of two parts: the design of the computing subsystem and the design of the communication subsystem. As the basic operators o in architecture A are conducted in nested loops, the loop optimizations, in particular the loop tiling, are widely studied and used in the design of the computing subsystem in FPGAs [25], [26]. In addition, with the consideration of the large amount of data (i.e., intermediate data and weights), and the limited on-chip buffer in FPGA, it is infeasible to put all data on FPGA. Therefore, data are moved between the off-chip and on-chip memories. As such, the communication bandwidth for moving each type of data needs to be determined in the design phase.

As a whole, the accelerator design is defined as $D = \langle T_m, T_n, T_r, T_c, I_b, W_b, O_b \rangle$, containing the loop tiling design $\langle T_m, T_n, T_r, T_c \rangle$ and bandwidth allocation $\langle I_b, O_b, W_b \rangle$. Specifically, for an operator o_k associated to a pair of nodes $v_i \rightarrow v_j$ in an architecture, T_m, T_n, T_r, T_c are the tiling parameters on output feature maps (OFM) channels ch_j , input feature maps (IFM) channels ch_i , rows r_i , and columns c_i ; while $\langle I_b, O_b, W_b \rangle$ are the bandwidth allocated for moving IFM (i.e., v_i), OFM (i.e., v_j), and weights (i.e., f_k). For a design D and an architecture A , the latency of each operator, say o_k , can be determined with ③ iDesign tool. Then, the summation of all operators will be the latency of A , denoted as $lat(A)$.

iii) iSearch: two-step exploration

In iSearch, the first step is to select a set of candidate backbone architectures to be optimized. Given a neural architecture and an FPGA, it has already been well studied to obtain the best accelerator design D , as in [25]. Based on the design, HotNAS can generate the search space iSpace (Section III ②). iSearch will select models from the model zoo to be the backbone architecture, which will be the starting point of HotNAS, as shown in the top of Figure 3. The selection process is based on a Monte Carlo test, where we are given a timing constraint TC and the search space iSpace. We can prune the models whose minimum latency in the test fails to meet TC . The feasible architectures will be sorted in terms of a weighted reward (will be introduced in Formula 18) in terms of the minimal latency and original accuracy. Then, Top-K architectures will be selected as a starting point, where K is a user-defined variable.

Now, iSearch gets into the second step to conduct the neural architecture search based on these selected models to make them meet the given timing constraint with high accuracy. iSpace tool will provide search spaces for iSearch, including the filter patterning P , channel cutting C , quantization Q , filter expansion X , and hardware design H . All these search spaces are coupled with each other. In iSearch tool, we develop a reinforcement learning based optimizer to simultaneously explore all these spaces. Kindly note that other optimization techniques such as evolutionary algorithms [27] can be easily plugged into the iSearch tool. For better understanding, we will

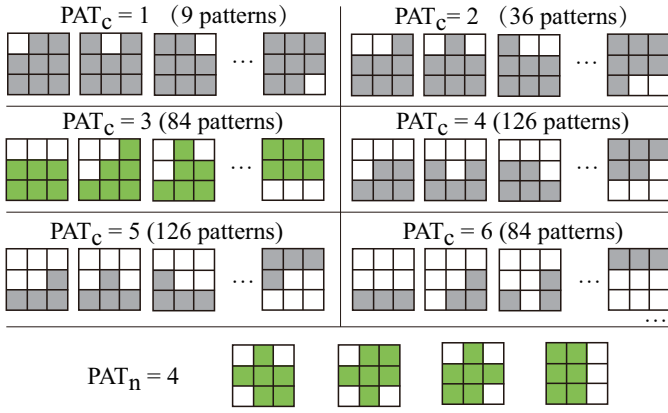


Figure 4. Patterns pruning for 3×3 filters with different pruning categories (PAT_c); each big square represent a pattern, and each blank square in a pattern indicates the weights in the corresponding position to be pruned. The last row shows the number of patterns (PAT_n) and the specific patterns to be selected for pruning filters.

present the details of the optimizer at the end of this section.

iv) Problem definition:

Based on all the above definitions, we formally define the architecture-hardware co-search optimization problem as follows: given a model zoo M , a specific FPGA FP , the accelerator design D_i of model A_i in M on FP , a target timing constraint T , and the baseline accuracy $acc_baseline$ we are going to determine:

- S : selection of architectures from zoo M , denoted as A_0 ;
- P, C, X, Q : tuning architecture hyperparameters of A_0 ;
- H : tuning hardware design hyperparameters on D_0 .

such that a new architecture A'_0 with competitive accuracy over $acc_baseline$ can be identified, while A'_0 under hardware design D'_0 can meet the timing constraint T .

② iSpace: An Integrated Search Space

iSpace links the compression technique with the neural architecture search and hardware design. In this work, we consider three model compression techniques: *i*) pattern pruning; *ii*) channel pruning; and *iii*) quantization. In the neural architecture search space, we consider *iv*) filter expansion; for hardware design space, we mainly consider *v*) communication bandwidth allocation and loop tiling, because FPGA accelerator design has typical templates which provides the above two kinds of hyperparameters in the design phase. Kindly note that, HotNAS is an open framework, with designers having the flexibility to modify or add new search parameters in terms of design needs. As an example, dataflow and loop order can be further integrated into hardware design space when it comes to ASIC design. However, this is out the scope of this work.

i) P: Pattern Pruning

The first search space is the pattern pruning space, which prunes the filter in the neural architecture A . A pattern is defined as a mask matrix $Mat[x][y]$; $Mat[x][y] = 0$ indicating that the weights at position $\langle x, y \rangle$ will be pruned, while $Mat[x][y] = 1$ indicates that the weights will remain. According to the number of 0 in $Mat[x][y]$, we can classify the pattern into different categories, and we use PAT_c to indicate the number of 0 in the pattern, as shown in Figure 4.

Among all patterns, one category will be selected for pruning. Each pattern category is further composed of many patterns; for instance, there are 84 potential patterns in the category of $PAT_c = 3$, as shown in Figure 4. For the hardware implementation, it simply cannot apply so many patterns as this will result in a large number of multiplexers in hardware implementation, making the design inefficient. In consequence, we will select a small number of patterns from the selected category, denoted as PAT_n . Figure 4 gives the example of the pattern pruning space for 3×3 filter, which selects the category of $PAT_c = 3$ and applies $PAT_n = 4$ patterns among 84 candidates.

The selected patterns will be applied for a set of filters. As demonstrated in [19], by applying the Euclidean norm, we can specify one pattern for each kernel in a filter, i.e., the determination of $p(f_i)$ (see the definition in ① iSearch). However, when implementing pattern pruning on hardware, the following two questions needing to be answered: (1) How many kernels in a filter will be pruned by each type of pattern. (2) Whether all layers need to be pruned or which layers will be pruned. For the first question, it is related to the tiling design parameters. In a tile, if multiple types of patterns are applied, it will break the execution pipeline and pattern pruning cannot improve performance at all. This will be shown in ③ iDesign. For the second question, applying patterns for the layers whose performance bottleneck is at communication, it will not help in improving performance but may reduce accuracy. Details will be illustrated in ④ iDetect.

ii) C: Channel Pruning

Unlike pattern pruning that changes the structure, the neural architecture will not be changed, with the channel pruning modifying the number of channels for a node $v_i \in V$ in architecture A . The left figure in Figure 5 shows the channel pruning, where CUT_n represents the number of channels to be cut off. We take $CUT_n = 2$ in this example. There are three consecutive nodes $v_i \rightarrow v_j \rightarrow v_k$, and we perform the channel pruning on v_j . In this figure, the grey channels in v_j indicate the ones to be cut off. A ripple effect is taken to both filters of $f_{i \rightarrow j}$ and $f_{j \rightarrow k}$. However, as the channel pruning may easily result in the accuracy drop since features are directly removed, we carefully formulate its search space for channel pruning only if the performance bottlenecks cannot be alleviated by other techniques (details in ④ iDetect).

iii) Q: Quantization

Quantization is another model compression technique. In general, the original model applies the data type of 32-bit floating-point, and we can convert it to the 16-bit fixed point without accuracy loss. Such a fixed point representation is composed of two parts, the integer and fraction parts represented by $\langle I, F \rangle$. For a given pre-trained architecture A , we can get the maximum and minimum parameters of one operator. Then, we can analyze the number of bits required by integer part I . Since the integer part is the most-significant bits, we will keep its bit-width, and further squeeze the fraction part F only, denoted as $Quan_f$ as shown in the right part of Figure 5. As will show in ④ iDetect, not all layers need to perform quantization, since it cannot alleviate specific types

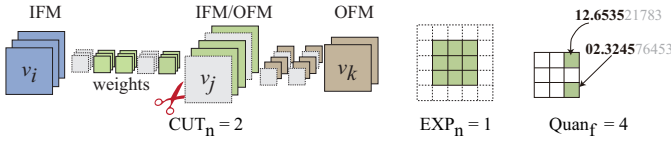


Figure 5. Three architecture search spaces: (left) channel cutting with ratio parameter CUT_r , (middle) kernel expansion with size parameter EXP_s , (right) weight quantization with fraction parameter $Quan_f$.

of performance bottlenecks.

iv) E: Filter Expansion

The previous three search spaces belong to model compression; while filter expansion belongs to neural architecture search space. This is motivated by the following two aspects: (1) many state-of-the-art neural architectures identified by NAS contains larger sized filters, and (2) for specific layers, the increase of filter sizes will not add latency overhead. This will be shown in ④ iDetect. We define EXP_n as the expansion factor on a filter, as shown in the middle part of Figure 5.

Furthermore, we have the following theorem to guarantee that the accuracy will not be reduced by expanding the kernel.

Theorem 1: Given a pre-trained model $A = \langle V, E, r, c, ch, o, f, para, acc \rangle$, for any operator o_i on edge e_i , the expansion on filter f_i with factor EXP_n will not decrease the accuracy, if the initial weights of the newly added weights on f_i are set to 0, and o_i is padded by EXP_n .

The proof of the above property is straightforward, since all computations remain the same when we increase the kernel size and padding with extra 0s. With the guarantee of no accuracy loss, the expanded kernel makes it possible to increase accuracy after a fine-tuned process.

v) H: Hardware Design Space

Finally, after the modifications to architectures, the original hardware design identified by the optimization algorithms may not be the optimal one. In iSpace, we also provide flexibility to modify the hardware design and build the hardware design space. In particular, according to the existing performance bottleneck, we create a search space to adjust bandwidth-related design hyperparameters: $\langle I_b, O_b, W_b \rangle$, and computation-related design hyperparameters, $\langle T_m, T_n, T_r, T_c \rangle$.

③ iDesign: Compression-Aware Performance Model

Figure 6 demonstrates the overview of system design, where the left-hand part is the off-chip memory to hold IFM, OFM, and weight; while the right-hand part is the on-chip accelerator design that implements both conventional convolution and depthwise convolution using on-chip computing resource (e.g., DSPs). In the accelerator design, say conventional convolution, a set of multiplication-and-accumulation are computed in parallel. Such a design has been using in many research works [25], [28], [29]; however, it still lacks a systematic model to efficiently support depthwise convolution and different compression techniques. In the following text, we will first overview the performance model of the conventional convolution [28], and then we revise the performance model to support depthwise convolution and compression.

First, we introduce the computing accelerator part. Let \mathbb{D} be the number of DSPs in the given FPGA, and K be the size

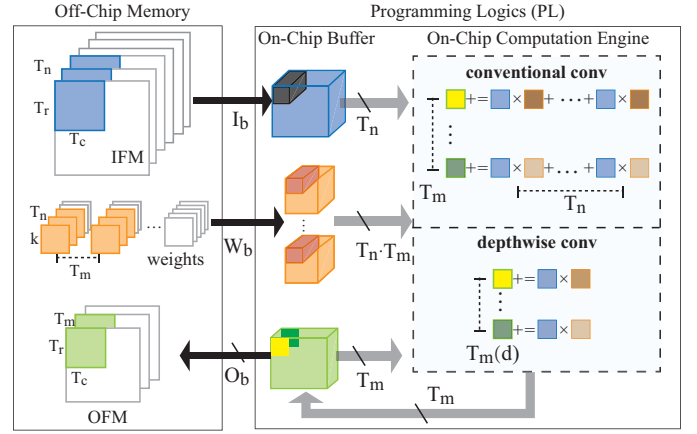


Figure 6. Illustration of accelerator architecture and design parameters: (left) off-chip memory to hold intermediate data and weights; (right) on-chip buffer and accelerator design.

of the filter. As shown in the right-hand part in Figure 6, the conventional convolution involves $T_m \times T_n$ multiplication and additions (MAC). For 16-bit data, each MAC needs one DSP. In addition, to consume all data in on-chip buffers, it needs to repeat $K \cdot K \cdot T_r \cdot T_c$ times for computation; and the pipeline initial interval (II) is optimized to 1 cycle. Therefore, we have the following constraints on computing resources and latency.

$$T_m \times T_n \leq \mathbb{D} \quad (1)$$

$$tComp = K \cdot K \cdot T_r \cdot T_c \times 1 \quad (2)$$

where $tComp$ is the latency of computation for all data provided by the on-chip buffer.

Second, the size of the on-chip buffer is limited by \mathbb{B} . There are three types of data in communication: IFM, OFM, and weights. We need to determine the on-chip buffer size for each type of data, denoted as bI , bO , bW , which can be easily obtained from the left part in Figure 6. Kindly note that the size of one on-chip buffer (BRAM) is limited, say 18K for ZCU102. For the dimension of data that needs to be accessed in parallel (e.g., channels of IFM, i.e., T_n), they need to be placed in different BRAMs. Hence, the amount of data without a parallel requirement (e.g., T_r and T_c in IFM) is divided by 18K. Finally, the size of the buffer is equal to 2 times tile size, where 2 indicates the double buffer utilized to hide communication by computation. We have the following constraints.

$$bI = 2 \times T_n \times \lceil T_r \cdot T_c \cdot bit_I / 18K \rceil \quad (3)$$

$$bO = 2 \times T_m \times \lceil T_r \cdot T_c \cdot bit_O / 18K \rceil \quad (4)$$

$$bW = 2 \times T_m \times T_n \times \lceil K \cdot K \cdot bit_W / 18K \rceil \quad (5)$$

$$bI + bO + bW \leq \mathbb{B} \quad (6)$$

where bit_I , bit_W , bit_O are the bit-width of the data type used for IFM, weights, and OFM.

Third, based on the buffer size and the bandwidth (I_b , W_b , O_b) allocated for each type of data buffer, we can get the communication latency (tI_{mem} , tW_{mem} , tO_{mem}) as follows.

$$tI_{mem} = \lceil T_n \cdot T_r \cdot T_c \cdot bit_I / I_b \rceil \quad (7)$$

$$tW_{mem} = \lceil T_m \cdot T_n \cdot K \cdot K \cdot bit_W / W_b \rceil \quad (8)$$

$$tO_{mem} = \lceil T_m \cdot T_r \cdot bit_O \cdot T_c / O_b \rceil \quad (9)$$

$$(I_b + W_b + O_b) \leq \mathbb{W} \quad (10)$$

where \mathbb{W} is the maximum bandwidth between off-chip memory and on-chip memory.

Finally, based on the above formulations, we can derive the latency model. Let M , N , R , C be the number of OFM channels, IFM channels, rows, and columns of the convolution layer. We have the following models.

$$Lat_1 = \max\{tComp, tI_{mem}, tW_{mem}\} \quad (11)$$

$$Lat_2 = \max\{\lceil \frac{N}{T_n} \rceil \cdot Lat_1, tO_{mem}\} \quad (12)$$

$$Lat = \lceil \frac{R}{T_r} \rceil \times \lceil \frac{C}{T_c} \rceil \times \lceil \frac{M}{T_m} \rceil \times Lat_2 + (tO_{mem} + Lat_1) \quad (13)$$

Since OFM is reused and stay in on-chip, it will be flushed to off-chip memory when IFM and weights are loaded for $\lceil \frac{N}{T_n} \rceil$ times. Lat_1 indicates the latency of computation, loading IFM, loading weights to be fired once, and Lat_2 indicates the latency of OFM to be flushed to off-chip memory. Finally, for one layer, OFM is flushed to off-chip memory for $B \times \lceil \frac{R}{T_r} \rceil \times \lceil \frac{C}{T_c} \rceil \times \lceil \frac{M}{T_m} \rceil$ times, and we have the total latency Lat .

For the model of depthwise convolution, we only need to modify T_m in the above formulas to be $T_m(d)$ and T_n to be 1. Kindly note that we consider the real-time scenario where the batch size is 1, and therefore, the communication subsystem (including on-chip buffer model Formulas 3 to 6, and off-chip memory access model Formulas 7 to 9) of two types of convolutions are shared. However, the accelerators are independent; therefore, we revise Formula 1 as follows.

$$T_m \times T_n + T_m(d) \leq \mathbb{D} \quad (14)$$

④ iDetect: Performance Bottleneck Detector

Based on the iDesign, we have several observations for the techniques introduced in ② iSpace, and we propose iDetect tool to analyze these search spaces in turn. Kindly note that all operators in a neural architecture will be implemented on one board and reuse these resources. Before discussing each search space, we first present the following corollary to detect the performance bottleneck of a layer based on iDesign.

Property 2: Given a layer and design parameters, we can detect the performance bottlenecks by considering Lat_1 and Lat_2 as follows:

- **O:** if Lat_2 is dominated by tO_{mem} , the performance bottleneck is on transmitting OFM data, otherwise,
- **I:** if Lat_1 is dominated by tI_{mem} , the performance bottleneck is on transmitting IFM data,
- **W:** if Lat_1 is dominated by tW_{mem} , the performance bottleneck is on transmitting weights,
- **C:** if Lat_1 is dominated by $tComp$, we have fully utilized the involved computation resource.

i) Pattern pruning can reduce computation time

Now, we are ready to answer the (1) question left in ② i): the number of kernels pruned by each type of pattern is

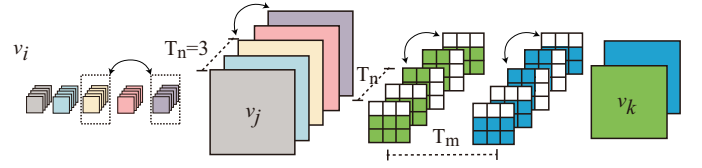


Figure 7. Reorder the input feature maps to make the pattern pruning take effects in reducing the computation latency.

coupled with the tiling factor T_m and T_n . As we can see from Figure 6, the data movement from on-chip weight buffer to the accelerator is conducted in a pixel-wise way. As a result, it requires $K \times K$ iterations to traverse the whole filter. To enable the effect of , we need to make sure that all patterns in one data tile are the same, as such we can skip these pruned weights in the outer loop to reduce the computation time. In this way, we can modify Formula 2 as follows.

$$tComp = (K \cdot K - PAT_n) \cdot T_r \cdot T_c \quad (15)$$

where PAT_n is the number of 0s in the pattern mask.

Next, since the pattern selection for kernels is based on the Euclidean norm, it cannot guarantee all patterns for same type of data tiles. We propose the input feature map reorder method to solve this problem. As shown in Figure 7, we can change the third and fifth channels for the filter used in the operator $o_{j,k}$. Correspondingly, we need to switch the feature map in node v_j . It will also affect the operator from v_i to v_j , where we need to switch the third and fourth filters. In this way, we can make the pattern pruning take effects and reduce the computation latency.

From Formulas 5 and 8, it may appear that pattern pruning can also reduce the on-chip buffer size and latency of loading weights. However, for buffer size, all layers reuse this buffer, and it cannot be specialized for one layer; while for loading weights, the pattern pruning will lead the loading procedure from sequential memory access to random access, as a result the latency may be even increased. Hence, we will keep the sequential memory access to guarantee performance.

Property 3: By applying the proposed reorder technique, pattern pruning can be employed to reduce the computing latency, but cannot reduce the latency of loading weights.

ii) Channel pruning can conditionally reduce latency

Channel pruning directly reduces the number of channels of feature maps in a node, and it can potentially reduce the latency. Let Cut_n be the number of channels cut on the feature maps of node v_i . When v_i acts as the input feature maps for an operator, we need to modify Formula 12 as follows:

$$Lat_2 = \max\{\lceil \frac{N - Cut_n}{T_n} \rceil \cdot Lat_1, tO_{mem}\} \quad (16)$$

Then, when v_i acts as the output feature maps for an operator, we revise Formula 13 as follows:

$$Lat = \lceil \frac{R}{T_r} \rceil \times \lceil \frac{C}{T_c} \rceil \times \lceil \frac{M - Cut_n}{T_m} \rceil \times Lat_2 + (tO_{mem} + Lat_1) \quad (17)$$

Property 4: Channel pruning can reduce the latency of a layer if and only if (1) $\lceil \frac{M - Cut_n}{T_m} \rceil \leq \lceil \frac{M}{T_m} \rceil$ or (2) $\lceil \frac{N - Cut_n}{T_n} \rceil < \lceil \frac{N}{T_n} \rceil$ and Lat_2 is not dominated by storing OFM data.

The above property indicates that pruning a small number of channels may not make an impact. As such, it guides the iSpace of channel pruning to take T_m or T_n as the step.

iii) *Quantization can reduce latency of loading weights*

Quantization is widely used in the neural network based FPGA implementations. It is demonstrated hybrid quantization can achieve good performance [30], where weights in different layers have different bit-widths. When we adopt such a hybrid approach, what benefits can be achieved? From Formula 8, we can see that the quantization can take effects in reducing the latency of loading weights. This can be implemented by composing multiple weights into one package. As with computing latency, since the initial interval is already optimized to 1 cycle as shown in Formula 2, the lower bit-width operations cannot further reduce clock cycles. Lower bit-width can reduce the number of computing resources and have the potential to achieve high clock frequency. However, when we consider an end-to-end implementation, the computing engine is shared by all layers. Therefore, the layer with the largest bit-width dominates the design performance.

Property 5: Quantization on a single layer can reduce latency of loading weights, but it may not reduce the computation latency if there exists another layer with larger bit-width.

⑤ *Optimizer: Search space exploration*

Finally, we introduce the RNN-based reinforcement learning optimizer employed in iSearch. As shown in the bottom part of Figure 3, an RNN controller is designed based on the created design space by the ② iSpace tool. Specifically, the controller is composed of a softmax classifier to predict hyperparameters for each search space in iSpace (e.g., quantization $Quan_f$ for a layer). The predicted hyperparameters will identify a specific neural network and hardware design, which can derive a reward in terms of accuracy and latency. The search process will optimize the controller by tuning its parameters θ_c to maximize the expectation of the reward. A policy gradient method will be employed to update parameters θ_c , aiming to predict better architectures over a series of episodes.

In each episode, the predicted hyperparameters can be regarded as actions. Based on the actions, the optimized neural architecture A and hardware design D can be derived. In order to update the controller for the next episode, we need to compute the reward according to the following procedures: (1) calculate latency lat of architecture A on design D by using the performance models proposed in ③ iDesign and ④ iDetect; (2) verify whether timing constraint T can be satisfied; if $lat > T$, we will directly calculate the reward without fine-tuning the model, otherwise, the reward is calculated based on accuracy and latency in the next step; (3) fine-tune architecture A to obtain accuracy acc on a hold-out dataset; since the model is pre-trained, we do not need to train the model from scratch; instead, we fine-tune the model for a small number of data batches (not epochs), say $\beta = 10$, to obtain acc . Finally, the calculation of reward is based on the following formula:

$$R(acc, lat) = \alpha \times r_{acc} + (1 - \alpha) \times r_{lat} \quad (18)$$

where α is a scaling parameter to control with the search is for higher accuracy (i.e., larger α) or lower latency (i.e., smaller

Table III
MODEL SELECTION BASED ON iDESIGN AND iDETECT TO CONDUCT 100 MONTE CARLO TEST IN iSPACE, WITH THE TIMING CONSTRAINT OF $T \leq 5ms$.

Models	Original Latency (ms)	Latency (ms) in Monte Carlo Tests			
		Min	Sat.	Max	Ave.
ProxylessNAS	5.83	4.47	✓	5.95	5.08
MnasNet	5.94	4.68	✓	7.88	5.21
ResNet	6.27	4.31	✓	5.9	5.05
MobileNet	6.29	5.07	×	7.88	6.2
FBNet	7.37	5.85	×	7.82	6.66

α). If $lat > T$ indicating that the timing constraint cannot be satisfied, we have $r_{acc} = -1$ and $r_{lat} = T - lat$; otherwise, we normalize r_{acc} and r_{lat} to the range from -1 to 1, as follows: $r_{acc} = \frac{acc - A_{min}}{A_{ori} - A_{min}} \times 2 - 1$ and $r_{lat} = \frac{T - lat}{T - T_{min}} \times 2 - 1$, where A_{ori} is the original accuracy of backbone architecture; T is the timing constraint; A_{min} and T_{min} are the lower bounds on accuracy and latency, which are involved for a better normalization.

Based on the reward function, the optimizer will iteratively work in two steps. First, the controller predicts a sample, and gets its reward R . Then, the Monte Carlo policy gradient algorithm [31] is employed to update the controller:

$$\nabla J(\theta) = \frac{1}{m} \sum_{k=1}^m \sum_{t=1}^T \gamma^{T-t} \nabla_{\theta} \log \pi_{\theta}(a_t | a_{(t-1):1}) (R_k - b) \quad (19)$$

where m is the batch size and T is the number of steps in each episode. Rewards are discounted at every step by an exponential factor γ and baseline b is the average exponential moving of rewards.

IV. EXPERIMENTAL RESULTS

The proposed HotNAS is evaluated on commonly used datasets, ImageNet [32] and CIFAR-10 with Xilinx ZCU102 board. In the following texts, we will first introduce the experimental setup. Then, we will compare HotNAS with the state-of-the-art models to show that HotNAS can achieve up to 5.79% top-1 accuracy gain with the same timing constraint. Next, we will visualize the results explored by HotNAS, followed by the design space exploration results to demonstrate the importance of co-exploring all design spaces in iSpace. Finally, we report results and detailed analysis on CIFAR-10, showing that HotNAS can achieve consistent improvement on different datasets.

A. Experimental setup

Model Zoo. For ImageNet dataset, we apply all models in torchvision, including AlexNet, VGGNet, ResNet, MobileNet-v2, Mnasnet, etc., as shown in Figure 1. We also include the FBNet [1] and ProxylessNAS [2] for comparison. In the experiments, we select a set of models to be optimized. According to iDesign and iDetect, we run Monte Carlo Tests to get statistic latency for 100 solutions in iSpace, as shown in Table III. We prune the models whose minimum latency cannot satisfy the timing constraints, say $T \leq 5$ in our settings.

Table IV
ON IMAGENET, COMPARISON OF THE STATE-OF-THE-ART NEURAL ARCHITECTURES WITH TIMING CONSTRAINTS OF 5ms.

Model	Type	Latency	Sat.	Param. (#)	Param. (S)	Top-1	Top-5	Top-1 Imp.	Top-5 Imp.	GPU Time
AlexNet	manually	2.02	✓	61.1M	122.20MB	56.52%	79.07%	-	-	-
<i>MnasNet 0.5 *</i>	<i>auto</i>	3.99	✓	2.22M	4.44MB	67.60%	87.50%	-	-	40,000H
SqueezeNet 1.0	manually	4.76	✓	1.25M	2.50MB	58.09%	80.42%	-	-	-
ProxylessNAS	auto	5.83	×	4.08M	8.16MB	74.59%	92.20%	-	-	200H
MnasNet	auto	5.94	×	4.38M	8.77MB	73.46%	91.51%	-	-	40,000H
Resnet	manually	6.27	×	11.69M	23.38MB	69.76%	89.08%	-	-	-
Co-Exploration [6]	auto	-	-	-	-	70.42%	90.53%	-	-	266H
HotNAS-Resnet(4ms)	auto	4.00	✓	10.99M	17.49MB	68.27%	88.21%	0.67%	0.71%	2H22M
HotNAS-Resnet	auto	4.22	✓	11.19M	17.90MB	69.14%	88.83%	1.54%	1.33%	2H01M
HotNAS-ProxylessNAS	auto	4.86	✓	4.38M	8.31MB	73.39%	91.47%	5.79%	3.97%	2H37M
HotNAS-Mnasnet	auto	4.99	✓	4.07M	6.56MB	73.24%	91.37%	5.64%	3.87%	1H50M

“*”: baseline; “*auto & manually*”: the model identified by NAS or human experts; “× & ✓”: violate or meet timing constraints.

Results of Co-Exploration is derived from [6], since the pre-trained model is not provided; The latency is not reported since it uses different hardware.

Kindly note that the maximum latency can be larger than the original model latency, because we change the hardware configuration during the search, which may reduce bandwidth for the data movement which is the performance bottleneck. Based on the results in Table III, we select ProxylessNAS (mobile), MnasNet 1.0 (depth multiplier of 1.0), and Resnet-18 (with 18 layers) [33] for optimization. We denote them as ProxylessNAS, Mnasnet, and Resnet, respectively.

For CIFAR-10 dataset, we collect the 4 sets of pre-trained models, including ResNet-18 [33], DenseNet-121 [34], MobileNet-v2 [35], and BiT [22], among which BiT achieves the state-of-the-art accuracy on CIFAR-10 dataset, which is built on top of existing neural networks. In our experiments, with the hardware performance consideration, we select the ResNet-50 based version for BiT, which provides a baseline with the accuracy of 97.07% and latency of 6.88ms. For a better presentation, we denote the above models as ResNet, DenseNet, MobileNet, and BiTNet, respectively.

iSearch. In iSearch component, we first need to determine parameters α and β . We set $\alpha = 0.7$ to generate the reward as shown in Formula 18, and set the number of batch size $\beta = 10$ to be used in the fine-tune phase. Furthermore, we will investigate the effects on performance made by different configurations of α and β on CIFAR-10. Second, we need to set the number of episodes for reinforcement learning; here, we set the maximum episode to be 2,000 which can guarantee the convergence of the controller. After running iSearch, we can obtain a set of architectures, and we will select the best architectures, i.e., the architecture with the highest accuracy under the given timing constraints.

iSpace. A new module that supports pattern pruning, channel pruning, filter expansion, and quantization is implemented in Pytorch by overriding the existing Conv2d module. During the iSearch process, the module can be customized for each layer in terms of the searched parameters, and automatically integrated into the model with the original weights.

iDesign. We apply Xilinx ZCU102 board with XCZU9EG

chip as the implementation hardware, which is composed of 600K logic cells, 32.1Mb on-chip buffers, 2,520 DSPs. For data movement between on-chip and off-chip memory, there are 4 HP ports with the bandwidth of 128 bits for each.

B. Results on ImageNet

i) Comparison with HotNAS

Table IV reports the comparison results of HotNAS with the existing state-of-the-art models. In the table, the column “Type” shows whether the model is identified by NAS or manually designed. The column “Sat.” shows whether the model satisfies the timing constraint of 5ms. Columns “Param. (#)” and “Param. (S)” reports the number of parameters and the size of parameters, respectively. Columns “Top-1”, “Top5”, “Top-1 Imp.”, and “Top-5 Imp.” are model accuracy and accuracy gain to the baseline model on ImageNet. Column “GPU Hours” shows the cost to identify the model for all models identified by NAS. Finally, the rows marked as bold are models identified by the proposed HotNAS.

From the results in Table IV, we have three important observations: (1) Directly plugging the existing models onto the target FPGA board will easily result in the latency to be violated; while the proposed HotNAS can guarantee to find the architectures to meet the latency constraints, meanwhile achieving high accuracy. (2) For the existing models that can directly satisfy the timing constraints, the highest top-1 accuracy and top-5 accuracy are merely 67.60% and 87.50%. In comparison, HotNAS can achieve 5.79% and 3.97% accuracy gain with 73.39% for top-1 and 91.47% for top-5. (3) The cost of the existing neural architecture search is extremely high, which is at least 200 GPU hours. In comparison, the proposed HotNAS only takes less than 3 GPU hours to identify the model. Furthermore, compared with the existing co-exploration method, the search time can be significantly reduced from 266 GPU hours with 2.97% Top-1 accuracy gain. All these observations clearly demonstrate the

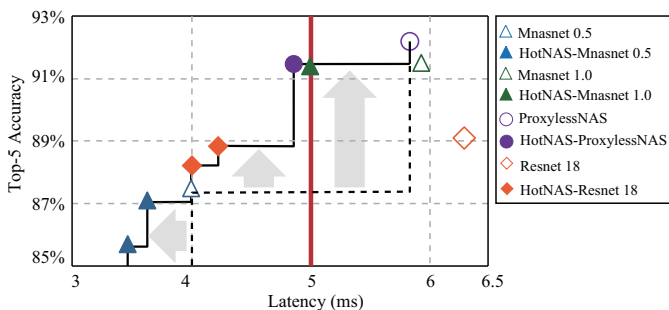


Figure 8. Pushing forward the Pareto frontier between latency and top-5 accuracy from the one built by the existing models to that by HotNAS under the timing constraint of 5ms.

superiority of HotNAS to obtain solutions with high accuracy and low search cost.

Besides, from the results, we can see that HotNAS performs good at reducing the latency while maintaining high accuracy. For Resnet18, HotNAS can reduce the latency from 6.27ms to 4.22ms with 32.70% reduction, while the top-5 accuracy loss is merely 0.25%; for ProxylessNAS, the latency reduction is 16.64% with only 0.53% top-5 accuracy loss; for Resnet18, these figures are 15.9% and 0.14%. We will have a detailed and visualized analysis in the latency reduction later in this section. A further observation from the above result is that the manually designed Resnet18 can achieve larger reductions in latency than the automatically identified ones. This is reasonable since the automatically designed architectures have already been used for optimizing for other platforms, while manually designed architectures may have more redundant parameters. This can also be observed by the reduction in both the number of parameters and the size of parameters.




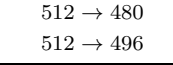

Figure 8 further shows the comparison of Pareto frontiers built by the existing models and HotNAS. In this figure, the x-axis and y-axis represent the latency and accuracy, respectively. The red line stands for the timing constraints. The solid points are solutions identified by HotNAS, while the hollow ones are the existing models. The arrows in this figure clearly demonstrate that HotNAS can significantly push forward the Pareto frontier between accuracy and latency in two directions: (1) vertical direction: improving accuracy; (2) horizontal direction: reducing latency. The results in this figure again demonstrate the efficiency and effectiveness of the proposed HotNAS.

ii) Results visualization

Table V shows the visualization results of HotNAS-Resnet18. For other resultant architectures, they have similar results, but the model is too large to demonstrate¹. In this table, column “iDetect” shows the performance bottleneck with the original design detected by HotNAS, and column “iSpace” shows the built search spaces for these the corresponding layers. The column “exploration results” show the detailed changes from the original architecture to the resultant model. Finally, the column “Red.” shows the latency reduction contributed by each search space.

It is clearly shown in this table that the proposed HotNAS can identify different types of performance bottleneck in

Table V
SOLUTION VISUALIZATION, USING THE MODIFIED LAYERS AND LATENCY REDUCTION IN HOTNAS-RESNET18 AS AN EXAMPLE

Layers/HW	iDetect	iSpace	Exploration Results	Red. (ms)
layer1[0].conv1			PATr=3, PATn=4	
layer1[0].conv2				
layer1[1].conv1				
layer1[1].conv2	C	Pattern		0.57
layer2[0].conv2				
layer2[1].conv1				
layer2[1].conv2				
layer4[0].conv1			512 → 480	
layer4[1].conv1	I	Channel	512 → 496	0.15
layer4[0].conv1	-			
layer4[1].conv1	-			
layer4[0].conv2		Quant.	[1, 15] → [1, 7]	1.01
layer4[1].conv2	W			
I_b			18 → 20	
W_b	-	HW	6 → 5	0.32
Total				2.05

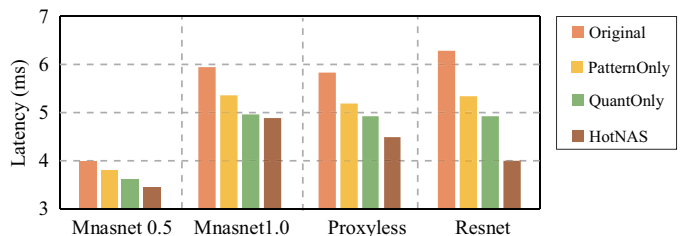


Figure 9. Comparison results in the latency reduction among three techniques using the original architecture as baseline: (1) PatternOnly: apply pattern pruning [19], (2) QuantOnly: apply hybrid quantization [30], (3) HotNAS.

the architecture, and apply the matched techniques to alleviate the performance bottlenecks. Specifically, pattern pruning identifies 4 patterns in pattern category $PAT_r = 3$, and achieves 0.57ms latency reduction. Channel pruning, quantization, and hardware modifications achieve a reduction of 0.15ms, 1.01ms, and 0.32ms, respectively. As a whole, the reduction is 2.05ms, from 6.27ms to 4.22ms, as shown in Table IV. Kindly note that since the latency of loading IFM and loading weights are quite close for layer 4, iSpace creates search spaces for both channel pruning and quantization.

iii) No space in iSpace can be dispensed

There are a lot of existing techniques that focus on devising a specific technique for model compression. We compare with the two most effective methods using pattern pruning only [19], denoted by PatternOnly; and hybrid quantization [19], denoted by QuantOnly. However, as discussed in this iDetect, no technique can cover all kinds of performance bottlenecks. Results in Figure 9 verify this claim. Kindly note that the hardware space is kept for all techniques for a fair comparison. In this figure, the x-axis is the backbone architecture, and the y-axis is the latency that can be achieved with the same accuracy constraint. The baseline is the original neural architecture without optimization.

Results in Figure 9 clearly demonstrate that without fully

¹This project will be open-source after the blind review, and all models with fine-tuned parameters can be accessed online.

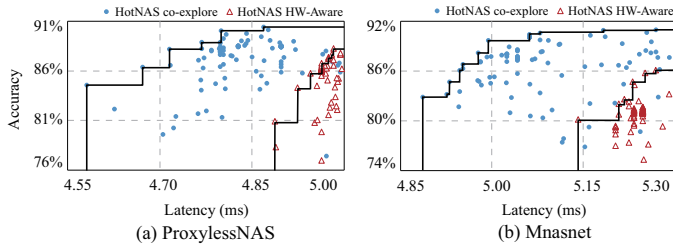


Figure 10. Illustration of importance of Co-Exploration.

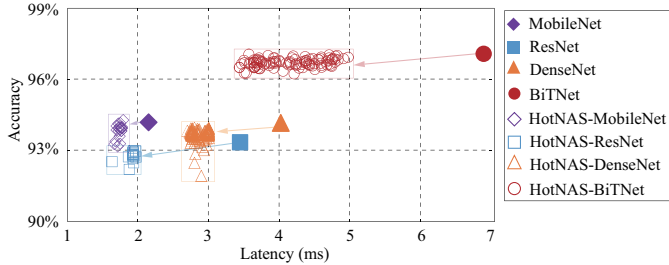


Figure 11. On CIFAR-10, exploration results of HotNAS for 4 pre-trained models with different target latency constraints.

considering the performance bottleneck and apply only one technique for optimization will lead to inferior solutions. Taking Resnet as an example, PatternOnly can reduce the time from 6.27ms to 5.34ms, and QuantOnly can further reduce it to 4.92ms. By a full consideration of all kinds of bottlenecks, HotNAS can achieve the architecture with 4ms, which achieves 25.09% and 18.71% latency reductions compared with PatternOnly and QuantOnly, respectively. Results from this group of experiments emphasize the needs of an automatic tool to analyze the model, detect the performance bottlenecks, and alleviate each kind of bottlenecks using the correct technique.

In addition to the compression techniques, the hardware exploration is also important, this is shown from the results in Figure 10. We compare the exploration results of HotNAS with the consideration of hardware modification space (“HotNAS w/ HW”) and results of that without hardware modification (“HotNAS w/o HW”). Figures 10(a) and 10(b) demonstrate the results for ProxylessNAS and Mnasnet, respectively. In these figures, the x-axis and y-axis are the latency and top-5 accuracy, respectively. Each blue dot represents a solution explored by “HotNAS w/ HW”, while each red triangle represents that explored by “HotNAS w/o HW”. The solutions with either higher accuracy or lower latency among each group of results form the Pareto frontiers.

As we can see from Figure 10, the proposed HotNAS can significantly push forward the Pareto frontiers. Specifically, for ProxylessNAS, the smallest latency of solutions explored by “HotNAS w/ HW” is 4.58ms, while that of “HotNAS w/o HW” is 4.89ms; for accuracy, “HotNAS w/ HW” can achieve 2.15% accuracy gain with latency reduction of 0.11ms. These results emphasize the importance of conducting co-exploration in neural architecture search.

C. Results on CIFAR-10

i). Pushing forward accuracy-latency Pareto frontier

Table VI
ON CIFAR-10, COMPARISON OF THE BASELINE MODELS AND THE BEST SOLUTIONS EXPLORED BY HOTNAS AFTER FINE-TUNING

Model	Accuracy			Latency (ms)		
	baseline	HotNAS	comp.	baseline	HotNAS	impr.
ResNet	93.33%	93.36%	+0.03%	3.44	1.93	43.90%
DenseNet	94.14%	94.19%	+0.05%	4.01	2.87	28.55%
MobileNet	94.17%	94.27%	+0.10%	2.14	1.79	16.74%
BiTNet	97.07%	97.13%	+0.06%	6.88	3.56	48.26%

HotNAS can consistently push forward the accuracy-latency Pareto frontier for different datasets. On CIFAR-10 dataset, we achieve similar results as ImageNet dataset. Figure 11 reports the exploration results for four pre-trained models. In this figure, x-axis and y-axis stand for latency and accuracy, respectively. The solid shapes represent the performance baseline models, while others represent the explored results. In order to evaluate the scalability of proposed HotNAS, unlike the previous experiment applying a uniform target latency, we set an individual target latency for each model according to the baseline latency. Specifically, the target latency constraints are 2ms, 3ms, 1.8ms, and 5 ms for ResNet, DenseNet, MobileNet, BiTNet, respectively.

Results in this figure show that all architectures identified by HotNAS can satisfy latency constraint, while achieving similar accuracy with the baseline architectures. As a result, the accuracy-latency Pareto frontier can be significantly pushed forward, likewise that for ImageNet. Another interesting observation is that BiTNet has a loose timing constraint, and HotNAS can find a wider range of results with high accuracy.

Table VI reports the detailed comparison of the best architectures identified by HotNAS over the baseline model. The best architecture is selected based on the architectures with the highest accuracy while satisfying the timing constraint. Then, we fine-tune the selected architecture for 10 epochs to obtain the final accuracy. The accuracy and latency for the original model and the one identified by HotNAS are reported in Columns “baseline” and “HotNAS” under Columns “Accuracy” and “Latency”.

From Table VI, it is clear to see that HotNAS can efficiently reduce the latency which achieving accuracy gain on CIFAR-10 dataset. Specifically, for ResNet, HotNAS identifies the solution with 43.90% latency reduction and 0.03% accuracy gain; these figures are 28.55% and 0.05% for DenseNet; 16.74% and 0.10% for MobileNet; 48.26% and 0.06% for BitNet. The above results demonstrate the efficiency and effectiveness of HotNAS.

ii). Exploration with different configurations

There are two hyperparameters in the RNN-based optimizer: β for the batch size of fine-tuning in the search process; α for the weights in the reward formulation. In the following, we will test different settings on both.

First, we apply two settings on β : (1) $\beta = 195$ for “1-epoch-search” which will fine-tune the identified architecture using the whole training set; (2) $\beta = 10$ for “fast-search” which only use a portion of dataset as in ImageNet experiments; Table VII reports the results. We can see that fast-search can achieve

Table VII
ON CIFAR-10, COMPARISON OF DIFFERENT SETTING IN HOTNAS ON FINE-TUNE BATCH SIZE β DURING THE SEARCH PROCESS; $\beta = 195$ FOR 1-EPOCH-SEARCH AND $\beta = 10$ FOR FAST-SEARCH

Model	1-epoch-search			fast-search		
	Accuracy	Latency	GPU Time	Accuracy	Latency	GPU Time
ResNet	93.36%	1.93	7M21S	92.74%	1.84	3M26S
DenseNet	94.08%	2.79	55M26S	94.19%	2.87	12M04S
MobileNet	94.27%	1.79	20M15S	94.21%	1.79	4M26S
BiTNet	97.13%	3.56	2H20M	97.04%	3.84	18M44S

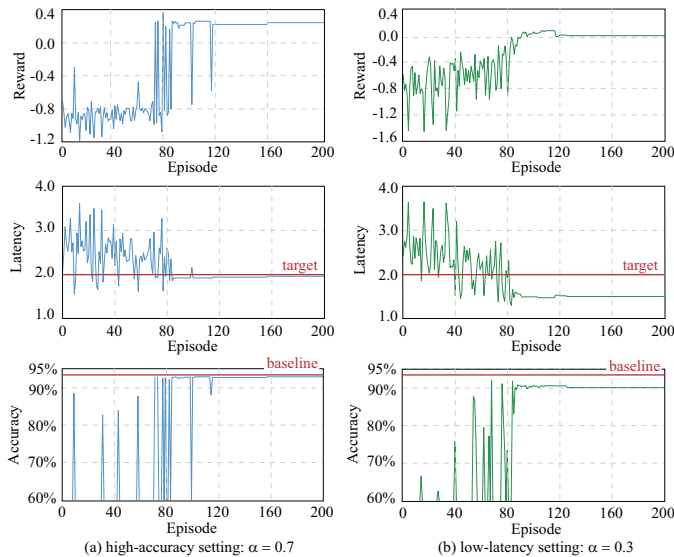


Figure 12. On CIFAR-10, comparison of different settings on scaling parameters α in optimizing ResNet.

competitive accuracy compare with 1-epoch-search; in particular, for DenseNet, fast-search achieves 0.11% higher accuracy. In addition, for all models, fast-search can find solutions in 20 minutes. These results demonstrate the efficiency of HotNAS.

Figure 12 reports the results of two settings on α : (1) $\alpha = 0.7$ aiming at higher accuracy as shown in Figure 12(a), (2) $\alpha = 0.3$ aiming at lower latency as shown in Figure 12(b). We report the comparison results of ResNet on CIFAR-10, and we observe similar results for the other models and datasets. In the experiments, we trace the reward, latency, and accuracy of solutions in each episode, which are reported in Figure 12(a) and (b), respectively. For latency, the red line shows the target latency, i.e., $\leq 2.0ms$; while for accuracy, the red line shows the baseline accuracy, which is 93.33%.

From Figure 12, we can see that the search processes are converged after 160 and 120 episodes for the high-accuracy setting and the low-latency setting, respectively. At the convergence, the latency of solutions identified by low-latency setting is lower than the high-accuracy setting; more interesting, the high-accuracy setting finds solutions with latency near to the threshold $2ms$. For accuracy, we can see that high-accuracy setting finds solutions with higher accuracy, which is almost the same with the baseline accuracy. As shown in the results in Table VI, after a fine-tuned process, the accuracy can even higher than the baseline. One more thing noted by the accuracy results is that there are several episodes having no

accuracy. This is because the latency cannot be satisfied, and we terminate the training procedure to accelerate the search process.

All the above results show that HotNAS provides flexibility for designers to better optimize neural architecture and hardware design according to their varied demands.

V. CONCLUSION

In this work, we identify the last mile problem in current neural architecture search and hardware accelerator design and propose the HotNAS toolset to solve the problem. Instead of search architectures from scratch, we propose to stand on the shoulders of the existing models to conduct an incremental hardware-aware neural architecture search. In HotNAS, four components work collaboratively to (1) identify the hardware performance bottleneck by iDetect, (2) build search spaces iSpace in terms of results from iDetect, (3) co-design the neural architecture and hardware accelerator by iSearch with the performance model provided by iDesign. Experimental results on ImageNet dataset demonstrate that HotNAS can guarantee the resultant system to meet timing specifications, while achieving over 5.6% top-1 and over 3.8% top-5 accuracy gain, compared with the state-of-the-art models.

REFERENCES

- [1] B. Wu, X. Dai, P. Zhang, Y. Wang, F. Sun, Y. Wu, Y. Tian, P. Vajda, Y. Jia, and K. Keutzer, "Fbnet: Hardware-aware efficient convnet design via differentiable neural architecture search," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 10734–10742, 2019.
- [2] H. Cai, L. Zhu, and S. Han, "Proxylessnas: Direct neural architecture search on target task and hardware," *arXiv preprint arXiv:1812.00332*, 2018.
- [3] M. Tan, B. Chen, R. Pang, V. Vasudevan, M. Sandler, A. Howard, and Q. V. Le, "Mnasnet: Platform-aware neural architecture search for mobile," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2820–2828, 2019.
- [4] C. Hao, X. Zhang, Y. Li, S. Huang, J. Xiong, K. Rupnow, W.-m. Hwu, and D. Chen, "Fpga/dnn co-design: An efficient design methodology for iot intelligence on the edge," in *2019 56th ACM/IEEE Design Automation Conference (DAC)*, pp. 1–6, IEEE, 2019.
- [5] W. Jiang, X. Zhang, E. H.-M. 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*, pp. 1–6, 2019.
- [6] W. Jiang, L. Yang, E. H.-M. Sha, Q. Zhuge, S. Gu, S. Dasgupta, Y. Shi, and J. Hu, "Hardware/software co-exploration of neural architectures," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 2020.
- [7] 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 Transactions on Computers*, 2020.
- [8] L. Yang, Z. Yan, M. Li, H. Kwon, L. Lai, T. Krishna, V. Chandra, W. Jiang, and Y. Shi, "Co-exploration of neural architectures and heterogeneous asic accelerator designs targeting multiple tasks," *arXiv preprint arXiv:2002.04116*, 2020.
- [9] L. Yang, W. Jiang, W. Liu, H. Edwin, Y. Shi, and J. Hu, "Co-exploring neural architecture and network-on-chip design for real-time artificial intelligence," in *2020 25th Asia and South Pacific Design Automation Conference (ASP-DAC)*, pp. 85–90, IEEE, 2020.
- [10] X. Zhang, C. Hao, Y. Li, Y. Chen, J. Xiong, W.-m. Hwu, and D. Chen, "A bi-directional co-design approach to enable deep learning on iot devices," *arXiv preprint arXiv:1905.08369*, 2019.
- [11] Y. Li, C. Hao, X. Zhang, X. Liu, Y. Chen, J. Xiong, W.-m. Hwu, and D. Chen, "Edd: Efficient differentiable dnn architecture and implementation co-search for embedded ai solutions," *arXiv preprint arXiv:2005.02563*, 2020.

- [12] S. Bian, W. Jiang, Q. Lu, Y. Shi, and T. Sato, "Nass: Optimizing secure inference via neural architecture search," *arXiv preprint arXiv:2001.11854*, 2020.
- [13] W. Jiang, E. H.-M. Sha, Q. Zhuge, and X. Chen, "Optimal functional-unit assignment and buffer placement for probabilistic pipelines," in *Proceedings of the Eleventh IEEE/ACM/IFIP International Conference on Hardware/Software Codesign and System Synthesis*, pp. 1–10, 2016.
- [14] E. Strubell, A. Ganesh, and A. McCallum, "Energy and policy considerations for deep learning in nlp," *arXiv preprint arXiv:1906.02243*, 2019.
- [15] Q. Lu, W. Jiang, X. Xu, Y. Shi, and J. Hu, "On neural architecture search for resource-constrained hardware platforms," *arXiv preprint arXiv:1911.00105*, 2019.
- [16] S. Han, H. Mao, and W. J. Dally, "Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding," *arXiv preprint arXiv:1510.00149*, 2015.
- [17] H. Mao, S. Han, J. Pool, W. Li, X. Liu, Y. Wang, and W. J. Dally, "Exploring the regularity of sparse structure in convolutional neural networks," *arXiv preprint arXiv:1705.08922*, 2017.
- [18] W. Wen, C. Wu, Y. Wang, Y. Chen, and H. Li, "Learning structured sparsity in deep neural networks," in *Advances in neural information processing systems*, pp. 2074–2082, 2016.
- [19] X. Ma, F.-M. Guo, W. Niu, X. Lin, J. Tang, K. Ma, B. Ren, and Y. Wang, "Pconv: The missing but desirable sparsity in dnn weight pruning for real-time execution on mobile devices," in *AAAI*, pp. 5117–5124, 2020.
- [20] W. Yawen, W. Zhepeng, J. Zhenge, S. Yiyu, and H. Jingtong, "Intermittent inference with nonuniformly compressed multi-exit neural network for energy harvesting powered devices," *arXiv preprint arXiv:2004.11293*, 2020.
- [21] H. Liu, K. Simonyan, and Y. Yang, "Darts: Differentiable architecture search," *arXiv preprint arXiv:1806.09055*, 2018.
- [22] A. Kolesnikov, L. Beyer, X. Zhai, J. Puigcerver, J. Yung, S. Gelly, and N. Houlsby, "Big transfer (bit): General visual representation learning," *arXiv preprint arXiv:1912.11370*, 2019.
- [23] C. Ying, A. Klein, E. Christiansen, E. Real, K. Murphy, and F. Hutter, "Nas-bench-101: Towards reproducible neural architecture search," in *International Conference on Machine Learning*, pp. 7105–7114, 2019.
- [24] X. Dong and Y. Yang, "Nas-bench-102: Extending the scope of reproducible neural architecture search," *arXiv preprint arXiv:2001.00326*, 2020.
- [25] C. Zhang, P. Li, G. Sun, Y. Guan, B. Xiao, and J. Cong, "Optimizing fpga-based accelerator design for deep convolutional neural networks," in *Proceedings of the 2015 ACM/SIGDA international symposium on field-programmable gate arrays*, pp. 161–170, 2015.
- [26] C. Zhang, D. Wu, J. Sun, G. Sun, G. Luo, and J. Cong, "Energy-efficient cnn implementation on a deeply pipelined fpga cluster," in *Proceedings of the 2016 International Symposium on Low Power Electronics and Design*, pp. 326–331, 2016.
- [27] E. Real, S. Moore, A. Selle, S. Saxena, Y. L. Suematsu, J. Tan, Q. Le, and A. Kurakin, "Large-scale evolution of image classifiers," *arXiv preprint arXiv:1703.01041*, 2017.
- [28] W. Jiang, E. H.-M. Sha, X. Zhang, L. Yang, Q. Zhuge, Y. Shi, and J. Hu, "Achieving super-linear speedup across multi-fpga for real-time dnn inference," *ACM Transactions on Embedded Computing Systems (TECS)*, vol. 18, no. 5s, pp. 1–23, 2019.
- [29] W. Jiang, E. H.-M. Sha, Q. Zhuge, L. Yang, X. Chen, and J. Hu, "Heterogeneous fpga-based cost-optimal design for timing-constrained cnns," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 37, no. 11, pp. 2542–2554, 2018.
- [30] K. Wang, Z. Liu, Y. Lin, J. Lin, and S. Han, "Haq: Hardware-aware automated quantization with mixed precision," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 8612–8620, 2019.
- [31] R. J. Williams, "Simple statistical gradient-following algorithms for connectionist reinforcement learning," *Machine learning*, vol. 8, no. 3-4, pp. 229–256, 1992.
- [32] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255, Ieee, 2009.
- [33] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- [34] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4700–4708, 2017.
- [35] 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*, pp. 4510–4520, 2018.