



# A Comprehensive Survey on Hardware-Aware Neural Architecture Search

Hadjer Benmeziane, Kaoutar El Maghraoui, Hamza Ouarnoughi, Smail Niar,  
Martin Wistuba, Naigang Wang

## ► To cite this version:

Hadjer Benmeziane, Kaoutar El Maghraoui, Hamza Ouarnoughi, Smail Niar, Martin Wistuba, et al..  
A Comprehensive Survey on Hardware-Aware Neural Architecture Search. 2021. hal-03269441

**HAL Id: hal-03269441**

**<https://uphf.hal.science/hal-03269441>**

Submitted on 24 Jun 2021

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# A Comprehensive Survey on Hardware-Aware Neural Architecture Search

Hadjer Benmeziiane, Kaoutar El Maghraoui, *IEEE Member*, Hamza Ouarnoughi, Smail Niar, *IEEE Senior Member*, Martin Wistuba, Naigang Wang, *IEEE Member*

This work has been submitted to the IEEE for possible publication. Copyright may be transferred without notice, after which this version may no longer be accessible.

**Abstract**—Neural Architecture Search (NAS) methods have been growing in popularity. These techniques have been fundamental to automate and speed up the time consuming and error-prone process of synthesizing novel Deep Learning (DL) architectures. NAS has been extensively studied in the past few years. Arguably their most significant impact has been in image classification and object detection tasks where the state of the art results have been obtained. Despite the significant success achieved to date, applying NAS to real-world problems still poses significant challenges and is not widely practical. In general, the synthesized Convolution Neural Network (CNN) architectures are too complex to be deployed in resource-limited platforms, such as IoT, mobile, and embedded systems. One solution growing in popularity is to use multi-objective optimization algorithms in the NAS search strategy by taking into account execution latency, energy consumption, memory footprint, etc. This kind of NAS, called hardware-aware NAS (HW-NAS), makes searching the most efficient architecture more complicated and opens several questions. In this survey, we provide a detailed review of existing HW-NAS research and categorize them according to four key dimensions: the search space, the search strategy, the acceleration technique, and the hardware cost estimation strategies. We further discuss the challenges and limitations of existing approaches and potential future directions. This is the first survey paper focusing on hardware-aware NAS. We hope it serves as a valuable reference for the various techniques and algorithms discussed and paves the road for future research towards hardware-aware NAS.

**Index Terms**—AutoML, Hardware-aware Neural Architecture Search, Neural Network Acceleration, Edge Intelligence.

## I. INTRODUCTION

DEEP Learning (DL) systems are revolutionizing technology around us across many domains such as computer vision [1], [2], [3], [4], speech processing [5], [6], [7] and Natural Language Processing (NLP) [8], [9], [10]. These breakthroughs would not have been possible without the availability of big data, the tremendous growth in computational power, advances in hardware acceleration, and the recent algorithmic advancements. However, designing accurate neural networks is challenging due to:

- The variety of data types and tasks that require different neural architectural designs and optimizations.
- The vast amount of hardware platforms which makes it difficult to design one globally efficient architecture.

For instance, certain problems require task-specific models, e.g. EfficientNet [11] for image classification and ResNest [12] for semantic segmentation, instance segmentation and object detection. These networks differ on the proper configuration of their architectures and their hyperparameters. The hyperparameters here refer to the pre-defined properties related to the architecture or the training algorithm.

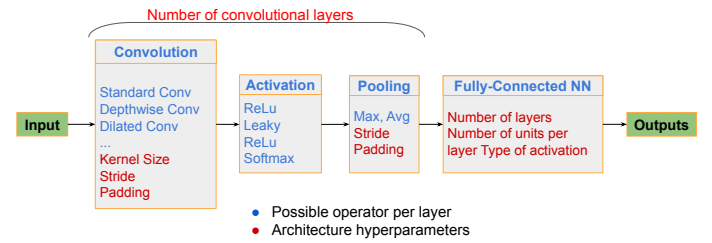


Fig. 1. Generic CNN architecture. For each layer an operator is chosen among a pre-defined list (convolution, dilated convolution, depthwise convolution, maxpooling, batch\_normalization...)

In general, the neural network architecture can be formalized as a Directed Acyclic Graph (DAG) where each node corresponds to an operator applied to the set of its parent nodes [13]. Example operators are convolution, pooling, activation, and self-attention. Linking these operators together gives rise to different architectures. A key aspect of designing a well-performing deep neural network is deciding the type and number of nodes and how to compose and link them. Additionally, the architectural hyperparameters (e.g., stride and channel number in a convolution, etc.) and the training hyperparameters (e.g., learning rate, number of epochs, momentum, etc.) are also important contributors to the overall performance. Figure 1 shows an illustration of some architectural choices for the type of convolutional neural network.

According to this representation, DL models can contain hundreds of layers and millions or even billions of parameters. These models are either handcrafted by repetitive experimentation or modified from a handful of existing models. These models have also been growing in size and complexity. All

H. Benmeziiane and H. Ouarnoughi and S. Niar are with Université Polytechnique Hauts-de-France, LAMIH/CNRS, Valenciennes, France. E-mail: (firstname.lastname@uphf.fr)

K. El Maghraoui and Naigang Wang are with IBM T. J. Watson Research Center, Yorktown Heights, NY 10598, USA. Email: (kelmaghr@us.ibm.com, nwang@us.ibm.com)

Martin Wistuba is with IBM Research AI, IBM Technology Campus, Dublin, Ireland. Email: (martin.wistuba@ibm.com)

of this renders handcrafting deep neural networks a complex task that is time-consuming, error-prone and requires deep expertise and mathematical intuitions. Thus, in recent years, it is not surprising that techniques to automatically design efficient architectures, or NAS for "Neural Architecture Search", for a given dataset or task, have surged in popularity.

In figure 2, we compare several deep learning models for the image classification task. Each dot in the plot corresponds to a given DL architecture that has been used for image classification. The dot size correlates with the size of the corresponding neural network in terms of the number of parameters. A quick look at the graph reveals the trend to design larger models to better Top 1 accuracy. However, the size is not necessarily correlated with better accuracy. There have been several efforts to conceive more efficient and smaller networks to achieve comparable Top 1 accuracy performance. We compare four classes of models: Handcrafted, Efficient handcrafted, NAS, and HW-NAS. Generally, throughout the years the handcrafted models rank high up in terms of accuracy but are much more complex in terms of the depth of the architecture or the large number of parameters. For instance, ViT-H [15], which is the state-of-the-art model as of December 2020, has over 600 million parameters and 32 layers. In the top right quadrant of the figure 2 (around the same region as most of the recently handcrafted models), we find some of the models that are automatically created by different NAS techniques. These latter techniques focus only on improving the model's accuracy without paying attention to the efficiency of the model in terms of its size and latency. Therefore, these NAS models are still large, with the number of parameters ranging between 100M and 500M.

Since 2015, we have noticed the rise of efficient handcrafted models. These models rely on compression methods (see section III-A1) to decrease the model's size while trying to maintain the same accuracy. MobileNet-V2 [16] and Inception [17] are good examples where the number of parameters is between 20M and 80M. This paper focuses on the class of hardware or platform-aware NAS techniques: *HW-NAS*. This class encompasses work that aims to tweak NAS algorithms and adapt them to find efficient DL models optimized for a target hardware device. HW-NAS began to appear in 2017 and since then achieved state of the art (SOTA) results in resource-constrained environments with Once-for-all (OFA) [18] for example.

A conventional NAS process requires the definition of three main components: the search space, the search strategy, and the evaluation methodology, as illustrated in figure 3. First, the *search space* (1) draws from a set of neural network architectures to define the neural network operators and how they are connected to form a valid network. It determines how the architectures can be formed and what architectures are allowed. For example, NASNet [19] introduced a fixed macro architecture where the search consists of finding the appropriate operators used within each block from a set of 12 specified operators. This search space is explored by a *search strategy* (2) which samples a population of network architectures' candidates. It evaluates the accuracy of the model using a specific *evaluation methodology* (3). The mea-

sured accuracy will then guide the search strategy to converge towards promising architectures in the search space. Generally, the evaluation component will train the architecture on the desired dataset, which often takes a considerable time. Many NAS algorithms have incorporated several techniques to speed up the training process described in Section VII-C.

NAS has proven its efficiency by proposing several SOTA models in Object Detection [20] and Image Classification [21]. However, these models are often composed of millions of parameters and require billions of floating-point operations (FLOP). This causes a large memory footprint and big FLOP preventing their usage in resource-constrained environments. Additionally, these models might require specific hardware (e.g. GPUs, TPUs, etc.) to allow their deployment in a reasonable time or real-time applications.

Since the inception of the very first NAS algorithm with reinforcement learning [22], there has been tremendous work to study efficient NAS. More recently, integrating hardware awareness in the search loop (i.e. HW-NAS) has attracted several researchers and has opened up interesting new research directions. Some HW-NAS efforts have demonstrated SOTA results and have balanced the trade-off between accuracy and hardware efficiency. For example, FBNet [23] has achieved SOTA results on ImageNet by using an objective function that minimizes both the cross-entropy error that leads to better accuracy and the latency that penalizes networks that are too slow.

This paper provides a detailed overview of existing HW-NAS research efforts and categorizes them according to their goals and problem formulation. There has been no survey dedicated to hardware-aware NAS in the literature to the best of our knowledge. With this survey, we provide a concise review of the NAS variants that focus on precision and hardware awareness. We hope that it will be beneficial to those with a basic understanding of NAS interested in (i) getting a comprehensible but accessible overview of how the multi-objective problem is solved and (ii) a succinct reference and guide of existing HW-NAS techniques.

Our review is divided into several sections. In section II, we highlight the relevant NAS surveys and explain the focus and importance of our survey compared to them. In section III, we present a brief overview of the efficient DL methods and where HW-NAS is situated among them. We also explain some of the concepts needed to understand the components of the NAS process. In section IV, we classify the different HW-NAS works based on their goals and target platforms. Based on the taxonomy, we introduce and define two search spaces: Architecture Search Space and Hardware Search Space in section V. In section VI, we formally define the HW-NAS optimization problem and how the multi-objective function is formulated. The following section *Search Strategies* is divided into three subsections. We first discuss, in section VII-A, the various algorithms used to explore the search space, including reinforcement learning and evolutionary algorithms. In section VII-B, we give insights on how NAS algorithms take care of the non-differentiable variables. And in section VII-C, we list the methods to speed up the search process and avoid training each sampled architecture. We dedicate section

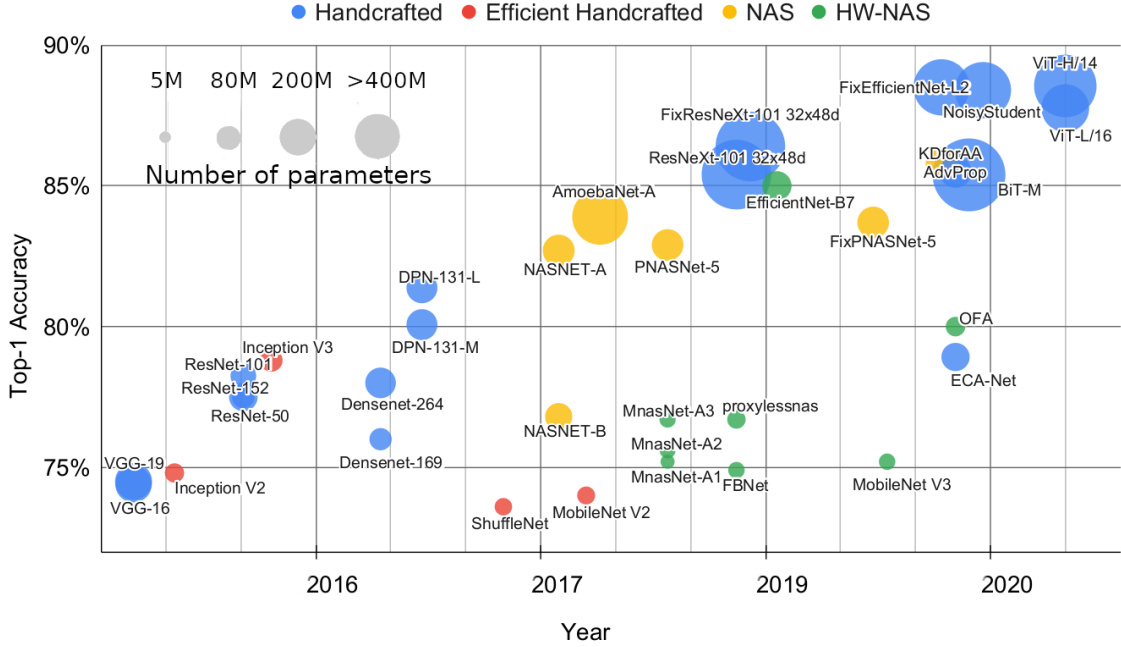


Fig. 2. Accuracy of various CNN models on ImageNet for Image Classification task with the number of parameters. Inspired by [14]

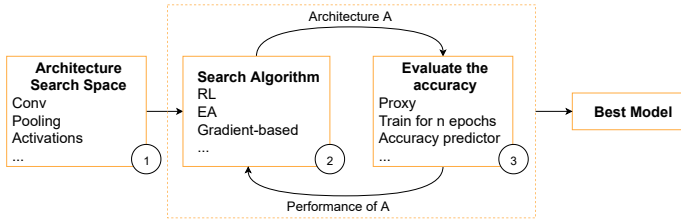


Fig. 3. Overview of conventional NAS components

VIII to the study of the hardware metrics used by HW-NAS and the methods used to measure them. In section IX, we introduce other considerations that either tackle the hardware efficiency objective, model compressions such as automatic quantization and pruning or add other objectives to the NAS like the robustness against adversarial attacks. We discuss the industrial adaptation of NAS by comparing several tools proposed by the commercial and research communities. We conclude in sections XI and XII by highlighting the challenges and limitations of current HW-NAS and discuss possible future directions. We summarized the referenced body of work in the table: <https://tinyurl.com/y6458skt>.

## II. EXISTING SURVEYS ON NAS

In the past few years, a large number of research papers have surveyed the state of the art neural architecture search algorithms. Two of the surveys we have examined stood out. The first survey, published in 2018 [24], describes the end-to-end design of a NAS algorithm. It details the components (i.e., Search Space, Search Strategy and Performance Estimation

Strategy) and gives a clear taxonomy of the optimization methods. In 2019, a similar survey was published by Wistuba et al. [25], which describes the different components of NAS algorithms with a focus on the optimization methods. It also includes a section that describes a high-level multi-objective neural architecture search and how it can be formulated and solved. While existing NAS surveys have succeeded at giving a thorough understanding of how general NAS algorithms work and explaining the various concepts that allow NAS techniques to achieve SOTA results, they fell short in including enough details about HW-aware NAS algorithms and the emerging efforts that strive to render NAS-generated models practical and amenable to deployment in today's constrained hardware platforms. Zhang et al. [26] surveyed a subset of hardware-aware NAS algorithms that only focus on FPGA platforms. There is a growing effort investigating the intersection of NAS and hardware platforms as depicted in Figure 4, where we can see the increasing number of research papers that are focusing on hardware-aware NAS. This work is a first attempt at filling the gap we have identified in existing NAS surveys. We aim to provide a detailed overview and analysis of existing HW-NAS research efforts. Table I shows the structure of the study and serves as a guide to the various sections and their content.

**Focus of this paper** In this survey, we provide a comprehensive overview of the architecture and hardware search spaces including how the HW-NAS problem is formulated, the search strategies, the hardware cost estimation methods, and the plethora of hardware platforms that have been targeted. We compare the used multi-objective search strategies such as reinforcement learning or evolutionary algorithm, along with some techniques that use non-differentiable parameters which



Background	Efficient Deep Learning	Model Compression HW-NAS Code Transformations	Section III-A1 Section III-A2 Section III-A3
	Algorithms	Reinforcement Learning Evolutionary Algorithm	Section III-B1 Section III-B2
Taxonomy of HW-NAS	Classification of HW-NAS based on their Goals		Section IV
Search Spaces	Architecture Search Space		Section V-A
	Hardware Search Space		Section V-B
HW-NAS Problem Formulation	Single-Objective Optimization	Two-Stage Search Constrained Optimization	Section VI-A1 Section VI-A2
		Scalarization NSGA-II	Section VI-B1 Section VI-B2
	Multi-Objective Optimization		
Search Strategies	Search Algorithm	Reinforcement Learning Evolutionary Algorithm Gradient-based Methods Bayesian Optimization & Random Search	Section VII-A1 Section VII-A2 Section VII-A3 Section VII-A4
		Non Differentiable Techniques	Over-parameterized Networks & their training Section VII-B
		Runtime Performance Optimization Strategies	Early Stopping Hot Start Proxy Datasets Accuracy Prediction Models Section VII-C
Hardware Cost Estimation Methods	Hardware Constraints Collection Techniques	Real-time measurements Lookup Table Analytical Estimation Prediction Models	Section VIII
Other Considerations for HW-NAS	Automatic Mixed-Precision Quantization		Section IX-A
	Automatic Pruning		Section IX-B
	Security and Reliability		Section IX-C
Discussions	Industrial Adaptation of NAS		Section X
	Challenges & Limitations	Benchmarking & Reproducibility	Section XI-A
		Transferability over Tasks	Section XI-B
		Transferability over HW Platforms	Section XI-C
	Outlook & Future Directions		Section XI-D

TABLE I  
CONTENT GUIDANCE FOR THIS SURVEY

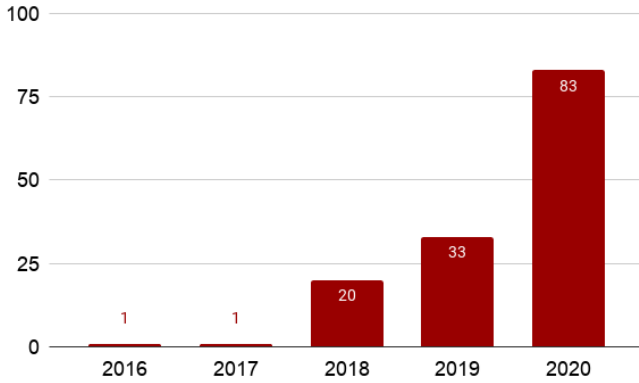


Fig. 4. Number of papers describing HW-NAS by Dec 2020. The top 5 conferences and journals are: NeurIPS, ECCV, IEEE Transactions on Pattern Analysis and Machine Intelligence, IJCNN, and MICCAI.

allow the use of an over-parameterized network. We also study different methods that collect the hardware cost metrics used as part of search optimization function. It is worth mentioning that this paper does not include Training Hyperparameter Search Optimization methods (THPO), and solely focuses on the operator selection and architectural hyperparameters.

### III. BACKGROUND

In this section, we define key terms needed to grasp the ideas surrounding HW-NAS and the concepts discussed in the

following sections.

#### A. Methodologies for Efficient Deep Learning

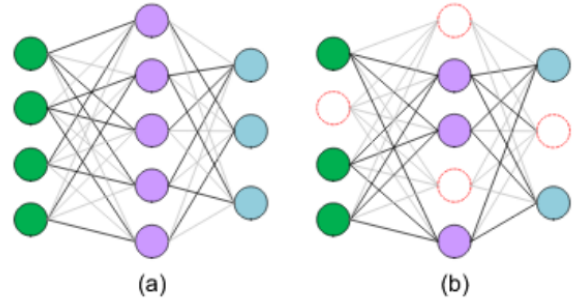


Fig. 5. Illustration of sparsification. (a) Weight Pruning. (b) Neuron Pruning. Gray lines correspond to pruned vertices (i.e. weights) and white nodes correspond to pruned neurons. Source [14]

The significant advances and breakthroughs of deep learning, that have propelled from academic and industrial research labs' circles, are now electrifying the computing industry and transforming the world. Deep learning is now being largely used to solve real-world problems. As deep learning is computationally demanding, most of the deployments happen in the cloud or on-premises data centers. However, with the arrival of powerful and low-energy consumption Internet of Things (IoT) devices and the growing need to take action in real or near-real-time, deep learning computations are increasingly

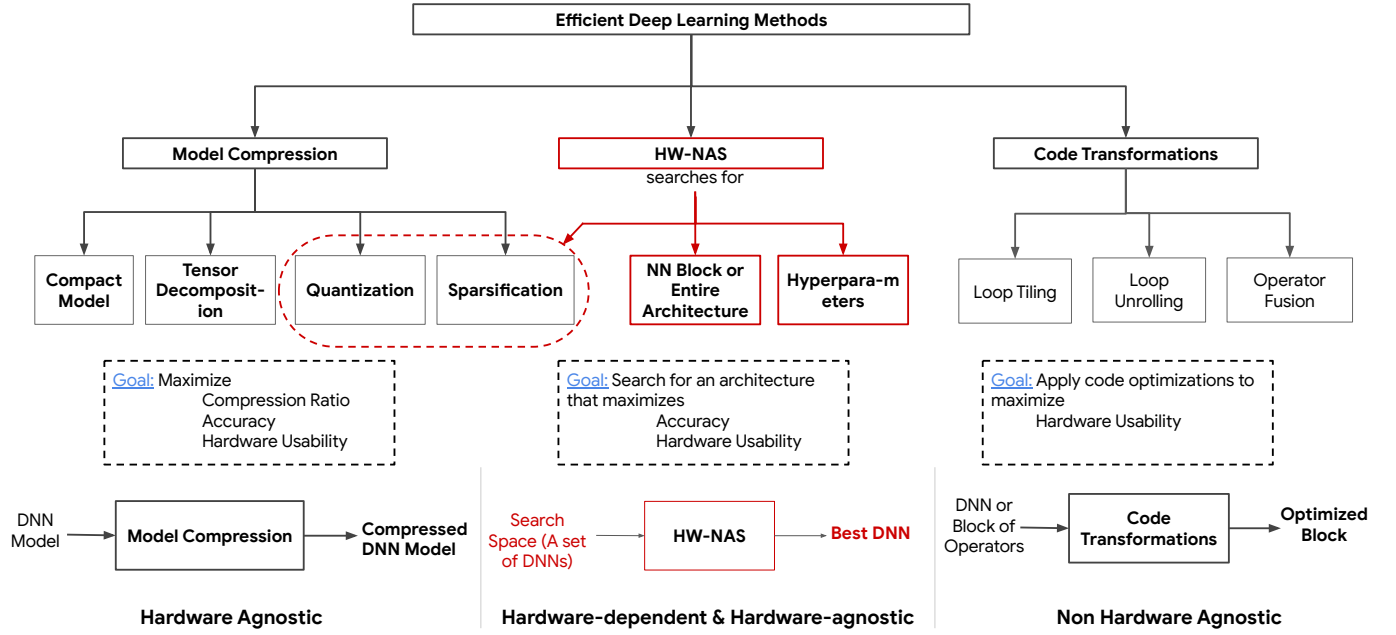


Fig. 6. Overview of efficient deep learning techniques

moving to the edge. Edge devices pose several challenges in this context, as they are constrained with limited energy and computational power. For example, autonomous driving cars depend on real-time object detection of the environment. They cannot tolerate the additional latency of sending the data to the cloud, processing the data and then sending it back to the edge device. Moreover, the compute capacity at the edge is significantly low, which does not match the increasing DL models complexity. This has motivated the research community to find innovative ways to reduce the DL models' size, their required number of floating operations, and their inference latency. This section presents an overview of the efficient deep learning techniques and where HW-NAS is situated among them.

Figure 6 illustrates the taxonomy of different techniques used to optimize deep learning models.

1) *Model Compression*: These methods aim to take one standard deep learning model such as ResNet or AlexNet and apply some optimizations that will decrease the model size and the number of FLOPs, i.e., increase the compression ratio, while trying to maintain the same accuracy. Relevant surveys [14], [27] on model compression classify the optimizations into these four classes:

- **Compact Model**: This technique modifies the standard operations used in DNNs. In a CNN, the standard convolution is replaced by more flexible convolution arithmetics that expand the number of feature maps and decrease the number of parameters such as dilated convolution [28] or separable depthwise convolution [29]. In an RNN, cells like S-LSTM [30], or JANET [31] simplify the gates and decrease the number of parameters compared to a regular LSTM.
- **Tensor Decomposition**: a tensor is the fundamental data

structure used in machine learning. It can represent vectors, matrices and even n-dimensional arrays. Therefore, shrinking the tensors allows accelerating DNNs and reducing their size. Tensor decomposition is an extension of the matrix decomposition techniques used in mathematical settings. Equation 1 formulates the matrix decomposition system where the number of parameters of A and B combined is smaller than the number of parameters of M.

$$M = AB \text{ with } M \in \mathbb{R}^{m \times n}, A \in \mathbb{R}^{m \times r}, B \in \mathbb{R}^{r \times n} \quad (1)$$

Note that tensor networks including hierarchical tensor representation (HT) [32] and tensor train decomposition (TT) [33] achieve higher compression rates in a fully-connected network since they usually contain more redundancy.

- **Quantization**: In deep learning, quantization [34] refers to converting data objects from 32-floating point to lower precision or a fixed point integer or even binary. These data objects can be the weights of the layer, the activations (the input data's internal representation), the error value, the gradient values and the weight update. Each method differs with the chosen number of bitwidth and the data objects that are quantized.
  - **Network Sparsification or pruning** [35], [36] attempts to compress the model by pruning some weights (edges) or operations (nodes). Usually, the decision of pruning is taken based on its importance, which is directly the weight values or learned via an attention layer.
- 2) *HW-NAS*: Another efficient deep learning technique is HW-NAS. In HW-NAS, we search for the architecture that maximizes the accuracy and hardware usability among a set of architectures. Note that some HW-NAS can be considered

under the model compression techniques as they search for the best bitwidth or the best way to prune. We further detail the search for hyperparameters, NN Block or full architectures in Section 4.

3) *Code Transformations*: An alternative approach that is gaining more attraction these recent years is to apply some code transformations that optimize the DNNs on the operator level [37]. These transformations are hardware-specific and require a compiler to apply the right transformation for the right hardware platform automatically.

## B. Search Algorithms

In this section, we give a high-level explanation of reinforcement learning and evolutionary algorithms; two of the most used search algorithms in NAS. The explanations are concise to ensure the core ideas are accessible to a broad audience, and so that the entire survey can be read end-to-end easily.

1) *Reinforcement Learning*: Reinforcement learning (RL) is a technique to allow an agent to learn by trial and error. The agent takes actions and interacts with an environment to maximize the accumulative reward. It is usually modelled as a Markov Decision Process (MDP).

*"Reinforcement Learning is learning what to do - how to map situations to actions -to maximize a numerical reward signal. The learner is not told which actions to take, but instead must discover which actions to yield the most reward by trying them."* [38]



Fig. 7. The agent-environment interaction in a Markov Decision Process. Source: [38]

As illustrated in figure 7, at each step, the agent observes the state of the environment sends and receives a reward for its previous action. It then selects its next action. The reward guides the agent to improve its policy such that better actions are chosen in the future. The policy of an agent is the algorithm that allows it to choose between multiple actions.

There are several variants and algorithms to train and update the policy of an agent. The basic one is the *sample-eval-update* method. This method is an iterative process where the agent samples one action, evaluates it using a Q-Learning [39] policy and updates the decision process accordingly. Another update method is *REINFORCE* [40]. It is a Monte-Carlo variant of policy gradients. The agent collects a set of actions (called trajectory)  $\tau$  using its current policy, and uses it to update the policy parameter. Since one full trajectory must be completed to construct a sample space, REINFORCE is updated in an off-policy way.

2) *Evolutionary Algorithms*: Evolutionary algorithms (EA) are optimization techniques that have three main characteristics:

- **Population-based**: EA maintain an entire set of candidate solutions. Each solution corresponds to a unique point in the search space of the problem. This set of solutions is called a population.
- **Fitness-oriented**: EA assign a fitness score to each solution, reflecting the quality of a solution.
- **Generations**: EA generate new populations by applying a series of transformations to the current population. These transformations are called mutations or crossover operations. Mutations transform one solution into another one, while crossover operations merge two solutions into a new one.

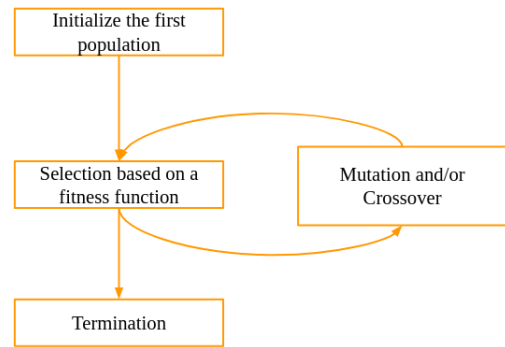


Fig. 8. Overview of the evolutionary algorithm steps.

A genetic algorithm is a type of evolutionary algorithms that encodes the individuals into numerical vectors, called chromosomes. The chromosome is represented as a set of parameters that defines a particular individual. A selection criterion is used to select a set of candidate individuals of which the fittest are mutated and recombined by crossover to create the next generation. This algorithm is the simplest version of the evolutionary methods as the chromosomes are numerical vectors and the mutations can be simple permutations.

## IV. TAXONOMY OF HW-NAS

Unlike conventional NAS, where the goal is to find the best architecture that maximizes model accuracy, hardware-aware NAS (HW-NAS) has multiple goals and multiple views of the problem. We can classify these goals into three categories (See figure 9 from left to right) :

- **Single Target, Fixed Configuration**: Most of existing HW-NAS fall under this category. The goal is to find the best architecture in terms of accuracy and hardware efficiency for one single target hardware. Consequently, if a new hardware platform has to be used for the NAS, we need to rerun the whole process and feed it the right values to calculate the new hardware's cost. These methods generally define the problem as a constrained or multi-objective optimization problem [53], [47], [23]. Within this category, two approaches are adopted:

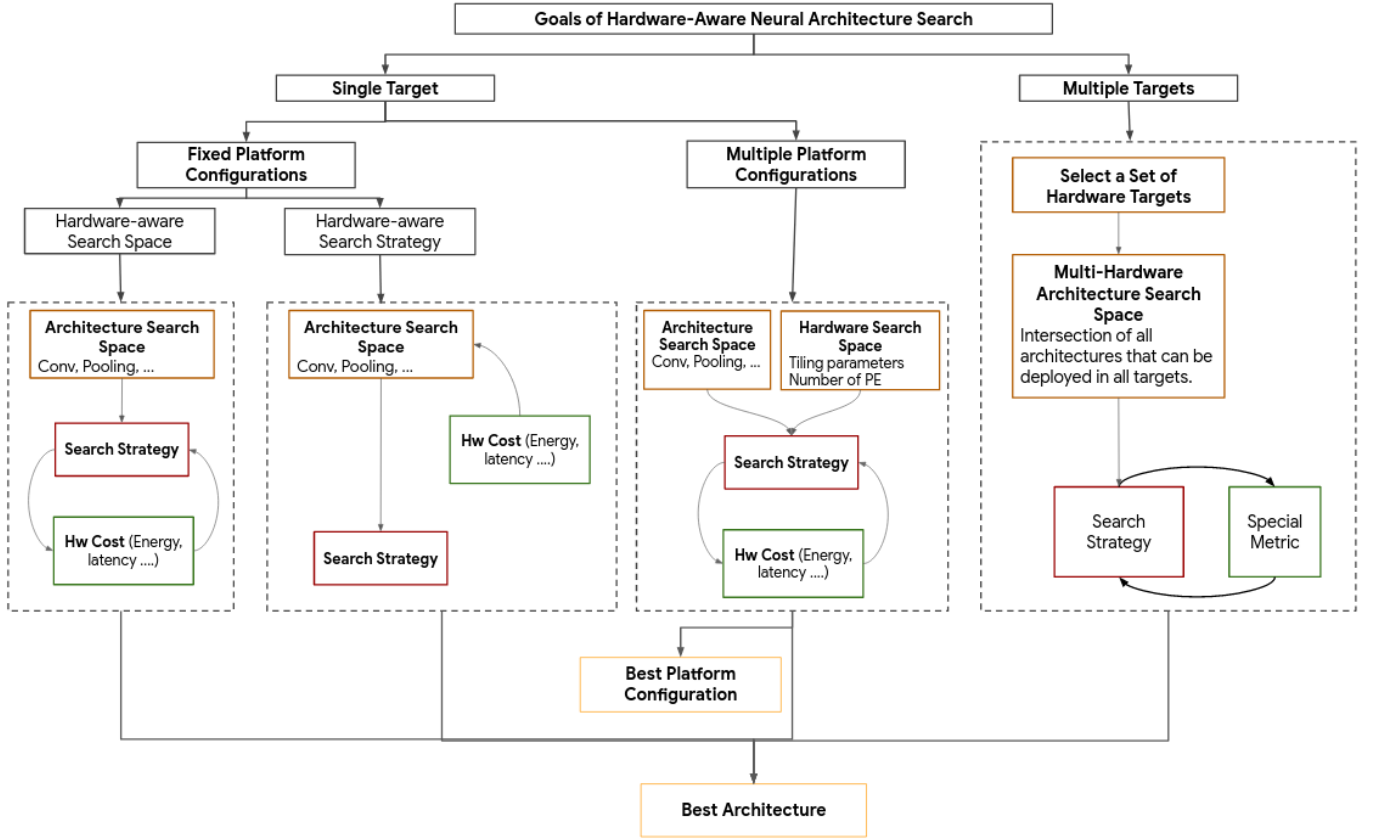


Fig. 9. Overview of different hardware-aware NAS designs.

## Category &amp; References

**Single Target, Fixed Configuration**

## Hardware-aware Search Strategy

[41], [42], [43], [44], [45], [46], [47], [48], [49], [50], [51], [52], [53], [23], [54], [55], [56], [57], [58], [59], [60], [61], [62], [63], [11], [64], [65], [66], [67], [68], [69], [70], [71], [72], [73], [74], [75], [76], [77], [78], [79], [80], [81], [82], [83], [84], [85], [86], [87], [88], [89], [90], [91], [92], [93], [94], [95], [96], [97], [98], [99], [100], [101], [102], [103], [104], [105], [106], [107], [108], [109], [110], [111], [112], [113]

## Hardware-aware Search Space

[72], [82], [114], [18], [115]

**Single Target, Multiple Configurations**

[116], [117], [118], [119], [120], [121], [122], [123], [124], [125], [126]

**Multiple Targets**

[127], [128]

TABLE II

TAXONOMY OF HARDWARE-AWARE NEURAL ARCHITECTURE SEARCH GOALS

- *Hardware-aware search strategy* where the search is defined as a multi-objective optimization problem. While searching for the best architecture, the search algorithm calls the traditional evaluator component to get the accuracy of the generated architecture but also a special evaluator that measures the hardware cost metric (e.g., latency, memory usage, energy consumption). Both model accuracy and hardware cost guide the search and enable the NAS to find the

most efficient architecture.

- On the other hand, the *Hardware-aware Search Space* approach uses a restricted pool of architectures. Before the search, we either measure the operators' performance on the target platform or we define a set of rules that will refine the search space; eliminate all the architectures' operators that do not perform well on the target hardware. For example, HURRICANE [72] uses different operator choices for three types of mobile processors: Hexagon DSP, ARM CPU and Myriad Vision Processing Unit (VPU). Accumulated domain knowledge from prior experimentation on a given hardware platform help narrow down the search space. For instance, they do not use depthwise convolutions for CPU, squeeze and excitation mechanisms for VPU and they do not lower the kernel sizes for a DSP. Such gathered empirical information helps to define three different search spaces according to the targeted hardware platform. Note that after defining the search space with these constraints, the search strategy is similar to the one used by conventional NAS, which means that the search is solely based on the accuracy of the architecture and no other hardware metric is incorporated.

- **Single Target, Multiple Configurations:** the goal of this category is not only to get the most optimal ar-

chitecture that gets the best accuracy but also to get an optimal architecture with latency guaranteed to meet the target hardware specification. For example, the authors of FNAS [116] define a new hardware search space containing the different FPGA specifications (e.g., tiling configurations). They also use a performance abstraction model to measure the latency of the searched neural architectures without doing any training. This allows them to quickly prune architectures that do not meet the target hardware specifications. In [129], the authors use the same approach for ASICs and define a hardware search space that contains various ASIC templates.

- **Multiple Targets:** In this third category, the goal is to find the best architecture when given a set of hardware platforms to optimize for. In other words, we try to find a single model that performs relatively well across different hardware platforms. This approach is the most favourable choice, especially in mobile development as it provides more portability. This problem was tackled by [127], [128] by defining a multi-hardware search space. The search space contains the intersection of all the architectures that can be deployed in the different targets. Note that, targeting multiple hardware specifications at once is harder as the best model for a GPU, can be very different to the best model for a CPU (i.e., for GPUs wider models are more appropriate while for CPUs deeper models are).

## V. SEARCH SPACES

Two different search spaces have been adopted in the literature to define the search strategies used in HW-NAS: the *Architecture Search Space* and the *Hardware Search Space*.

### A. Architecture Search Space

The Architecture Search Space is a set of feasible architectures from which we want to find an architecture with high performance. Generally, it defines a set of basic network operators and how these operators can be connected to construct the computation graph of the model. We distinguish two approaches to design an architecture search space:

- 1) *Hyperparameter Optimization for a fixed architecture:* In this approach a neural architecture is given including its operator choices. The objective is limited to optimizing the architecture hyperparameters (e.g., number of channels, stride, kernel size).
- 2) *True Architecture Search Space:* The search space allows the optimizer to choose connections between operations and to change the type of operation.

Both approaches have their advantages and disadvantages but it is worth mentioning that although former approach reduces the search space size, it requires considerable human expertise to design the search space and introduces a strong bias. Whereas the latter approach decreases the human bias but considerably increases the search space size and hence the search time.

Generally, in the latter approach, we distinguish three types (See figure 10):

- **Layer-wise Search Space**, where the whole model is generated from a pool of operators. FBNet Search Space [23], for example, consists of a layer-wise search space with a fixed macro architecture which determines the number of layers and dimensions of each layer where the first and last three layers have fixed operators. The remaining layers need to be optimized.
- **Cell-based Search Space**, where the model is constructed from repeating fixed architecture patterns called blocks or cells. A cell is often a small acyclic graph that represents some feature transformation. The cell-based approach relies on the observation that many effective handcrafted architectures are designed by repeating a set of cells. These structures are typically stacked and repeated a number of time to form larger and deeper architectures. This search space focuses on discovering the architecture of specific cells that can be combined to assemble the entire neural network. Although cell-based search spaces are intuitively efficient to look for the best model in terms of accuracy, they lack flexibility when it comes to hardware specialization [47], [23].
- **Hierarchical Search Space**, works in 3 steps: First the cells are defined and then bigger blocks containing a defined number of cells are constructed. Finally the whole model is designed using the generated cells. MNASNet [47] is a good example of this category of search spaces. The authors define a factorized hierarchical search space that allows more flexibility compared to a cell-based search space. This allows them to reduce the size of the total search space compared to the global search space.

In existing NAS research works, the authors define a macro-architecture that generally determines the type of networks considered in the search space. When considering CNNs, the macro architecture is usually identical to the one shown in figure 1. Therefore, many works [53], [23], [47], [127], [130] differ in the number of layers, the set of operations and the possible hyperparameters values. Recently, the scope of network type is changing. For instance, NASCaps [100] changes their macro-architecture to allow the definition of capsules. Capsules network [131] are basically cell-based CNNs where each cell (or capsule) can contain a different CNN architecture.

Other works like [66], [99] focus on transformers and define their macro-architecture as a transformer model. The search consists of finding the number of attention heads and their internal operations. When dealing with the hyperparameters only, the macro architecture can define a variety of network types. Authors in [75], [105] mix different definitions, transformers + CNN and transformers + RNN respectively. They define a set of hyperparameters that encompasses the pre-defined parameters for different network types at the same time.

Lately, more work [53], [64] have been considering the use of over-parameterized networks (i.e. supernetworks) to speedup the NAS algorithms. These networks consist of adding architectural learnable weights that select the appropriate operator at the right place. Note that these techniques have been applied to transformers as well [132].



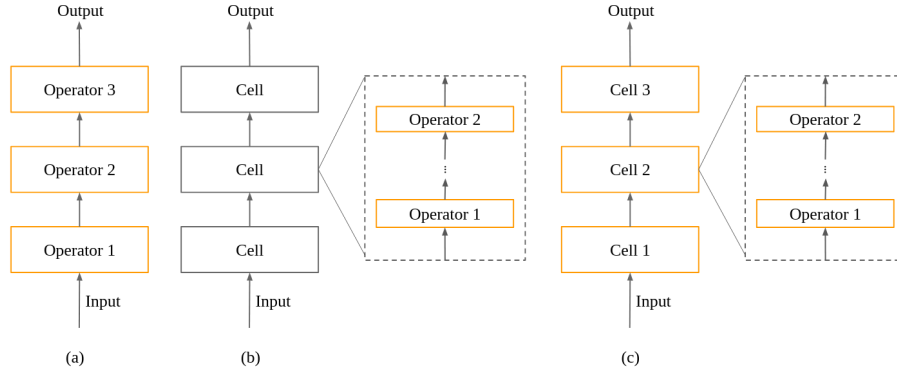


Fig. 10. Architecture search spaces types. (a) Global search space, (b) Cell-based search space, and (c) Hierarchical search space. In orange the operators considered during the search.

Finally, in some research efforts, the pool of operators/architectures is refined with only the models that are efficient in the targeted hardware [128], [127]. The search space size is considerably reduced by omitting all the architectures that cannot be deployed.

### B. Hardware Search Space (HSS)

Some HW-NAS methods include a HSS component which generates different hardware specifications and optimizations by applying different algorithmic transformations to fit the hardware design. This operation is done before evaluating the model. Although the co-exploration is effective, it increases the search space time complexity significantly. If we take FPGAs as an example, their design space may include: IP instance categories, IP reuse strategies, quantization schemes, parallel factors, data transfer behaviours, tiling parameters, and buffer sizes. It is arguably impossible to consider all these options as part of the search space due to the added search computation cost. Therefore, many existing strategies limit themselves to only few options.

Hardware Search Space (HSS) can be further categorized as follows:

- **Parameter-based:** The search space is formalized by a set of different parameter configurations. Given a specific data set, **FNAS** [116] finds the best performing model, along with the optimization parameters needed for it to be deployed in a typical FPGA chip for deep learning. Their HSS consists of four tiling parameters for the convolutions. **FNASS** [117] extends **FNAS** by adding more optimization parameters such as loop unrolling. The authors in [118], [119] used a multi-FPGA hardware search space. The search consists of dividing the architecture into pipeline stages that can be assigned to an FPGA according to its memory and DSP slices, in addition to applying an optimizer that adjusts the tiling parameters. Another example is presented in [120], where the adopted approach takes the global structure of an FPGA and adds all possible parameters to its hardware search space including the input buffer depth, memory interface width, filter size and ratio of the convolution engine. [121] searches the internal configuration of an

FPGA by generating simultaneously the architecture hyperparameters, the number of processing elements, and the size of the buffer. **FPGA/DNN** [133] proposes two components: *Auto-DNN* which performs hardware-aware DNN model search and *Auto-HLS* which generates a synthesizable C code of the FPGA accelerator for explored DNNs. Additional code optimizations such as buffer re-allocation and loop fusion on the resulting C-code are added to automate the hardware selection.

- **Template-based:** In this scenario, the search space is defined as a set of pre-configured templates. For example, **NASAIIC** [129] integrates NAS with Application-Specific Integrated Circuits (ASIC). Their hardware search space includes templates of several existing successful designs. The goal is to find the best model with the different possible parallelizations among all templates. In addition to the tiling parameters and bandwidth allocation, the authors in [122] define a set of FPGA platforms and the search finds a coupling of the architecture and FPGA platform that fits a set of pre-defined constraints (e.g., max latency 5ms)

In general, we can classify the targeted hardware platforms into 3 classes focusing on their memory and computation capabilities:

- **Server Processors:** this type of hardware can be found in cloud data centers, on premise data centers, edge servers, or supercomputers. They provide abundant computational resources and can vary from CPUs, GPUs, FPGAs and ASICs. When available, machine learning researchers focus on accuracy. Many NAS works consider looking for the best architecture in these devices without considering the hardware-constraints. Nevertheless, many HW-NAS works target server processors to speed up the training process and decrease the extensive resources needed to train a DL architecture and use it for inference.
- **Mobile Devices:** With the rise of mobile devices, the focus has shifted to enable fast and efficient deep learning on smartphones. As these devices are heavily constrained with respect to their memory and computational capabilities, the objective of ML researchers shift to assessing the trade-off between accuracy and efficiency. Many HW-



NAS algorithms target smartphones including FBNet [23] and ProxylessNAS [53] (refer to table III). Additionally, because smartphones usually contain system on chips with different types of processors, some research efforts [79] have started to explore ways to take advantage of these heterogeneous systems.

- **Tiny Devices:** The strong growth in use of microcontrollers and IoT applications gave rise to TinyML [134]. TinyML refers to all machine learning algorithms dedicated to tiny devices, i.e., capable of on-device inference at extremely low power. One relevant HW-NAS method that targets tiny devices is MCUNet [123], which includes an efficient neural architecture search called TinyNAS. TinyNAS optimizes the search space and handles a variety of different constraints (e.g., device, latency, energy, memory) under low search costs. Thanks to the efficient search, MCUNet is the first to achieve >70% ImageNet top-1 accuracy on an off-the-shelf commercial microcontroller.

Targeted HW	References
CPU	[54], [72], [74], [76], [93], [120], [79], [80], [135], [73]
GPU	[43], [46], [51], [52], [48], [49], [50], [55], [56], [59], [58], [81], [73], [69], [74], [75], [83], [93], [90], [136], [121], [137], [61], [62], [67], [68], [70], [71], [95], [96], [97], [102], [99], [107], [108], [109], [113]
MCU	[123], [63], [111], [138], [23], [89]
Mobile Devices	[47], [53], [23], [57], [59], [66], [18], [60], [87], [85], [86], [110], [84], [80], [99]
ASICs	[72], [118], [119], [91], [103], [106], [126], [125], [100]
Edge TPU	[99], [85], [82]

TABLE III

CLASSIFICATION OF HW-NAS BASED ON THEIR TARGETED HARDWARE.

### C. Current Hardware-NAS Trends

Figure 11 shows the different types of platforms that have been targeted by HW-NAS in the literature. In total, we have studied 126 original hardware-aware NAS papers. By target, we mean the platform that the architecture is optimized for. Usually the search algorithm is executed in a powerful machine, but that is not the purpose of our study. “*No Specific Target*” means that the HW-NAS incorporates hardware agnostic constraints into the objective function such as the number of parameters or the number of FLOPs. In the figure, the tag “*Multiple*” means multiple types of processing elements have been used in the HW platform. Table III gives the list of references per targeted hardware.

In the figure, we note that the number of research papers targeting GPUs and CPUs has more or less remained constant. However, we can clearly see that FPGAs and ASICs are gaining popularity over the last 3 years. This is consistent with the increasing number of deep learning edge applications. Two recent interesting works are [127], [128] both of which target multiple hardware platforms at once.

In figure 12, we illustrate the different DNN operations that compose the architecture search space. First, we divide the CNN into two groups, *standard CNN* which only utilizes

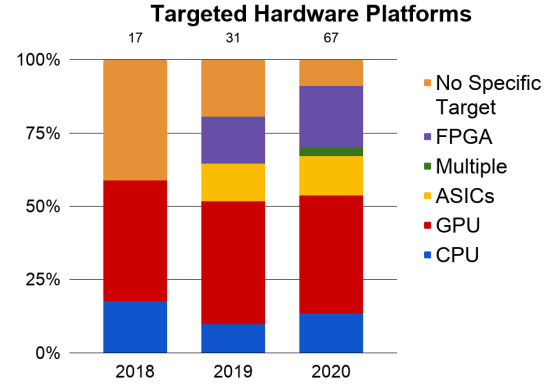


Fig. 11. Statistics on targeted platforms

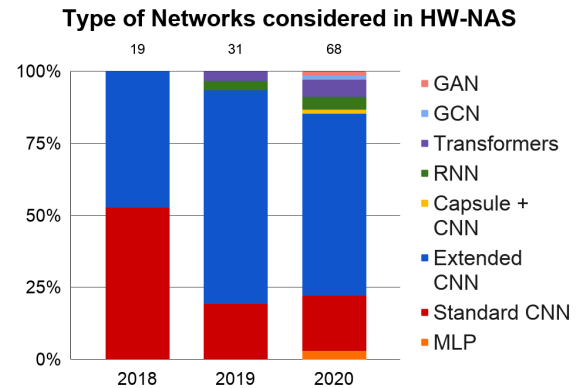


Fig. 12. Statistics about the type of networks described by the HW-NAS search spaces

a standard convolution and *extended CNN* which involves special convolution operations such as the depthwise separable convolution or grouped convolutions. NAS has been mostly dominated by convolutional neural networks as shown in the figure. However, recent works have started explore more operators by incorporating capsule networks [100], transformers [139], and GANs [94].

## VI. HARDWARE-AWARE NAS PROBLEM FORMULATION

Neural Architecture Search (NAS) is the task of finding a well-performing architecture for a given dataset. It is cast as an optimization problem over a set of decisions that define different components of deep neural networks (i.e., layers, hyperparameters). This optimization problem can simply be seen as formulated in equation 2.

$$\max_{\alpha \in A} f(\alpha, \delta) \quad (2)$$

We denote the space of all feasible architectures as  $A$  (also called search space). The optimization method is looking for the architecture  $\alpha$  that maximizes the performance metric denoted by  $f$  for a given dataset  $\delta$ . In this context,  $f$  can simply be the accuracy of the model.

Although it is important to find networks that provide high accuracy, these NAS algorithms tend to give complex

models that cannot be deployed on many hardware devices. To overcome this problem, practitioners consider other objectives, such as the number of model parameters, the number of floating-point operations, and device-specific statistics like the latency or the energy consumption of the model. Different formulations were used to incorporate the hardware-aware objectives within the optimization problem of neural architecture search. We classify these approaches into two classes, single and multi-objective optimization. The single objective optimization can be further classified as two-stage or constrained optimization. Similarly, the multi-objective optimization approach can be further classified as single or multi-objective optimizations. Please refer to figure 13 for a summary of these approaches. These 2 classes are further detailed with examples from the literature in the following sections:

#### A. Single-Objective Optimization

In this class, the search is realized considering only one objective to maximize, i.e the accuracy. Most of the existing work in the literature [47], [53], [78], [122], [120], that tackle the hardware-aware neural architecture search, try to formulate the multi-objective optimization problem into a single objective to better apply strategies like reinforcement learning or gradient-based methods. We can divide this class into two different approaches: Two-stage optimization and constrained optimization.

1) *Two-Stage optimization*: In this first category, we retain the original formulation of the NAS problem and then specify the model for deployment. This approach is suboptimal as the final architecture proposed by the NAS is not always the one that gives the best performances on the hardware device. Two-stage optimization consists of applying NAS methods to obtain the best-performing architecture and then in a second stage, specialize this architecture for deployment on a target hardware platform. This specialization performs a series of optimization to fit the hardware requirements. [62] applies a reinforcement learning agent to find the best quantization bitwidth and pruning level after selecting the most accurate model. [18] searches over a pre-trained and selected architecture to find the most efficient one in terms of latency and energy consumption.

2) *Constrained optimization*: In this approach, the hardware-aware characterizations are considered as constraints in the original NAS formulation. The constraints take the form of thresholds to be respected. For example, inference time, energy consumption, memory occupation, etc. The conditions are added as constraints to the optimization problem to enforce requirements like fewer parameters or faster inference time. The threshold and the trade-off between different constraints can be adapted to practical requirements. For such cases, the single-objective optimization problem defined in Equation 2 turns into a constrained optimization problem defined by

$$\begin{aligned} & \max_{\alpha \in A} f(\alpha, \delta) \\ & \text{subject to } g_i(\alpha) \leq T_i \quad \forall i \in I \end{aligned} \quad (3)$$

Here,  $g_i$  corresponds to the different constraints taken into account (e.g., latency, memory, energy consumption) and  $T_i$  denotes the respective threshold. As most of the optimization methods used by NAS (i.e. reinforcement learning and evolutionary algorithms) were designed for unconstrained optimization problems, this formulation is hard to be adopted directly. Therefore, many researchers turned to penalty methods to transform the equation into a single objective function that contains the hardware constraints as well as the accuracy measurement [47], [23], [53]. For example, MNASNet [47], uses equation (4), where  $f$  is the accuracy measurement function, LAT is the latency of the model and  $T$  is the threshold. They use a learnable parameter  $w$  to control the effect of the hardware constraints on the global objective function.

$$\max_{\alpha \in A} f(\alpha) \cdot [LAT(\alpha)/T]^w \quad (4)$$

ProxylessNAS [53] uses a loss function that comprises of the cross-entropy (CE) loss and hardware-aware constraints.

$$\mathcal{L} = \mathcal{L}_{CE} + \lambda_1 \|w\|^2 + \lambda_2 E[\text{latency}] \quad (5)$$

Equation 5 illustrates the loss calculated by the reinforcement learning agent used by ProxylessNAS.  $\lambda_1$  and  $\lambda_2$  are learnable parameters that adjust the effect of the efficiency of the overall loss. Specifically, a policy is learned that decides whether to add, remove or keep a layer as well as whether to alter its number of filters.

#### B. Multi-Objective Optimization

Another approach is to handle multiple fronts in the formalism of a multi-objective optimization problem defined as:

$$\max_{\alpha \in A} f_1(\alpha, \delta), f_2(\alpha, \delta), \dots, f_n(\alpha, \delta) \quad (6)$$

In this scenario, there is often no single optimal solution that simultaneously maximizes every objective function. Additionally, some objective functions can be conflicting. For instance, trying to minimize the number of parameters while aiming at maximize the accuracy. Therefore, in these situations, the task boils down to finding Pareto-optimal solutions. These techniques can:

- 1) transform the problem to a single-objective optimization using scalarization method, also called weighted sum method or,
- 2) solve the multi-objective optimization problem using dedicated heuristics or meta-heuristics such as genetic algorithms or tabu search [137]. In general, this second approach provides not one optimal solution but a set of solutions that form the optimal Pareto front of the multi-objective optimization problem.

1) *Scalarization Methods*: One way to solve the multi-objective optimization problem is to use a scalarization approach. Equation 7 formulates this method. We use a parameterized aggregation function  $h$  to transform the multi-objective optimization problem into a single-objective optimization problem.

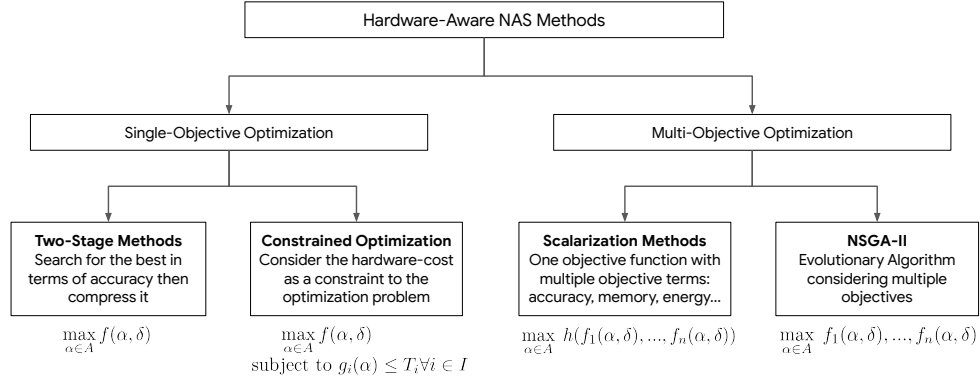


Fig. 13. HW-NAS problem formulations.

$$\max_{\alpha \in A} h(f_1(\alpha, \delta), f_2(\alpha, \delta), \dots, f_n(\alpha, \delta)) \quad (7)$$

The function  $h$  can be a weighted sum, a weighted exponential sum, a weighted min-max or a weighted product. However, in most situations, not all Pareto optimal solutions can be found in solving this problem with a fixed setting of the weights. Therefore, the problem is solved for multiple values of the vector  $w$  which requires multiple optimization runs. To mitigate the cost of having multiple run, researchers commonly use a set of fixed weights according to the desired trade-off between the objectives and the practitioners' preferences. Thanks to the scalarization, the problem becomes a single-objective optimization problem which can be solved by any optimizing methods discussed in Section VII. For example, [46] proposes to use the weighted sum as the objective function. The proposed formulation of this function is described by equation 8.  $ACC$  refers to the accuracy metric,  $E$  refers to the energy consumed by the architecture  $\alpha$  and  $w$  is a learned parameter to adjust the effect of the energy on the reward function.

$$\max_{\alpha \in A} w \cdot ACC(\alpha, \delta) - (1 - w) \cdot E(\alpha) \quad (8)$$

2) *NSGA-II*: An alternative approach is to use the elitist evolutionary algorithm NSGA-II [140] which claims that a linear combination of objectives is suboptimal. Thus, it is important to cast the problem as a multi-objective optimization problem, where a series of models is found along the Pareto front of multiple objectives such as accuracy, computational cost or inference time, and number of parameters. In this algorithm, the architectures are divided into fronts based on their dominance. The architecture in the  $i$ -th front is only dominated by all the architectures in the  $1, \dots, i-1$  fronts. Within each front, the architectures are prioritized by the crowding distance, which is computed by the sum of all the neighborhood distances across all the objectives. Hardware-aware NAS works [136], [50], [141] have been using NSGA-II algorithm to ensure the exploration of diverse architectures in the search space. Moreover, NSGANet [50] uses Bayesian Optimization to profit from search history. MoreMNAS [141] uses a hybrid search strategy combining NSGA-II with reinforcement learning to regulate arbitrary mutations.

## VII. SEARCH STRATEGIES

In this section, we will review all the search strategies used in hardware-aware NAS and how the hardware cost is integrated into each one of them. Following the algorithm presented in 1, the search strategy must define:

- An accuracy evaluation method: usually it is the actual training of the architecture on a subset of the dataset provided by the user  $D_{train}$ . However, training each architecture while searching is computationally expensive and takes several GPU hours. Therefore, some speedup methods that try to decrease the training time and estimate the accuracy without effectively training an architecture are used and presented in section VII-C.
- A hardware cost evaluation method: the hardware metrics need to be measured either by a real-time execution of each architecture on the targeted platform or an estimation method. The approaches used are presented in section VIII.
- A search algorithm: The search algorithm defines how the architectures are sampled from the search space, either with a generator or at random, and how the search explores this search space by updating the sampling strategy to better fit the best models according to an objective function. See the next section VII-A.

---

### Algorithm 1: A generalized pipeline of HW-NAS

---

**Input:** The Search Space  $S$ , a Dataset  $D$ , Accuracy Evaluation Methods  $E_{acc}(S, w(S), D)$ , Hardware Evaluation Methods  $E_{hw}(S, w(s), D)$ , an Aggregation function  $f$

**Result:** Optimal Deep Learning Architecture  
**repeat**

Sample an architecture  $s$  from the search space  $S$ ;  
Train  $s$  on  $D$  to obtain the trained weights  $w(s)$ ;  
Compute  $acc = E_{acc}(s, w(s), D)$ ; Compute  $hw\_cost = E_{hw}(s, w(s), D)$ ;  
**if**  $f(acc, hw\_cost)$  surpass the best seen value **then**  
| Update best architecture  $< s, w(s)$ ;  
| Update sampling algorithm with  $f(acc, hw\_cost)$ ;  
**until** convergence or time limit is reached;

---

Strategy	References
<b>Evolutionary Algorithm</b>	[44], [50], [54], [55], [117], [68], [61], [73], [88], [87], [91], [93], [96], [92], [123], [100], [124], [106], [109]
<b>Reinforcement Learning</b>	[45], [46], [47], [49], [51], [114], [116], [58], [60], [62], [118], [67], [66], [70], [119], [76], [77], [78], [129], [120], [82], [85], [121], [89], [90], [94], [135], [122], [127], [105], [108]
<b>Gradient-based Methods</b>	[142], [53], [23], [57], [143], [59], [64], [18], [49], [74], [81], [83], [84], [86], [72], [128], [95], [97], [99], [101], [102], [103], [104], [125], [107], [110], [111], [126]
<b>Bayesian Optimization</b>	[52], [63], [75], [112], [144]
<b>Random Search</b>	[115]
<b>Hybrid</b>	EA + RL [141], [145] EA + Bayesian Optimization [113] EA + Gradient-based[69]

TABLE IV  
SEARCH STRATEGIES OF HW-NAS

### A. Search Algorithm

1) *Reinforcement Learning*: Most HW-NAS methods use reinforcement learning to search for the best architecture [47], [26], [129], [120], [122], [146], [119], [78] because the NAS problem is easily modeled as a Markov Decision Process. The RL controller samples an architecture from the search space and is rewarded according to its accuracy and hardware cost. The agent will then adjust its weights to generate better models. Different works differ on how they represent the agent's policy (set of actions) and how they optimize it.

Using reinforcement learning, **MNASNet** [47] tries to find the Pareto optimal solution of the objective function described in equation 4. It uses a *sample-eval-update* loop to train its RNN controller. To generate a block, the controller will first choose two hidden states (i.e., outputs of previous blocks) as inputs. Then, it will select an operation to apply to each one of them. Finally, it selects a combination method (e.g., addition or concatenation) to obtain the final output of the block. This implementation has been defined by NASNet [147] before. Once a model is sampled, it is trained on the target task to get its accuracy and deployed on real phones to get its latency. The system computes the reward value and adjusts the controller parameters accordingly.

A similar approach was used by **FPNet** [78] but the RNN controller predicts only the architectural hyperparameters (i.e., number of filters, filter height, filter width, stride height, stride width, etc.) while keeping The macro architecture fixed.

**Codesign-NAS** [120] proposes to use reinforcement learning to explore both architecture and hardware search spaces. The authors investigated three RL-based search strategies and use the REINFORCE algorithm to improve the accuracy and efficiency of image classification on FPGA. The first strategy consists of combining the two search spaces and updating the CNN and accelerator options at the same time. The second method uses two different controllers, one to learn the CNN architecture and another one to select the best FPGA options. The last one is a conventional NAS where they separately search for the best model in terms of accuracy then as a

completely separate step look for the most efficient model. An expected result is that this latter strategy gives bad results when it comes to the constrained environment. An interesting result is that the second method (i.e., phase search) seems to be the most promising achieving higher rewards in most of their experiments.

2) *Evolutionary Algorithm*: Another popular strategy in conventional NAS [22], [148] is the use of evolutionary algorithms. Generally, neuro-evolutionary NAS evolves a population of models, sample some models to generate offsprings by applying some mutations (recombination is not used in neuro-evolutionary NAS), and finally evaluate the fitness of the offsprings and update the new generation by adding the best ones to the population. When it comes to integrating the hardware constraints to the NAS algorithms, some research efforts have used evolutionary algorithms [132], [123], [100], [72].

**TinyNAS** [123] uses a hardware-aware search space approach. It first optimizes the search space to fit the tiny and diverse resource constraints and then performs an evolution search algorithm within the optimized search space to find the most accurate model. The evolution search is performed on a trained super network that contains all the possible sub-networks. First, they sample 100 satisfying networks that fit the resource constraints. Then, just like conventional NAS, they measure the validation accuracy of each model, mutate the offsprings and update the new generation. This process is repeated for 30 iterations.

**Once-for-all** [18] proposes a special technique to train the supernet using gradient-based methods called **progressive shrinking**. A supernet is another name for the over-parameterized network. In their version, the supernet includes the highest values for the width, depth and channel number in each convolution. The process begins by training for larger width, depth and channel numbers and then fine-tune the network for the smaller sizes. Once the supernet is trained, they specialize the network using an evolutionary algorithm for different hardware settings with a combined cost of accuracy and latency.

**HAT** [132] is interesting because to our knowledge it is the only work that is trying to search for efficient transformers, targeting NLP tasks. The authors perform an evolutionary search with hardware latency constraints on a SuperTransformer. This means the engine only adds SubTransformers with latency smaller than the hardware constraint to the population.

**NASCaps** [100] proposes a NAS framework that generates Capsule Networks along with CNN. The proposed framework uses a multi-objective Genetic Algorithm (based on NSGA-II, see section 4.2) to pick the Pareto optimal solutions. The two key operations are the *crossover* and *mutation*. In the crossover, they define the splitting point by ensuring that the generated DNN is made up of at least on initial convolution layer and a minimum of two capsules. No standard convolution is placed between two capsule layers. The mutation is performed by randomly choosing one of the layer descriptors from the candidate network and modifying one of the main parameters of the selected layer.

3) *Gradient-Based Methods*: Arguably the most promising search strategy promising in terms of results, Gradient-based methods are increasingly used by hardware-aware NAS [53], [23], [128], [64], [149] and NAS generally. Running the search separated from the evaluation requires a lot of time and computation. Therefore, a common idea is to have a supernet that can emulate any child model in the search space. This means that different parts of the graph share weights between their common edges. This idea of *weight sharing* has the advantage of considerably reducing the search time. Gradient-based methods train the supernet to simultaneously get the architecture parameters and weights. This technique has been initiated by DARTS [150].

**ProxylessNAS** [53] is one of the papers pioneering this method. It defines a supernet with binary architecture parameters (i.e., 1 implies that the operator is selected in the layer and 0 otherwise). The loss function combines the cross-entropy and the latency to better update the weights and architecture parameters. An approach similar to BinaryConnect [151] is used to update the binary architecture parameters using an approximation of the gradient w.r.t architecture parameters.

**FBNet** [23] proposes to use differentiable neural architecture search to discover hardware-aware efficient CNNs. They also use a combination of the cross-entropy and latency to train their supernet using stochastic gradient descent. Rather than using binary parameters, they relax the problem of finding the best architecture to finding a distribution that yields to the best model.

**SqueezeNAS** [64] focuses on semantic segmentation and uses a method very similar to FBNet [23]. **HTAS** [128] uses a gradient-based method to find the best width and depth for their transformable CNNs.

**XNAS** [149] proposes to use the prediction with expert advice theory [152] for the selection. It leverages the Exponentiated-Gradient algorithm (EG) [153] rather than the classical gradient descent which prevents the decay of architecture weights to promote the selection of arbitrary architectures.

4) *Random Search & Bayesian Optimization*: The most convenient and easiest search strategy to implement is the random search strategy. Generally, random search and Bayesian optimization are used for hyperparameter optimization. Therefore, most of the existing works that have adopted random search optimize the architectural hyperparameters within a fixed macro architecture. They argue that it is more important to design the architecture search space for the targeted hardware platform than to complicate the search strategy and incorporate the hardware constraints in the objective function. NASNet [19], for example, tried both methods (i.e., reinforcement learning and random search). They found that with reinforcement learning the results are slightly better (Top-1 accuracy on CIFAR-10, 0.912 - 0.925).

Li et al. [154] investigate the use of random search on two standard NAS benchmarks (i.e., PTB and CIFAR-10). They use an approach that is similar to ProxylessNAS [53] which allows them to train a single network at a time and thus reduce the memory footprint with weight sharing. As a result, they show that random search with early-stopping is a competitive NAS baseline as it outperforms ENAS [155].

Furthermore, random search with weight-sharing outperforms random search with early stopping, achieving SOTA results on PTB. A more recent investigation [59] showed that random search cost time is not negligible and comparable to NAS methods.

## B. Non-differentiable Hardware Constraints

In many of these search strategies, the authors use an over-parameterized network. Therefore, the loss function must be differentiable w.r.t the architectural parameters. Also, we need to check that the incorporated hardware cost is differentiable. Several methods have been used to make the gradient computation over discrete variables possible. In this section, we review some extensions to the differentiable architecture search space [150].

a) *Gumbel Softmax*: [156] One way to relax the discrete variables is to use the Gumbel softmax function. It helps insert some random noise following the Gumbel distribution so that the gradient computation is possible. This technique was used by FBNet [23], [157].

b) *Estimated Continuous Function*: ProxylessNAS [53] mimics the concept of BinaryConnect [151]. They approximately estimate the gradient w.r.t the architecture parameters using the gradient w.r.t the binary gates. To reduce the computational cost, they compare the gates two-by-two by factorizing the task of using one out of N paths into multiple binary selection tasks.

c) *REINFORCE algorithm*: An alternative approach to BinaryConnect is also proposed by ProxylessNAS [53]. They utilized REINFORCE to train the binarized weights. Furthermore, they combine the gradient-based update rule to the REINFORCE updates to form a new general update rule for the architecture parameters.

## C. Runtime Performance Optimization Strategies

In this section, we discuss the various methods that are used to speed up the NAS algorithms. We are not going to mention weight sharing as it was described in section VII-A3, but we note that it is another way to accelerate the process tied to the search strategy used.

While NAS methods are efficient at finding state of the art architectures, their search cost is extremely high. Undoubtedly, the most time-consuming component of the NAS process is the training of each model to obtain the validation accuracy and thus making an evaluation. This training alone requires hours for just one architecture on a GPU. For example, training ResNet-50 requires 29 hours on 8 Tesla P100 GPUs [158]. For NAS, we want to estimate which architecture(s) is the most accurate without needing the exact accuracy. Therefore, we will review the accuracy estimation methods.

a) *Early Stopping*: Early stopping is a technique that is used by several NAS works [47], [53], [154]. The idea is to train the models for few epochs (usually five epochs) and to take the validation accuracy of this premature model as an approximation of the performance of the fully trained model.

*b) Hot start:* HotNAS [122] proposes the idea of hot starting. The search heuristics does not start from a random model but it starts with an efficient model. Starting from ProxylessNAS [53] and MNASNet [47] models, they obtained models with better latency and slightly lower accuracy: 91.47 and 92.2 for HotNAS and ProxylessNAS respectively. In addition to obtaining efficient models, the NAS search is faster as we start from a suboptimal model.

*c) Proxy datasets:* Using proxy datasets like CIFAR-10 helps train the model with fewer data elements and then fine-tune it to be used for higher resolutions and more complex tasks.

*d) Accuracy Prediction Models:* Another approach to speed up the training process is to directly predict the accuracy based on the architecture specifications and dataset characteristics. The most used prediction models are the following.

**Peephole** [159] is a framework to predict network performance before training based on its architecture. It encodes the architectures into vectors and inputs them onto an LSTM layer. The result is concatenated to the number of epochs and passed to an MLP which outputs the predicted accuracy. Similarly, PNAS [160] uses an RNN model to handle variable-sized inputs which consist of multiple LSTM layers. The final LSTM hidden state goes through a fully-connected layer and sigmoid layer to regress the validation accuracy.

Another approach by [161] uses a simple regression model (SRM regressor) with features like the number of weights and the number of layers as well as the hyperparameters, to predict the learning curves during the training process.

**NeuNetS** [55] defines a special component called *TAP*, which stands for Train-less Accuracy Predictor. *TAP* is designed to perform fast and reliable CNN accuracy predictions.

Recently, [162] uses a semi-supervised performance predictor. To locate latent feature representations, the architecture graph representation is used by an auto-encoder, which is then fine-tuned using a graph similarity calculation. Finally, a graph convolution neural network is used to output the final performance.

## VIII. HARDWARE COST ESTIMATION MODELS

An important component in HW-NAS is the hardware cost measurements. First of all, many metrics have been used in order to evaluate the hardware efficiency of an architecture including the number of FLOPs, the number of parameters, the latency or execution time of the inference, the energy consumption, the memory footprint, the area of the hardware platform, etc.

*a) FLOPs & Model Size:* The first HW-NAS approaches that were published in 2016 and 2017 [41], [42] use the number of parameters and number of FLOPs as an objective function to minimize. These techniques assume that the number of operations is positively correlated to the execution time. However, recent works have proved that two models can have the same number of FLOPs but different latencies [163], [72], [132]. For example, NASNet-A and MobileNetV1 have roughly similar number of FLOPs, yet, NASNet-A can have slower latency due to the hardware-unfriendly structure. Therefore, using FLOPs as the hardware cost metric is not efficient

	Real-time measurements	Lookup Table	MLP	XGBoost	Analytical Estimation
RMSE	0	4.32	3.83	3.8	5.6
Search Time (s)	172200	74869	31713	30531	33123
Search Time Speedups	1	2.3	5.43	5.64	5.2

TABLE V  
COMPARISON OF HARDWARE COST MEASUREMENT METHODS. THE HARDWARE COST IS MEASURED ON TESLA K80 GPU.

and may return suboptimal models. On the other hand, using the model size represented by the number of parameters allows to reduce the memory footprint and tends to be considered as an automatic method to search for compressed models [91].

*b) Latency:* Searching for low-latency architectures at inference time is crucial for real world application such as autonomous driving and traffic control. Moreover, resource-limited devices have latency constraints. Thus, a lot of works consider the latency in their objective function and search for the trade-off between inference time and accuracy.

*c) Energy Consumption:* Energy is usually profiled by the provided hardware platform profilers such as nvprof by NVIDIA. The energy can be formalized either as the peak power consumption or the average power, both metrics are used by different HW-NAS works including [46], [164], [165].

*d) Area:* Another metric that interest chip manufacturers is its area. The goal is to get the smallest chip possible that could run the best model. [129] uses MAESTRO [166] to explore the area and power consumption and search for the best model and best ASICs templates within a set of pre-defined ones. The area of the circuit is also a good indicator of the static power consumption. These two values are correlated.

*e) Memory footprint:* We can get the model size by calculating the number of parameters it needs to learn but a more efficient way is to profile how much memory it uses while running; this is the memory footprint. Having a low memory footprint is important for edge devices. These devices are not able to run models with high memory footprint. To reduce the memory footprint in edges devices, techniques like those presented in Section IX are applied.

Table VI summarizes all the methods used to estimate the latency, energy consumption in different NAS methods.

Real-world measurements provides a high accuracy in measuring the hardware efficiency of an architecture. MnasNet [47] uses this method in the exploration. It achieves 75.2% top-1 accuracy with 78ms latency on a Pixel phone platform, which is 1.8x faster than MobileNetV2 with 0.5% higher accuracy. However, this method considerably slows down the search algorithm by averaging hundreds of runs to get precise measurements. Additionally, this strategy is not scalable and requires that all the hardware platforms are available. This solution could be costly and needs a lot of mobile devices and software engineering work. That's why many works tend to use a prediction model [53], [72], [118], [119], [26], [133]



Method	How the method is achieved ?	Hardware Cost Metric	References
Real-time measurements	The sampled model is executed on the hardware target while searching.	Latency	MNASNet [47] NetAdapt [164] Z. Guo et al. [57] MCUNet [123]
		Energy	NetAdapt [164] MONAS [46] C. Gong et al. [165]
Lookup Table Models	A lookup table is created beforehand and filled with each operator latency on the targeted hardware. Once the search starts, the system will calculate the overall cost from the lookup table.	Latency	FBNet [23] HotNAS [122]
Analytical Estimation	Compute a rough estimate using the processing time, the stall time, and the starting time.	Latency	FNAS [26] NASCaps [100] A. Anderson et al. [167] Q. Lu et al. [119]
		Energy	NASCaps [100]
		Memory footprint	NASCaps [100]
		Area	NASAIC [129]
Prediction Model	Build a ML model to predict the cost using architecture and dataset features.	Latency	proxylessNAS [53] NA-SAIC [129] NeuNets [55] LEMONADE [46]

TABLE VI  
SUMMARY OF HARDWARE COST ESTIMATION METHODS

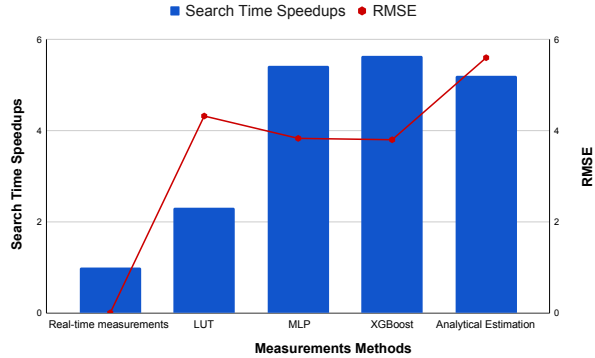


Fig. 14. Comparison of hardware cost measurement methods. LUT stands for Look Up Table. The speedups are calculated w.r.t the real-time measurements. The exact statistics are displayed in table V.

or a pre-collected lookup table [23], [120], [122] or computing an analytical estimation [26], [100].

In ProxylessNAS [53], the authors have developed three latency prediction models for three different platforms: a mobile phone (Google Pixel 1), a GPU (NVIDIA V100) and a CPU (Intel Xeon). To build their mobile latency predictors, they use the type of operators, input and output feature map sizes and other architectural hyperparameters as features. The real values for Pixel 1 phone have been measured with Tensorflow-Lite as software. On ImageNet, their model achieves 3.1% better top-1 accuracy than MobileNetV2, while being 1.2x faster with measured GPU latency.

In NASCaps [100], the functional behavior of a given specialized CNN and CapsNet hardware accelerator is modeled at a high level, to quickly estimate the memory usage, energy consumption, and latency. The HW platform, the ASIC CapsAcc in the paper, is described at the RTL-level. Using a VLSI CAD tool and the RTL specifications, energy, memory, latency costs of elementary operations are measured.

Elementary operations' cost correspond for example to the number of cycles and energy required to execute a layer. These values are then multiplied by the number of occurrences of each operation in the architecture and accumulated to obtain the total cost for the latency and energy.

Although these techniques are efficient, they require hardware experts to build the models. For the lookup table method, for example, the researcher needs to dedicate a lot of time to optimize the code of each operator/architecture in the targeted hardware, which requires compilation knowledge. Similarly, to build the best model predictor, the researcher needs expert knowledge to select the best features and verify the results. Therefore, these methods impose a barrier to non-hardware experts.

In order to fairly compare the accuracy of each method, we computed the latency of each architecture in NAS-Bench-101 [168] and compare it to the real-time measurements, according to three collecting methods: lookup table, prediction model and analytical estimation. For the lookup table, we calculated the latency of each operator used in the cell of the benchmark including Identity, Conv3x3BnRelu, Conv1x1BnRelu, MaxPool3x3, BottleneckConv3x3, BottleneckConv5x5, and MaxPool3x3Conv1x1. When a cell is generated, we sum the latency of the constructing operators and we get the latency of the whole cell. For the prediction model, we used two different models: a simple MLP and XGBoost, both trained on the real-time measurements of a portion of the benchmark (training set). We choose these two methods because they are both used by popular HW-NAS in [53] and [23] respectively. Lastly, for the analytical estimation we computed the number of MAC for the cell and multiply that by the latency of one multiply-add tensor instruction.

This experiment was run on a Tesla K80 GPU, the prediction model MLP had to run for 50 epochs with early stopping. NAS-Bench-101 defines more than 400k cells, we have done our test on 165,580 cells. The search algorithm used to calculate the search time is an evolutionary algorithm based

on the validation accuracy given by the benchmark and the latency measured by the different methods. Table V presents the results of the accuracy values of different methods. As expected, the analytical estimation does not produce good results compared to the prediction models or the lookup table method. The prediction models even with a simple XGBoost give the best results and accelerate the search more than 5 times compared to the real-time measurements.

## IX. OTHER CONSIDERATIONS FOR HARDWARE-AWARE NAS

Prior and Parallel to the hardware-aware NAS efforts, researchers have been working on reducing the memory footprint of deep learning models and execution time to facilitate the efficient deployment and design hardware-friendly models. Two main methods have been used, namely handcrafting new operators that are more efficient such as separable convolutions [29], grouped convolutions [4], and applying deep learning optimizations such as quantization [169] and pruning [170]. This latter method is automated by several NAS works to compress the model and make it possible to execute on different hardware accelerators. Moreover, reducing the number of learnable parameters makes the training faster. For instance, SAL [171] reduces the number of parameters of ResNet-56 from 1.22M to 0.36M without a big degradation of the model's accuracy; 0.6% decrease.

In this section, we are going to define the two most used deep learning optimizations and review the NAS works that focus on searching the right optimization parameters.

### A. Automatic Mixed-Precision Quantization

In the deep learning compression field, quantization is one of the most important methods. Starting with BinaryConnect [151] in 2015 to binarize CNNs weights, it is now implemented in many deep learning software such as PyTorch or Tensorflow. The idea is that the weights of the deep learning model do not have to be represented in full 32-bits precision and can be represented in 8-bits precision, or even binary precision in some cases, without significantly decreasing the model's accuracy. This idea was extended to the activations and weights [172]. Recently, mixed-precision quantization that applies different bitwidth values for different layers in the same network is more commonly used.

**HAQ** [62] implemented a dedicated reinforcement learning agent that learns to assign the right bitwidth to each layer. Their goal is to specialize the architecture to a specific hardware platform by incorporating hardware constraints and accuracy into the reward function. For each layer, the agent takes two decisions: one for the weights and another one for the activation.

**Single-Path** [57] proposes an evolutionary algorithm that searches for the mixed-precision quantization policy. But the search costs a huge amount of time and processing as the space for mixed-precision is enormous.

**BP-NAS** [173] cast the mixed-precision quantization problem as a constrained optimization. In addition, it proposes a new constraint that encourages the model to focus on

valid architectures while imposing a large punishment to the quantization outside the valid domain.

Although, these works present interesting solutions to the mixed-precision problem, the huge search space and computational cost of the learning process still represent a real challenge. Furthermore, mixed-precision representation requires specific MAC architectures with scalable-precision. This places a cap on the power efficiency, according to this study [174]. With this in mind, [175] proposes a uniform quantization search algorithm called neural channel expansion (NCE). NCE expands the number of channel of a layer when it is more sensitive to the quantization error while maintaining the same precision level.

### B. Automatic Pruning

The second prominent compression method is pruning. Pruning methods eliminate some neurons or connections according to a defined criterion to reduce the number of parameters in the architecture. Generally, we evaluate the importance of the neurons, we prune the least important ones, and finally, we fine-tune the network.

**AMC** [62] proposes an automatic model compression that looks for the optimal sparsity for each layer during pruning. The authors trained a reinforcement learning agent to predict the best sparsity for given hardware. The reward function includes accuracy and FLOPs after pruning the architecture.

X. Dong and Y. Yang in [176] propose to prune the over-parameterized network without performance damage. They directly search within their NAS process for a network with a flexible channel and layer sizes. **ABCPruner** [177] uses artificial bee colony algorithm to efficiently find the optimal pruned structure. Another worth mentioning paper is **Partial Order Pruning** [178] which proposes a hardware-aware NAS that prunes the search space with a partial order assumption to look for the best speed and accuracy trade-off.

### C. Security and Reliability Considerations

Other NAS methods try to address safety-critical issues by discovering architectures that are robust against adversarial attacks [179], [180]. RAS [179] formulates the robustness as the sum of the accuracies on a bunch of adversarial samples. This robust evaluation makes it easier for the evolutionary algorithm to select better architecture in the population and apply different mutations (e.g., add a block, remove a block, add a connection...). Along with [181], they reveal these observations: first, the more dense the architecture is the more robust it is, second, under computational budget, adding convolution operations to direct connection edge is effective, and finally, flow of solution procedure (FSP) matrix is a good indicator of network robustness. Please note here, that none of the NAS methods that consider security and reliability in the search place are hardware-aware.

## X. INDUSTRIAL ADOPTION OF NAS

One salient aspect that is not discussed in most AutoML surveys is its applications in the industrial domain. In this

section, we compare different tools that apply NAS principles in order to let users build a specific model for their datasets.

- 1) **Auto-Keras** is a framework that is based on the deep learning framework *keras*. By far, auto-keras is the most used autoML tool with features allowing to explore a wide range of neural network types (CNN, RNN and MLP). The search runs in parallel on CPUs and GPUs, with an adaptive search strategy for different GPU memory limits. Auto-keras builds its search space by applying a set of network morphism operations which keep the functionality of a neural network while changing its neural architecture.
- 2) **Microsoft NNI** (Neural Network Intelligence) is a general AutoML toolkit to help users tune their machine learning models (e.g., hyperparameters), design neural network architecture. One property it focuses on is the efficiency of their automation. For example, it leverages early feedback to speedup the tuning procedure. The proposed NAS framework helps the user define their own super network. For example, one can specify multiple operators for one single layer including depthwise convolution, dilated convolution, maxpooling, and the NNI will automatically find the best candidate based on a gradient-based optimization. On the other hand, the toolkit is also customizable for NAS researchers and enable them to modify the search algorithm and/or objective function.
- 3) **Google AutoML** The idea with Google AutoML is to create a meta-model capable of learning a way to design, generate and propose architectures for a given dataset. The model will eventually achieve an accuracy on the data set that will be used as a reward to guide the reinforcement learning agent to explore and find the best architecture from the search space. Google AutoML generalize the use of machine learning to everyone by offering an easy-to-use toolkit to upload the data set and let the system find the appropriate model. The AutoML toolkit includes AutoML Vision that manipulates image data set and tasks and AutoML NLP that processes text datasets and focuses on the translation task.
- 4) **IBM NeuNets** NeuNetS runs on the IBM cloud environment. The framework is designed to minimize human interaction and automate the whole machine learning pipeline, from the pre-processing methods to the algorithm search and model deployment. The process involves three components TAPAS, NCEvolve and HDMS. TAPAS synthesizes the neural network by offering a smart accuracy predictor based on the dataset characterizations and architecture features. NCEvolve searches for top-performing networks by minimizing the training time and resource needs. HDMS uses reinforcement learning to find and tune the architecture tailored for less common datasets. Finally, NeuNetS framework applies an evolutionary algorithm to tune some architectural hyperparameters such as the convolution filters.
- 5) **TPOT** [182] (Tree-based Pipeline Optimization Tool) is designed as a Python package that optimizes the machine learning pipeline using genetic programming. However, as it is based on scikit-learn framework, the package only allows classical machine learning algorithms such as random forest or decision trees and small MLPs and focuses only on tabular data sets. The package is still under development and may include more data types in the future.
- 6) **Microsoft Archai** Archai is a platform designed exclusively for Neural Architecture Search. It allows NAS researchers to easily mix and match between different techniques while ensuring reproducibility and fair comparison between the search algorithms and speedup techniques. The framework also considers multi-objective search optimizations by providing reports on model statistics such as number of parameters, flops, inference latency and model memory usage.
- 7) **dec.ai AutoNAC** dec.ai is a startup that aims to accelerate deep neural network inference on any hardware while preserving accuracy. AutoNAC is one of their commercial tools. As inputs, the framework needs a training and testing dataset, an already trained model and access to the hardware platform over which the model should be deployed. AutoNAC will then automatically compute a new low-latency model that preserves the accuracy. The idea behind AutoNAC is to construct small and fast models such that each model specializes in a small portion of the dataset, the basic strategy uses divide-and-conquer paradigm, where a complex problem is decomposed into a set of simpler sub-problems. For example, rather than building one complex model that classifies the data into the three classes defined by the data set, three simpler models that separate the classes two-by-two are built.
- 8) **Darwin** is a cloud-based service that utilizes genetic programming to iteratively select the best architectural hyperparameters and pass them to the next generation. Given a target dataset, Darwin tries to find a previously trained dataset that shares the same characteristics and initiate the search with the best model for that dataset. The framework also includes automation for other machine learning components such as the feature extraction by automatically detecting the data types and data pre-processing.

Each of these tools have their own advantages and disadvantages as well as different spaces to compete with. Microsoft NNI, Archai are designed to help NAS researchers compare, reproduce and experiment with multiple NAS search strategies and speedup methods. AutoKeras and TPOT are designed for the general machine learning researcher or engineer. Although they restrict their search space to a small set of architectural hyperparameters, the performance optimizations of their search are really attractive as an exchange for accuracy in this settings. In addition, both are open source which allow the user to run locally their search for privacy reasons. On the other hand, IBM Neunets, Google AutoML and dec.ai AutoNAC try to democratize machine learning and make it available and accessible to everyone regardless of their skill level and

Tool	Data Type			Model Type				Search Strategy					Speed up Methods			Compression Methods			Training HPO	Cloud-Based	Hardware-aware NAS	Training Framework
	Tabular	Images	Sequence	Classical ML	MLP	CNN	RNN	Transformers	Reinforcement Learning	Evolutionary Algorithm	Gradient-based	Bayesian Optimization	Hybrid	Limit Search Time	Hot Start	Proxy	Quantization	Pruning				
Auto-Keras	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	Keras
Microsoft NNI	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	Multiple
IBM NeuNets	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	No training: TAPAS
Google AutoML (+)	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	Tensorflow
TPOT (*)	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	Scikit-Learn
Microsoft Archai (*)	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	PyTorch
deci.ai AutoNAC (+)	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	Multiple
Darwin (+)	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	Tensorflow

Fig. 15. Comparison table of functionalities for NAS tools. (\*) still under development, (+) commercialized tools.

without writing any piece of code. It is a more ambitious goal which is the reason why they are using more limited architecture search spaces. However, their search algorithms consume a lot of computing power and time that’s why they are usually offered within a cloud service. The majority of the surveyed industrial NAS frameworks are not hardware-aware except for AutoNAC which given a defined model and its accuracy tries to find a faster architecture that doesn’t reduce the specified precision. Nonetheless, Microsoft NNI and Google AutoML apply some compression techniques to the found architecture to specialize it in the target hardware.

## XI. CHALLENGES AND LIMITATIONS

In this section, we lay down the main challenges and barriers that prevent us to unfold hardware-aware NAS’s full potential. Most of the challenges are also applicable to general NAS methods.

### A. Benchmarking and Reproducibility

An important challenge while working on NAS and HW-NAS is the reproducibility of the methods, which allows comparisons and concrete improvements. Unfortunately, due to the use of different search spaces, various training methods, and the required significant computational resources, reproducibility is a difficult step. This difficulty is even higher when it comes to HW-NAS considering the numerous possibilities of targeted hardware devices. To address this issue, many works [168], [183], [184], [139] have proposed different benchmarks and data sets that allow NAS researchers to:

- Eliminate the cost of generating a search space by querying directly a tabular data set.
- Evaluate different search strategies on the same search space which allows a fair comparison between them.
- Provide data sets to be used by accuracy predictor models and hardware cost models.
- Open HW-NAS research to non-hardware experts by proposing datasets that contain the hardware-related metrics. These metrics are usually obtained after optimization of the operators and software that require specific hardware knowledge.

In this section, we will review each NAS benchmark and highlight their strengths and weaknesses.

**NAS-Bench-101** [168]<sup>1</sup> is a tabular dataset that maps 432k unique architectures to their relative training accuracy, validation accuracy, testing accuracy as well as training time and the number of trained parameters. Each architecture is trained for various numbers of epochs 4, 12, 36, 108. The architectures follow a fixed macro architecture, with searchable cells stacked (See section V-A). Each cell is constructed with up to 7 layers that can include 3 types of operations: 3x3 conv, 1x1 conv, and 3x3 max-pooling. NAS-Bench-101 allows flexibility by allowing different layers to be used in the stacked cells. The search space is constrained by limiting the number of edges to 9 and the number of nodes to 7 including the input and output nodes for each cell. However, not including operations such as separable convolution or dilated convolutions makes the resulting models parameter-heavy, which is not hardware-friendly. Another downside of this benchmark is the inability to use one-shot optimizer NAS, this is due to the tabular format. To enable the evaluation of weight-sharing NAS methods, two benchmarks have been released NAS-Bench-201 [183]<sup>2</sup> and NAS-Bench-1Shot1 [184]<sup>3</sup>.

**NAS-Bench-201** represents 15,625 architectures using a fixed cell-based macro architecture. Similar to NAS-Bench-101, it uses a predefined set of operations including conv 3x3, conv 1x1, Avg pooling, skip connection and no operation label. Each architecture is trained on three different datasets with different complexities namely, CIFAR-10, CIFAR-100 and imagenet-16. The authors extended this benchmark and presented a year later NATS-Bench. **NATS-Bench** [185] increased the search space size by varying the sizes of their architectures.

Similarly, **NAS-Bench-1Shot1** presented a new reformulation to reuse the even much more extensive computation of the NAS-Bench-101 dataset (120 TPU years) to create three new one-shot search spaces with growing complexity

<sup>1</sup>Open sourced at <https://github.com/google-research/nasbench>

<sup>2</sup>Open sourced at <https://github.com/D-X-Y/AutoDL-Projects>

<sup>3</sup>Open sourced at <https://github.com/automl/nasbench-1shot1>

		NAS-Bench-101 [168]	NAS-Bench-201 [183]	NATS-Bench [185]	NAS-Bench-1shot1 [184]	NAS-Bench-NLP [139]	NAS-Bench-301[186] (DARTS)	HW-NAS-Bench [187]
Arch Search Space		Cell-based CNN	Cell-based CNN	Cell-based CNN	Super Network CNN	Cell-based LSTM	Super Network CNN	Cell-based CNN (NAS-Bench-201 + FBNet)
Size (number of architectures)		423k	15,625	15,625	363,648	14k	$10^{18}$	$46875 + 10^{21}$
Datasets	CIFAR-10	✓	✓	✓	✓		✓	✓
	CIFAR-100	✓	✓	✓	✓			✓
	ImageNet			✓	✓			✓
	PTB					✓		
	WikiText2					✓		
Metrics	Validation Accuracy	✓	✓	✓	✓	✓	✓	✓
	Training Time	✓	✓	✓	✓	✓	✓	✓
	Trained Parameters	✓	✓	✓		✓		✓
	FLOPs		✓	✓				✓
	Test Accuracy	✓	✓	✓	✓	✓	✓	✓
	Latency		✓	✓	✓		✓	✓
Energy								✓
Predictor Models							✓	✓
Hw Platforms	GPU	GTX 1080Ti	GTX 1080Ti	Not Mentioned	Not mentioned	Tesla V100-SXM2	Not Mentioned	Edge GPU Jetson TX2
	Edge TPU							Edge TPU Dev Board
	Smartphone							Pixel 3
	Raspberry Pi							Raspi 4
	FPGA							Xilinx ZC706
	ASICs							ASIC-Eyeriss

TABLE VII  
COMPARISON OF NAS BENCHMARKS.

containing 6240, 29160, and 363648 architectures. In order to use the expensive experimentation done by NAS-Bench-101, the authors created a one-shot architecture that contains all the discrete cell architectures defined by NAS-Bench-101. This allows the search algorithm to only train the supernet and get the performance of each path from the NAS-Bench-101. Therefore, their benchmark construction does not need any additional cost.

Alternatively, **NAS-Bench-NLP** defines a tabular benchmark for NLP tasks. Their cell-based search space is constructed based on an LSTM Macro-architecture borrowed from AWD-LSTM [188]. Each cell can contain up to 24 nodes with 3 hidden states and 3 linear input vectors. With these constraints and a set of operators composed by the most used activations in recurrent cells, they are able to build LSTM, GRU cells and many more. The architectures are trained on Penn Tree Bank (PTB) dataset [189] and a sub-sample of the best performing ones are also trained on WikiText-2 [190], which is a more realistic dataset for real-world NLP problems.

A major disadvantage of these benchmarks is the size and diversity of their search spaces. Indeed, as presented by the experimentation results we obtained in figure 16, on the small dataset NAS-Bench-201, a local search, which is the simplest optimization strategy, achieves state of the art results without a significant search time compared to other search strategies, except NAS without training [191]. In this experiment, we compare the different results obtained by different NAS approaches on the NAS-Bench-201. Most of them use the metrics provided by the benchmark along with fine-tuning the architecture to obtain a more accurate validation metric, except NAS without training [191]. The NAS

search without any further training achieves decent results within 17sec of search. By dividing the benchmarks into  $N$  mini-batches, they increase their training efficacy. The higher this number is ( $N$ ), the higher the over-fitting probability on the benchmark. Therefore, using small datasets with complex search algorithms does not yield any good results in terms of accuracy of the final architecture or efficiency of the search.

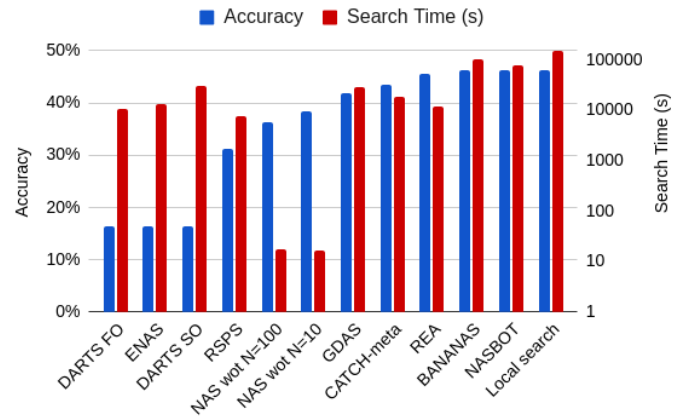


Fig. 16. Results of different search algorithms on NAS-Bench-201.

**NAS-Bench-301** is a much bigger benchmark that was designed to overcome the over-fitting problem on the architectures. It is based on DARTS [150] search space. Since DARTS search space is far too large to be exhaustively evaluated by real-world measurements, the authors built a surrogate model capable of predicting the various performance metrics. This model is trained on a subsample of the benchmark,

60k architectures whose performance have been measured on CPUs.

A more recent paper introduced the first hardware-aware NAS benchmark, **HW-NAS-Bench** [187]. This work extends the number of hardware metrics and records the latency, the energy consumed on 6 hardware devices including commercial edge devices, FPGA, and ASIC. Its search space is a combination of FBNet [23] search space and NAS-Bench-201.

Although these datasets provide a good start to test different search strategy, they lack a lot of important operators that can significantly change the resulting architecture for hardware-aware NAS. In addition, we note a growing number of different applications of NAS in various tasks such as image restoration [192], [193], [194], [195], semantic segmentation [130], [64], and medical segmentation [138], [196]. Therefore, there is a growing need for proper benchmarks for each of the diverse tasks.

### B. Transferability of the AI Models

A concept that is quite popular in deep learning is transferability. We train the model on a proxy dataset; that is usually smaller and simpler, then we fine-tune the trained weights on the targeted dataset to get a more efficient model. For example, we can train our image classifier on CIFAR-10 [147], which contains 10 classes of 32x32 pixels images, and then fine-tune our model on ImageNet [197], which is made of 14 million images and 20,000 classes. Another benefit of transferability is when the target dataset is not large enough to train a model on.

To enhance transferability, previous NAS works used cell-based search spaces. To transfer a model obtained from a cell-based search spaces, we just need to adjust the input sizes of the cells and deepen the network by adding more cells. However, stacking the same blocks seems to be not efficient when incorporating hardware constraints. As MNASNet [47] argued, restricting the cell diversity is critical for achieving high accuracy and low latency on mobile settings. Therefore, MNASNet uses a hierarchical search space that diversify the cells in the architecture as well as the operators within that cell.

Many NAS works [160], [149] have included dedicated evaluations of the transferability of their final model. PNAS [160] proved that CIFAR-10 classification is highly correlated to ImageNet classification, hence transferring a cell or an architecture from CIFAR-10 to ImageNet would produce good results. XNAS [149] transferred their final cell structure on 6 popular classification benchmarks surpassing other conventional NAS methods while taking account of the hardware constraints.

Another interesting concept is presented by NAT [198]. They leverage NAS process to directly find transferable weights (i.e., get rid of the fine-tuning stage). The key idea behind NAT is that they start from a supernet and adaptively modify it to obtain a task-specific super network. This latter approach can then be used directly to search for architectures within one task without the additional training cost. Under mobile scenarios, they demonstrated the efficacy of NAT on 11 benchmarks including ImageNet.

Overall, taking the transferability of the model into account remains a difficult challenge. Most solutions modify the search space to fit a certain task by leveraging a task-specific macro-architecture. Hardware-awareness adds another level to this challenge as the architecture needs to be flexible to adapt to multiple platforms. Note that, transfer learning generalizes the ability to use a model from one task to another. In this section, we only talked about NAS that transfers their models from one dataset to another within the same task, which is a sub definition of transfer learning. Unfortunately, as NAS tasks evolve around image classification mainly, the whole search needs to be executed all over again if we want an architecture for another task.

### C. Transferability of the Hardware-aware NAS across multiple Platforms

HW-NAS suffers from the conditional optimality due to the variety of existing devices. Ideally, we should design different architectures for different platforms. However in real situations, given the prohibitive cost of the search and the cost of training on multiple architectures, we often resort to designing one architecture and deploying it anywhere. Transferring a model from one platform to another or being able to produce hardware transferable models via the NAS process is an interesting challenge for HW-NAS. The main difficulty lies in the variety and complexity of the existing platforms. For example, different platforms might perform well on different types of convolutions. In the following section, we discuss two approaches. The first transfers a single-target HW-NAS to another target by modifying the measurement values. The second takes the final architecture of the NAS process and transforms it to fit another platform. Each approach has its pros and cons as discussed below.

*a) Transfer the entire NAS process:* this means that the whole NAS process needs to be re-executed to suit the new targeted hardware. In a single target HW-NAS, it is quite a challenge to transfer the entire process to another platform. Without mentioning the huge computational cost of retraining the entire process, the collection of the hardware constraints can be costly as well. When using real-world measurements, [62] ran the NAS search for three hardware platforms: CPU (Xeon E5-2640 v4), GPU (Tesla V100) and mobile phone (Google Pixel-1). However, using real-world measurements considerably slows down the search algorithm and requires the availability of the targeted hardware during the search time. On the other hand, using an analytical estimation requires expert knowledge for the different targeted platforms. When using other collection methods such as the lookup table or the prediction model, we'll need to collect data from the new platforms by running again the entire set of operators. To this end, HW-NAS tries to create a general measurement method. For example, Once-for-all [18] created a lookup table with the reported inference latency on each tested hardware platform (i.e., Samsung S7 Edge, Note8, Note10, Google Pixel1, Pixel2, LG G8, NVIDIA 1080Ti, V100 GPUs, Jetson TX2, Intel Xeon CPU, Xilinx ZU9EG, and ZU3EGFPGAs). According to the used measurement method, transferring the



NAS process to target another platform is increasingly difficult and not scalable.

*b) Transfer the final model:* An alternative approach is to find the best model for one hardware platform and then try to specialize it for another one. This was done by [18], [53], [149]. Usually, the specialization is done by compressing the model using quantization which enables the model to fit in tiny devices. However, the specialization is challenging because of the following reasons:

- **An operator may be efficient for one platform and less efficient in another:** In [127], the authors argued that separable convolutions give great results when ran on GPUs but perform badly on CPUs. Additionally, it is common that deeper networks perform well on CPUs while wider ones perform well on GPUs because of the possible parallelization. Table VIII demonstrates the comparison between the execution times of different operators on an Intel i7 CPU and NVIDIA TX2 GPU. The results reinforce our assumption that different operators's efficiency vary from one platform to another. Therefore, the best model is highly correlated to the platform we choose which makes design a HW-NAS that targets multiple platforms very challenging.
- **Limits of the compression methods:** We consider here the quantization and pruning. For these two methods, we know that theoretically the compression ratio have a threshold that cannot be surpassed. For example, quantizing a model implies encoding its activations and weights into the minimum possible bit length. Theoretically this length is 2 (one bit). Even without considering the trade-off between accuracy and model size, these methods have limits. This is why we need to start from a model that is already optimized to be able to deploy the model on tiny devices. For this end, [62] is composed of three components. The first components is a multi-objective NAS that searches for the best model in terms of accuracy and latency. Given the resultant model of the first component, the second component searches for pruning possibilities that would preserve the accuracy and decrease the model size. Finally, the last component takes care of quantizing the model with a mixed precision. [199] starts from a standard architectures such as VGG, ResNet and GoogleNet and casts the quantization as a neural architecture search problem. This work achieves a minimal loss of accuracy with appreciable memory savings. In addition to the limit for the model size, we also have a limit for the accuracy. The compressed model typically does not have a better accuracy than the pre-trained model we started from.

#### D. Outlook and Future Directions

In addition to the benchmarking, reproducibility and transferability challenges presented in the previous three subsections, HW-NAS suffers from the same limitations as NAS. The main one is the cost of the search algorithm. The current popular way to speed up the process is to use differentiable

Operator	Avg Latency on CPU (ms)	Avg Latency on GPU (ms)	Accuracy
Conv2d	1.33	0.85	0.67
Separable Depth-wise Convolution	2.05	0.54	0.64
Dilated convolution	1.36	0.835	0.56
Grouped Convolution	2.27	1.94	0.62
LSTM cell	14.93	2.45	0.57
GRU cell	9.32	2.53	0.65

TABLE VIII

COMPARISON BETWEEN DIFFERENT OPERATORS ON INTEL I7 CPU AND NVIDIA TX2 GPU. THE CONVOLUTION OPERATORS WERE USED TO CREATE A CNN MODEL THAT WAS TRAINED FOR IMAGE CLASSIFICATION ON IMAGENET. THE RNN CELLS WERE TRAINED ON TEXT CLASSIFICATION ON IMDB DATASET [200]. Results were obtained using PyTorch with number of samples 1000.

neural architecture search and create a supernetwork representing the whole search space. This supernetwork is trained once to get the weights of all the sub-networks, thus avoiding training each sampled architecture. This technique reduces the search time from several days to hours. However, a major disadvantage of this network is its restriction on the targeted task and domain. More research is needed to explore efficient ways and tricks to make HW-NAS more practical and more amenable to a diverse set of tasks and domains. This would especially be useful in commercial settings.

Given the increasingly fragmented hardware landscape, HW-NAS should investigate the possibilities to reduce the conditional optimality problem. Most existing NAS approaches design once and deploy to all. We have seen this clearly through the industrial adoption of NAS techniques which lacks hardware-awareness. Ideally, we should be able to design different architectures for different devices as no-one-model fits all. This problem is significantly exacerbated by the increasing hardware heterogeneity in AI-powered devices.

There is no doubt that the future of mobile and handheld devices is AI. AI-focused mobile chips from top manufacturing companies like Apple, Samsung, Huawei and others make their ways into the mainstream. These devices use SoCs that take advantage of multiple platforms (e.g., GPUs, CPUs and NPU in the same chip). However, this heterogeneity of platforms needs to be well understood to speed up the inference time of deep neural networks [201]. This motivates further the research community to simultaneously explore both the architecture search space and the hardware design space to identify the best neural architecture and hardware pairs that maximize both test accuracy and hardware efficiency. Such co-exploration will be critical to allow designing architectures that can be deployed efficiently on a variety of platforms: data center, edge, mobile, and embedded.

HW-NAS also suffers from a lack of benchmarks and public datasets that can help accelerate HW-NAS research and make it reproducible and accessible. HW-NAS-Bench [187] is a recent first attempt towards this direction. It tries to offer computed metrics on multiple commercial platforms including edge GPUs, FPGA and ASICs and make these measurements available to the research community. One of their main goals

is to democratize HW-NAS research to non-hardware experts.

More work is also needed to further co-explore neural architectures with quantization, pruning, and hardware design. Most of the existing efforts, such as mixed-precisions have been studied in the context of fixed architectures.

HW-NAS should also look at exploring neural architecture search with emerging computing paradigms such as in-memory-computing [202]. These new non-von-Neumann paradigms present novel solutions to AI computing based on emerging nano-devices called Phase-Change Memory (PCM) [203]. Existing HW-NAS solutions are dedicated to the conventional von-Neumann computing architecture, where the memory wall heavily limits performance. The optimization space spans multiple design points that range from device types, to circuit topologies, to device non-idealities and variations, to neural architectures.

While deep neural networks have clear commercial use cases, the next AI breakthrough may require an entirely different combination of algorithm, hardware and software. HW-NAS offers a paradigm that opens up the design space and pushes forward the Pareto frontier between hardware efficiency and model accuracy for efficient and improved hardware/software co-design, hence pushing AI to its next frontiers.

## XII. CONCLUSION

We have provided a comprehensive survey and systematic analysis of state-of-the-art hardware-aware neural architectural search. We reviewed several multi-objective strategies that aim to find the optimal architecture with the highest accuracy while reducing memory and computations costs. We provide a HW-NAS taxonomy and categorize existing approaches along four key dimensions: the search space, the search strategy, the acceleration technique, and the hardware cost estimation strategy. We also discussed other considerations in Hardware-aware NAS that include optimizations such as quantization and pruning. Finally, we presented a discussion on future directions that would benefit existing and future HW-NAS researchers. We believe that HW-NAS is an important direction that would enable hardware-software co-design and make deep learning solutions not only accessible to everyone but also sustainable.

Reference	Target Model	Arch Space	HSS	Strategy	Hw Cost Estimation Method	Target Hw	HWT	Speedup Technique	P	THPO	OS
[47]	Standard CNN	Factorized Hierarchical Layer-wise		Reinforcement learning	Real-world measurements	Pixel 1		early stopping	✓		✓
[53]	Extended CNN	Layer-wise		Gradient-based Optimization	ML Prediction Model	Pixel 1		Super Network		✓	✓
[23]	Extended CNN	Layer-wise		Gradient-based Optimization	Lookup Table Model	Samsung iPhone X	✓	Super Network	✓		✓
[123]	Standard CNN	Optimized hyperparameters Search space		Evolutionary Algorithm	Try-to-deploy	MCU		Super Network	✓		
[26]	Standard CNN	Hyper-parameters Search space	✓	Reinforcement learning	Analytical Estimation	FPGA	✓	Ignore not deployable models	✓	✓	
[129]	Extended CNN	Hyper-parameters	✓	Reinforcement learning	ML Prediction Model	ASIC	✓	Ignore not deployable models			
[100]	Capsule + Standard CNN	Layer-wise	✓	Evolutionary algorithm	Analytical Estimation	ASIC	✓	Early stopping	✓		✓
[127]	Standard CNN	Multi Hardware		Reinforcement Learning	ML Prediction Model	Multiple	✓	Super Network			
[122]	Extended CNN	Layer-wise	✓	reinforcement learning	Lookup Table Model	FPGA	✓	Hot Start			
[72]	Extended CNN	Optimized layer-wise		Two stage One-Shot Search	ML Prediction Model	DSP /VPU/CPU	✓	Super Network			
[78]	Extended CNN	Hierarchical Search Space	✓	Reinforcement Learning	Real-time measurements	FPGA		NOT Mentioned			
[120]	Standard CNN	NASBench[168]	✓	Reinforcement Learning	Lookup Table Model	FPGA		NOT Mentioned	✓		
[133]	Extended CNN	Optimized layer-wise	✓	Stochastic Coordinate Descent	ML Prediction Model	FPGA	✓	Accuracy prediction	✓		
[132]	Transformers	Encoder wise Search Space		Evolutionary Algorithm	ML Prediction Model	Raspberry-Pi/CPU/GPU	✓	Super Network			✓

TABLE IX

DETAILED OVERVIEW OF 15 MOST POPULAR HW-NAS.

HWT: HARDWARE TRANSFERABILITY, P: PROXY DATASETS,

THPO: TRAINING HYPERPARAMETER OPTIMIZATION, OS: OPEN SOURCE.

FOR A COMPLETE LIST OF ALL HW-NAS WORKS, PLEASE VISIT THE FOLLOWING LINK [HTTPS://TINYURL.COM/Y6458SKT](https://tinyurl.com/Y6458SKT)

## REFERENCES

- [1] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun, "Overfeat: Integrated recognition, localization and detection using convolutional networks," *arXiv preprint arXiv:1312.6229*, 2013.
- [2] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [3] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *European conference on computer vision*. Springer, 2014, pp. 818–833.
- [4] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [5] D. S. Park, Y. Zhang, Y. Jia, W. Han, C.-C. Chiu, B. Li, Y. Wu, and Q. V. Le, "Improved noisy student training for automatic speech recognition," *Interspeech 2020*, Oct 2020.
- [6] W. Xiong, J. Droppo, X. Huang, F. Seide, M. Seltzer, A. Stolcke, D. Yu, and G. Zweig, "The microsoft 2016 conversational speech recognition system," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017, pp. 5255–5259.
- [7] G. Hinton, L. Deng, D. Yu, G. E. Dahl, A.-r. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. N. Sainath *et al.*, "Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups," *IEEE Signal processing magazine*, vol. 29, no. 6, pp. 82–97, 2012.
- [8] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT*, J. Burstein, C. Doran, and T. Solorio, Eds., 2019, pp. 4171–4186.
- [9] R. Collobert, J. Weston, L. Bottou, M. Karlen, K. Kavukcuoglu, and P. Kuksa, "Natural language processing (almost) from scratch," *Journal of machine learning research*, vol. 12, pp. 2493–2537, 2011.
- [10] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey *et al.*, "Google's neural machine translation system: Bridging the gap between human and machine translation," *arXiv preprint arXiv:1609.08144*, 2016.
- [11] M. Tan and Q. Le, "EfficientNet: Rethinking model scaling for convolutional neural networks," in *Proceedings of the 36th International Conference on Machine Learning*, 2019, pp. 6105–6114.
- [12] H. Zhang, C. Wu, Z. Zhang, Y. Zhu, Z. Zhang, H. Lin, Y. Sun, T. He, J. Mueller, R. Manmatha *et al.*, "Resnest: Split-attention networks," *arXiv preprint arXiv:2004.08955*, 2020.
- [13] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016, <http://www.deeplearningbook.org>.
- [14] L. Deng, G. Li, S. Han, L. Shi, and Y. Xie, "Model compression and hardware acceleration for neural networks: A comprehensive survey," *Proceedings of the IEEE*, vol. 108, no. 4, pp. 485–532, 2020.
- [15] Anonymous, "An image is worth 16x16 words: Transformers for image recognition at scale," in *Submitted to International Conference on Learning Representations*, 2021, under review. [Online]. Available: <https://openreview.net/forum?id=YicbFdNTTy>
- [16] M. Sandler, A. G. Howard, M. Zhu, A. Zhmoginov, and L. Chen, "Inverted residuals and linear bottlenecks: Mobile networks for classification, detection and segmentation," *CoRR*, vol. abs/1801.04381, 2018. [Online]. Available: <http://arxiv.org/abs/1801.04381>
- [17] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, 2016, pp. 2818–2826. [Online]. Available: <https://doi.org/10.1109/CVPR.2016.308>
- [18] H. Cai, C. Gan, T. Wang, Z. Zhang, and S. Han, "Once-for-all: Train one network and specialize it for efficient deployment," in *8th International Conference on Learning Representations, ICLR*. OpenReview.net, 2020. [Online]. Available: <https://openreview.net/forum?id=HylxE1HKwS>
- [19] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le, "Learning transferable architectures for scalable image recognition," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8697–8710.
- [20] G. Ghiasi, T. Lin, R. Pang, and Q. V. Le, "NAS-FPN: learning scalable feature pyramid architecture for object detection," *CoRR*, vol. abs/1904.07392, 2019. [Online]. Available: <http://arxiv.org/abs/1904.07392>
- [21] A. Howard, M. Sandler, B. Chen, W. Wang, L. Chen, M. Tan, G. Chu, V. Vasudevan, Y. Zhu, R. Pang, H. Adam, and Q. Le, "Searching for mobilenetv3," in *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 1314–1324.
- [22] B. Zoph and Q. V. Le, "Neural architecture search with reinforcement learning," 2016. [Online]. Available: <http://arxiv.org/abs/1611.01578>
- [23] 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 *Conference on Computer Vision and Pattern Recognition, CVPR*. Computer Vision Foundation / IEEE, 2019, pp. 10 734–10 742.
- [24] T. Elsken, J. H. Metzen, F. Hutter *et al.*, "Neural architecture search: A survey," *Journal of Machine Learning Research*, vol. 20, no. 55, pp. 1–21, 2019.
- [25] M. Wistuba, A. Rawat, and T. Pedapati, "A survey on neural architecture search," *CoRR*, vol. abs/1905.01392, 2019. [Online]. Available: <http://arxiv.org/abs/1905.01392>
- [26] X. Zhang, W. Jiang, Y. Shi, and J. Hu, "When neural architecture search meets hardware implementation: from hardware awareness to co-design," in *IEEE Computer Society Annual Symposium on VLSI (ISVLSI)*, 2019, pp. 25–30.
- [27] Y. Cheng, D. Wang, P. Zhou, and T. Zhang, "A survey of model compression and acceleration for deep neural networks," *CoRR*, vol. abs/1710.09282, 2017. [Online]. Available: <http://arxiv.org/abs/1710.09282>
- [28] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," in *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2016. [Online]. Available: <http://arxiv.org/abs/1511.07122>
- [29] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [30] X. Zhu, P. Sobhani, and H. Guo, "Long short-term memory over tree structures," *CoRR*, vol. abs/1503.04881, 2015. [Online]. Available: <http://arxiv.org/abs/1503.04881>
- [31] J. van der Westhuizen and J. Lasenby, "The unreasonable effectiveness of the forget gate," *CoRR*, vol. abs/1804.04849, 2018. [Online]. Available: <http://arxiv.org/abs/1804.04849>
- [32] L. Song, M. Ishteva, A. P. Parikh, E. P. Xing, and H. Park, "Hierarchical tensor decomposition of latent tree graphical models," in *Proceedings of the 30th International Conference on Machine Learning, ICML*, vol. 28, 2013, pp. 334–342. [Online]. Available: <http://proceedings.mlr.press/v28/song13.html>
- [33] Y. Znyied, R. Boyer, A. L. F. de Almeida, and G. Favier, "A tt-based hierarchical framework for decomposing high-order tensors," *SIAM J. Sci. Comput.*, vol. 42, no. 2, pp. A822–A848, 2020.
- [34] S. Gupta, A. Agrawal, K. Gopalakrishnan, and P. Narayanan, "Deep learning with limited numerical precision," in *Proceedings of the 32nd International Conference on Machine Learning*, vol. 37, 2015, pp. 1737–1746.
- [35] S. Han, H. Mao, and W. J. Dally, "Deep compression: Compressing deep neural network with pruning, trained quantization and huffman coding," in *4th International Conference on Learning Representations, ICLR*, 2016. [Online]. Available: <http://arxiv.org/abs/1510.00149>
- [36] M. Zhu and S. Gupta, "To prune, or not to prune: exploring the efficacy of pruning for model compression," *CoRR*, vol. abs/1710.01878, 2017. [Online]. Available: <http://dblp.uni-trier.de/db/journals/corr/corr1710.html#abs-1710-01878>
- [37] M. Li, Y. Liu, X. Liu, Q. Sun, X. You, H. Yang, Z. Luan, L. Gan, G. Yang, and D. Qian, "The deep learning compiler: A comprehensive survey," *IEEE Transactions on Parallel and Distributed Systems*, vol. 32, no. 3, pp. 708–727, 2021.
- [38] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [39] B. Jang, M. Kim, G. Harerimana, and J. W. Kim, "Q-learning algorithms: A comprehensive classification and applications," *IEEE Access*, vol. 7, pp. 133 653–133 667, 2019.
- [40] R. J. Williams, "Simple statistical gradient-following algorithms for connectionist reinforcement learning," *Machine Learning*, vol. 8, no. 3–4, p. 229–256, 1992.
- [41] S. C. Smithson, G. Yang, W. J. Gross, and B. H. Meyer, "Neural networks designing neural networks: Multi-objective hyper-parameter optimization," *CoRR*, vol. abs/1611.02120, 2016. [Online]. Available: <http://arxiv.org/abs/1611.02120>
- [42] A. Gordon, E. Eban, O. Nachum, B. Chen, T. Yang, and E. Choi, "Morphnet: Fast & simple resource-constrained structure learning of deep networks," *CoRR*, vol. abs/1711.06798, 2017. [Online]. Available: <http://arxiv.org/abs/1711.06798>

- [43] J.-D. Dong, A.-C. Cheng, D.-C. Juan, W. Wei, and M. Sun, "Ppp-net: Platform-aware progressive search for pareto-optimal neural architectures," 2018. [Online]. Available: <https://openreview.net/forum?id=BINT3TAIM>
- [44] T. Elsken, J. H. Metzen, and F. Hutter, "Efficient multi-objective neural architecture search via lamarckian evolution," in *7th International Conference on Learning Representations, ICLR*. OpenReview.net, 2019. [Online]. Available: <https://openreview.net/forum?id=ByME42AqK7>
- [45] Y. Zhou, S. Ebrahimi, S. Ö. Arik, H. Yu, H. Liu, and G. Diamos, "Resource-efficient neural architect," *CoRR*, vol. abs/1806.07912, 2018. [Online]. Available: <http://arxiv.org/abs/1806.07912>
- [46] C. Hsu, S. Chang, D. Juan, J. Pan, Y. Chen, W. Wei, and S. Chang, "MONAS: multi-objective neural architecture search using reinforcement learning," *CoRR*, vol. abs/1806.10332, 2018. [Online]. Available: <http://arxiv.org/abs/1806.10332>
- [47] M. Tan, B. Chen, R. Pang, V. Vasudevan, and Q. V. Le, "Mnasnet: Platform-aware neural architecture search for mobile," *CoRR*, vol. abs/1807.11626, 2018. [Online]. Available: <http://arxiv.org/abs/1807.11626>
- [48] N. Mitschke, M. Heizmann, K. Noffz, and R. Wittmann, "Gradient based evolution to optimize the structure of convolutional neural networks," in *25th IEEE International Conference on Image Processing (ICIP)*, 2018, pp. 3438–3442.
- [49] V. Nekrasov, H. Chen, C. Shen, and I. D. Reid, "Fast neural architecture search of compact semantic segmentation models via auxiliary cells," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2019, pp. 9126–9135.
- [50] Z. Lu, I. Whalen, V. Boddeti, Y. Dhebar, K. Deb, E. Goodman, and W. Banzhaf, "Nsga-net: Neural architecture search using multi-objective genetic algorithm," in *Proceedings of the Genetic and Evolutionary Computation Conference*, 2019. [Online]. Available: <https://doi.org/10.1145/3321707.3321729>
- [51] A. Cheng, C. H. Lin, D. Juan, W. Wei, and M. Sun, "Instanas: Instance-aware neural architecture search," in *The Thirty-Fourth AAAI Conference on Artificial Intelligence*, 2020, pp. 3577–3584. [Online]. Available: <https://aaai.org/ojs/index.php/AAAI/article/view/5764>
- [52] L. Cai, A. Barneche, A. Herbout, C. S. Foo, J. Lin, V. R. Chandrasekhar, and M. M. S. Aly, "TEA-DNN: the quest for time-energy-accuracy co-optimized deep neural networks," in *IEEE/ACM International Symposium on Low Power Electronics and Design*, 2019, pp. 1–6. [Online]. Available: <https://doi.org/10.1109/ISLPED.2019.8824934>
- [53] H. Cai, L. Zhu, and S. Han, "Proxylessnas: Direct neural architecture search on target task and hardware," *CoRR*, vol. abs/1812.00332, 2018. [Online]. Available: <http://arxiv.org/abs/1812.00332>
- [54] X. Dai, P. Zhang, B. Wu, H. Yin, F. Sun, Y. Wang, M. Dukhan, Y. Hu, Y. Wu, Y. Jia, P. Vajda, M. Uyttendaele, and N. K. Jha, "Chamnet: Towards efficient network design through platform-aware model adaptation," in *Conference on Computer Vision and Pattern Recognition, CVPR*, 2019, pp. 11 398–11 407.
- [55] A. Sood, B. Elder, B. Herta, C. Xue, C. Bekas, A. C. I. Malossi, D. Saha, F. Scheidegger, G. Venkataraman, G. Thomas, G. Mariani, H. Strobelt, H. Samulowitz, M. Wistuba, M. Manica, M. R. Choudhury, R. Yan, R. Istrate, R. Puri, and T. Pedapati, "Neunets: An automated synthesis engine for neural network design," *CoRR*, vol. abs/1901.06261, 2019. [Online]. Available: <http://arxiv.org/abs/1901.06261>
- [56] X. Li, Y. Zhou, Z. Pan, and J. Feng, "Partial order pruning: For best speed/accuracy trade-off in neural architecture search," in *Conference on Computer Vision and Pattern Recognition, CVPR*. Computer Vision Foundation / IEEE, 2019, pp. 9145–9153.
- [57] Z. Guo, X. Zhang, H. Mu, W. Heng, Z. Liu, Y. Wei, and J. Sun, "Single path one-shot neural architecture search with uniform sampling," *arXiv preprint arXiv:1904.00420*, 2019.
- [58] J. Yu and T. S. Huang, "Network slimming by slimmable networks: Towards one-shot architecture search for channel numbers," *CoRR*, vol. abs/1903.11728, 2019. [Online]. Available: <http://arxiv.org/abs/1903.11728>
- [59] D. Stamoulis, R. Ding, D. Wang, D. Lymberopoulos, B. Priyantha, J. Liu, and D. Marculescu, "Single-path nas: Designing hardware-efficient convnets in less than 4 hours," in *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, 2019, pp. 481–497.
- [60] Y. Xiong, R. Mehta, and V. Singh, "Resource constrained neural network architecture search: Will a submodularity assumption help?" in *International Conference on Computer Vision, ICCV*, 2019, pp. 1901–1910.
- [61] M. Baldeon-Calisto and S. K. Lai-Yuen, "Adaresu-net: Multiobjective adaptive convolutional neural network for medical image segmentation," *Neurocomputing*, vol. 392, pp. 325 – 340, 2020.
- [62] S. Han, H. Cai, L. Zhu, J. Lin, K. Wang, Z. Liu, and Y. Lin, "Design Automation for Efficient Deep Learning Computing," *arXiv e-prints*, p. arXiv:1904.10616, Apr. 2019.
- [63] I. Fedorov, R. P. Adams, M. Mattina, and P. N. Whatmough, "Sparse: Sparse architecture search for cnns on resource-constrained microcontrollers," *CoRR*, vol. abs/1905.12107, 2019. [Online]. Available: <http://arxiv.org/abs/1905.12107>
- [64] A. Shaw, D. Hunter, F. Iandola, and S. Sidhu, "Squeezenas: Fast neural architecture search for faster semantic segmentation," 2019.
- [65] P. Liu, B. Wu, H. Ma, and M. Seok, "Memnas: Memory-efficient neural architecture search with grow-trim learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [66] Y. Guo, Y. Zheng, M. Tan, Q. Chen, J. Chen, P. Zhao, and J. Huang, "NAT: neural architecture transformer for accurate and compact architectures," in *Advances in Neural Information Processing Systems* 32, H. M. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. B. Fox, and R. Garnett, Eds., 2019, pp. 735–747. [Online]. Available: <http://papers.nips.cc/paper/8362-nat-neural-architecture-transformer-for-accurate-and-compact-architectures>
- [67] W. Bae, S. Lee, Y. Lee, B. Park, M. Chung, and K. Jung, "Resource optimized neural architecture search for 3d medical image segmentation," in *Medical Image Computing and Computer Assisted Intervention - MICCAI*, D. Shen, T. Liu, T. M. Peters, L. H. Staib, C. Essert, S. Zhou, P. Yap, and A. R. Khan, Eds., vol. 11765, 2019, pp. 228–236.
- [68] C. Schorn, T. Elsken, S. Vogel, A. Runge, A. Guntoro, and G. Ascheid, "Automated design of error-resilient and hardware-efficient deep neural networks," *CoRR*, vol. abs/1909.13844, 2019. [Online]. Available: <http://arxiv.org/abs/1909.13844>
- [69] D. Wang, M. Li, L. Wu, V. Chandra, and Q. Liu, "Energy-aware neural architecture optimization with fast splitting steepest descent," *CoRR*, vol. abs/1910.03103, 2019. [Online]. Available: <http://arxiv.org/abs/1910.03103>
- [70] T. Cassimon, S. Vanneste, S. Bosmans, S. Mercelis, and P. Hellinckx, "Using neural architecture search to optimize neural networks for embedded devices," in *Advances on P2P, Parallel, Grid, Cloud and Internet Computing*, L. Barolli, P. Hellinckx, and J. Natwichai, Eds., 2020.
- [71] B. Wang, Y. Sun, B. Xue, and M. Zhang, "Evolving deep neural networks by multi-objective particle swarm optimization for image classification," *CoRR*, vol. abs/1904.09035, 2019. [Online]. Available: <http://arxiv.org/abs/1904.09035>
- [72] L. Zhang, Y. Yang, Y. Jiang, W. Zhu, and Y. Liu, "Fast hardware-aware neural architecture search," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2020, pp. 2959–2967.
- [73] M. Loni, S. Sinaei, A. Zoljodi, M. Daneshmand, and M. Sjödin, "Deepmaker: A multi-objective optimization framework for deep neural networks in embedded systems," *Microprocessors and Microsystems*, vol. 73, p. 102989, 2020.
- [74] Y. Xu, L. Xie, X. Zhang, X. Chen, B. Shi, Q. Tian, and H. Xiong, "Latency-aware differentiable neural architecture search," 2020.
- [75] M. S. Iqbal, J. Su, L. Kotthoff, and P. Jamshidi, "Flexibo: Cost-aware multi-objective optimization of deep neural networks," 2020.
- [76] S. Bian, W. Jiang, Q. Lu, Y. Shi, and T. Sato, "NASS: optimizing secure inference via neural architecture search," in *ECAI - 24th European Conference on Artificial Intelligence*, G. D. Giacomo, A. Catalá, B. Dilkina, M. Milano, S. Barro, A. Bugarin, and J. Lang, Eds., vol. 325, 2020, pp. 1746–1753.
- [77] Z. Chen, F. Zhou, G. Trimponias, and Z. Li, "Multi-objective neural architecture search via non-stationary policy gradient," 2020.
- [78] Y. Yang, C. Wang, L. Gong, and X. Zhou, "Fpnet: Customized convolutional neural network for fpga platforms," in *International Conference on Field-Programmable Technology (ICFPT)*, 2019, pp. 399–402.
- [79] B. Lu, J. Yang, L. Y. Chen, and S. Ren, "Automating deep neural network model selection for edge inference," in *IEEE First International Conference on Cognitive Machine Intelligence (CogMI)*, 2019, pp. 184–193.
- [80] J. Fernandez-Marques, P. N. Whatmough, A. Mundy, and M. Mattina, "Searching for winograd-aware quantized networks," 2020.
- [81] S. Hong, D. Kim, and M. Choi, "Memory-efficient models for scene text recognition via neural architecture search," in *IEEE Winter Applications of Computer Vision Workshops (WACVW)*, 2020, pp. 183–191.

- [82] S. Gupta and B. Akin, "Accelerator-aware neural network design using autml," 2020.
- [83] S.-Y. Huang and W.-T. Chu, "Ponas: Progressive one-shot neural architecture search for very efficient deployment," *arXiv preprint arXiv:2003.05112*, 2020.
- [84] Y. Li, X. Jin, J. Mei, X. Lian, L. Yang, C. Xie, Q. Yu, Y. Zhou, S. Bai, and A. L. Yuille, "Neural architecture search for lightweight non-local networks," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 10 294–10 303.
- [85] Y. Xiong, H. Liu, S. Gupta, B. Akin, G. Bender, P.-J. Kindermans, M. Tan, V. Singh, and B. Chen, "MobileDets: Searching for Object Detection Architectures for Mobile Accelerators," *arXiv e-prints*, 2020.
- [86] Z. Yu, Y. Qin, X. Xu, C. Zhao, Z. Wang, Z. Lei, and G. Zhao, "Autofas: Searching lightweight networks for face anti-spoofing," in *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing*, 2020, pp. 996–1000.
- [87] X. Xia and W. Ding, "Hnas: Hierarchical neural architecture search on mobile devices," 2020.
- [88] H. Phan, Z. Liu, D. Huynh, M. Savvides, K. Cheng, and Z. Shen, "Binarizing mobilenet via evolution-based searching," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 13 417–13 426.
- [89] T. Cassimon, S. Vanneste, S. Bosmans, S. Mercelis, and P. Hellinckx, "Designing resource-constrained neural networks using neural architecture search targeting embedded devices," *Internet of Things*, p. 100234, 2020.
- [90] E. Lee and C. Lee, "Neuralscale: Efficient scaling of neurons for resource-constrained deep neural networks," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 1475–1484.
- [91] T. Wang, K. Wang, H. Cai, J. Lin, Z. Liu, H. Wang, Y. Lin, and S. Han, "Apq: Joint search for network architecture, pruning and quantization policy," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 2075–2084.
- [92] E. K. Wang, S. P. Xu, C. Chen, and N. Kumar, "Neural-architecture-search-based multiobjective cognitive automation system," *IEEE Systems Journal*, pp. 1–8, 2020.
- [93] A. Piergiovanni, A. Angelova, and M. S. Ryoo, "Tiny video networks," 2019.
- [94] R. Lee, Ł. Dudziak, M. Abdelfattah, S. I. Venieris, H. Kim, H. Wen, and N. D. Lane, "Journey towards tiny perceptual super-resolution," in *Computer Vision – ECCV 2020*, A. Vedaldi, H. Bischof, T. Brox, and J.-M. Frahm, Eds., 2020, pp. 85–102.
- [95] M. Ferianc, H. Fan, and M. Rodrigues, "Vinnas: Variational inference-based neural network architecture search," 2020.
- [96] Z. Lu, K. Deb, E. Goodman, W. Banzhaf, and V. N. Boddeti, "Nsganetv2: Evolutionary multi-objective surrogate-assisted neural architecture search," in *Computer Vision – ECCV 2020*, A. Vedaldi, H. Bischof, T. Brox, and J.-M. Frahm, Eds., 2020, pp. 35–51.
- [97] Y. Hu, X. Wu, and R. He, "Tf-nas: Rethinking three search freedoms of latency-constrained differentiable neural architecture search," 2020.
- [98] X. Wei, H. Chen, W. Liu, and Y. Xie, "Mixed-precision quantization for cnn-based remote sensing scene classification," *IEEE Geoscience and Remote Sensing Letters*, pp. 1–5, 2020.
- [99] H. Tsai, J. Ooi, C. Ferng, H. W. Chung, and J. Riesa, "Finding fast transformers: One-shot neural architecture search by component composition," *CoRR*, vol. abs/2008.06808, 2020. [Online]. Available: <https://arxiv.org/abs/2008.06808>
- [100] A. Marchisio, A. Massa, V. Mrazek, B. Bussolino, M. Martina, and M. Shafique, "Nascaps: A framework for neural architecture search to optimize the accuracy and hardware efficiency of convolutional capsule networks," in *ICCAD*, 2020.
- [101] D. Bruggemann, M. Kanakis, S. Georgoulis, and L. V. Gool, "Automated search for resource-efficient branched multi-task networks," 2020.
- [102] G. Li, M. Xu, S. Giancola, A. K. Thabet, and B. Ghanem, "LC-NAS: latency constrained neural architecture search for point cloud networks," *CoRR*, vol. abs/2008.10309, 2020. [Online]. Available: <https://arxiv.org/abs/2008.10309>
- [103] J. Lee, D. Kang, and S. Ha, "S3NAS: fast npu-aware neural architecture search methodology," *CoRR*, vol. abs/2009.02009, 2020. [Online]. Available: <https://arxiv.org/abs/2009.02009>
- [104] H. Chen, L. Zhuo, B. Zhang, X. Zheng, J. Liu, R. Ji, D. Doermann, and G. Guo, "Binarized neural architecture search for efficient object recognition," 2020.
- [105] W. Niu, Z. Kong, G. Yuan, W. Jiang, J. Guan, C. Ding, P. Zhao, S. Liu, B. Ren, and Y. Wang, "Real-time execution of large-scale language models on mobile," 2020.
- [106] C.-C. Wang, C.-T. Chiu, and J.-Y. Chang, "Efficientnet-elite: Extremely lightweight and efficient cnn models for edge devices by network candidate search," 2020.
- [107] Z. Yuan, X. Liu, B. Wu, and G. Sun, "Enas4d: Efficient multi-stage cnn architecture search for dynamic inference," *ArXiv*, vol. abs/2009.09182, 2020.
- [108] P. Achararit, M. A. Hanif, R. V. W. Putra, M. Shafique, and Y. Hara-Azumi, "Apnas: Accuracy-and-performance-aware neural architecture search for neural hardware accelerators," *IEEE Access*, vol. 8, pp. 165 319–165 334, 2020.
- [109] Z. Lu, I. Whalen, Y. Dhebar, K. Deb, E. Goodman, W. Banzhaf, and V. N. Boddeti, "Multi-objective evolutionary design of deep convolutional neural networks for image classification," *IEEE Transactions on Evolutionary Computation*, pp. 1–1, 2020.
- [110] J. G. López, A. Agudo, and F. Moreno-Noguer, "E-dnas: Differentiable neural architecture search for embedded systems,"
- [111] C. 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," 2020.
- [112] E. Liberis, Łukasz Dudziak, and N. D. Lane, "μnas: Constrained neural architecture search for microcontrollers," 2020.
- [113] Y. Yang, A. Nam, M. M. Nasr-Azadani, and T. Tung, "Resource-aware pareto-optimal automated machine learning platform," *CoRR*, vol. abs/2011.00073, 2020. [Online]. Available: <https://arxiv.org/abs/2011.00073>
- [114] K. A. Laube and A. Zell, "Shufflenasnets: Efficient cnn models through modified efficient neural architecture search," in *International Joint Conference on Neural Networks (IJCNN)*, 2019, pp. 1–6.
- [115] J. Li, W. Diao, X. Sun, Y. Feng, W. Zhang, Z. Chang, and K. Fu, "Automated and lightweight network design via random search for remote sensing image scene classification," *The International Archives of Photogrammetry, Remote Sensing and Spatial Information Sciences*, vol. 43, pp. 1217–1224, 2020.
- [116] W. Jiang, X. Zhang, E. Sha, L. Yang, Q. Zhuge, Y. Shi, and J. Hu, "Accuracy vs. efficiency: Achieving both through fpga-implementation aware neural architecture search," *56th ACM/IEEE Design Automation Conference (DAC)*, pp. 1–6, 2019.
- [117] Y. Yu, Y. Li, S. Che, N. K. Jha, and W. Zhang, "Software-defined design space exploration for an efficient AI accelerator architecture," *CoRR*, vol. abs/1903.07676, 2019. [Online]. Available: <http://arxiv.org/abs/1903.07676>
- [118] W. Jiang, L. Yang, E. H. Sha, Q. Zhuge, S. Gu, Y. Shi, and J. Hu, "Hardware/software co-exploration of neural architectures," *CoRR*, vol. abs/1907.04650, 2019. [Online]. Available: <http://arxiv.org/abs/1907.04650>
- [119] Q. Lu, W. Jiang, X. Xu, Y. Shi, and J. Hu, "On neural architecture search for resource-constrained hardware platforms," *CoRR*, vol. abs/1911.00105, 2019. [Online]. Available: <http://arxiv.org/abs/1911.00105>
- [120] M. S. Abdelfattah, Ł. Dudziak, T. Chau, R. Lee, H. Kim, and N. D. Lane, "Best of both worlds: Automl codesign of a cnn and its hardware accelerator," in *Proceedings of the 57th ACM/EDAC/IEEE Design Automation Conference*, 2020.
- [121] W. Chen, Y. Wang, S. Yang, C. Liu, and L. Zhang, "You only search once: A fast automation framework for single-stage dnn/accelerator co-design," in *Design, Automation & Test in Europe Conference & Exhibition, DATE*, 2020, pp. 1283–1286. [Online]. Available: <https://doi.org/10.23919/DATE48585.2020.9116474>
- [122] W. Jiang, L. Yang, S. Dasgupta, J. Hu, and Y. Shi, "Standing on the shoulders of giants: Hardware and neural architecture co-search with hot start," 2020.
- [123] J. Lin, W.-M. Chen, Y. Lin, J. Cohn, C. Gan, and S. Han, "Mcnunet: Tiny deep learning on iot devices," in *Advances in Neural Information Processing Systems (NeurIPS'20)*, 2020.
- [124] P. Colangelo, O. Segal, A. Speicher, and M. Margala, "Automl for multilayer perceptron and FPGA co-design," *CoRR*, vol. abs/2009.06156, 2020. [Online]. Available: <https://arxiv.org/abs/2009.06156>
- [125] K. Choi, D. Hong, H. Yoon, J. Yu, Y. Kim, and J. Lee, "Dance: Differentiable accelerator/network co-exploration," *arXiv preprint arXiv:2009.06237*, 2020.

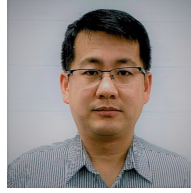


- [126] Y. Zhang, Y. Fu, W. Jiang, C. Li, H. You, M. Li, V. Chandra, and Y. Lin, "Dna: Differentiable network-accelerator co-search," *arXiv preprint arXiv:2010.14778*, 2020.
- [127] G. Chu, O. Arikian, G. Bender, W. Wang, A. Brighton, P. Kindermans, H. Liu, B. Akin, S. Gupta, and A. Howard, "Discovering multi-hardware mobile models via architecture search," *CoRR*, vol. abs/2008.08178, 2020. [Online]. Available: <https://arxiv.org/abs/2008.08178>
- [128] Y. Jiang, X. Wang, and W. Zhu, "Hardware-aware transformable architecture search with efficient search space," in *IEEE International Conference on Multimedia and Expo (ICME)*, 2020, pp. 1–6.
- [129] L. Yang, Z. Yan, M. Li, H. Kwon, W. Jiang, L. Lai, Y. Shi, T. Krishna, and V. Chandra, *Co-Exploration of Neural Architectures and Heterogeneous ASIC Accelerator Designs Targeting Multiple Tasks*, 2020.
- [130] C. Liu, L. Chen, F. Schroff, H. Adam, W. Hua, A. L. Yuille, and F. Li, "Auto-deeplab: Hierarchical neural architecture search for semantic image segmentation," in *Conference on Computer Vision and Pattern Recognition, CVPR*, 2019, pp. 82–92.
- [131] S. Sabour, N. Frosst, and G. E. Hinton, "Dynamic routing between capsules," in *Advances in Neural Information Processing Systems 30*, I. Guyon, U. von Luxburg, S. Bengio, H. M. Wallach, R. Fergus, S. V. N. Vishwanathan, and R. Garnett, Eds., 2017, pp. 3856–3866. [Online]. Available: <http://papers.nips.cc/paper/6975-dynamic-routing-between-capsules>
- [132] H. Wang, Z. Wu, Z. Liu, H. Cai, L. Zhu, C. Gan, and S. Han, "HAT: Hardware-aware transformers for efficient natural language processing," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 7675–7688.
- [133] 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 *Proceedings of the 56th Annual Design Automation Conference*, 2019.
- [134] "Why tinymml is a giant opportunity," <https://venturebeat.com/2020/01/11/why-tinymml-is-a-giant-opportunity/>.
- [135] C. Fu, H. Chen, Z. Yang, F. Koushanfar, Y. Tian, and J. Zhao, "Enhancing model parallelism in neural architecture search for multidevice system," *IEEE Micro*, vol. 40, no. 5, pp. 46–55, 2020.
- [136] Y.-H. Kim, B. Reddy, S. Yun, and C. Seo, "Nemo : Neuro-evolution with multiobjective optimization of deep neural network for speed and accuracy," 2017.
- [137] K. Hu, S. Tian, S. Guo, N. Li, L. Luo, and L. Wang, "Recurrent neural architecture search based on randomness-enhanced tabu algorithm," in *International Joint Conference on Neural Networks (IJCNN)*, 2020.
- [138] Y. Weng, T. Zhou, Y. Li, and X. Qiu, "Nas-unet: Neural architecture search for medical image segmentation," *IEEE Access*, vol. 7, pp. 44 247–44 257, 2019.
- [139] N. Klyuchnikov, I. Trofimov, E. Artemova, M. Salnikov, M. Fedorov, and E. Burnaev, "Nas-bench-nlp: Neural architecture search benchmark for natural language processing," 2020.
- [140] K. Deb, A. Pratap, S. Agarwal, and T. Meyarivan, "A fast and elitist multiobjective genetic algorithm: Nsga-ii," *IEEE transactions on evolutionary computation*, vol. 6, no. 2, pp. 182–197, 2002.
- [141] X. Chu, B. Zhang, R. Xu, and H. Ma, "Multi-objective reinforced evolution in mobile neural architecture search," *CoRR*, vol. abs/1901.01074, 2019. [Online]. Available: <http://arxiv.org/abs/1901.01074>
- [142] X. Zhang, Z. Huang, and N. Wang, "You only search once: Single shot neural architecture search via direct sparse optimization," *CoRR*, vol. abs/1811.01567, 2018. [Online]. Available: <http://arxiv.org/abs/1811.01567>
- [143] D. Stamoulis, R. Ding, D. Wang, D. Lymberopoulos, B. Priyantha, J. Liu, and D. Marculescu, "Single-path NAS: device-aware efficient convnet design," *CoRR*, vol. abs/1905.04159, 2019. [Online]. Available: <http://arxiv.org/abs/1905.04159>
- [144] M. Gennari do Nascimento, T. W. Costain, and V. A. Prisacariu, "Finding Non-Uniform Quantization Schemes using Multi-Task Gaussian Processes," *arXiv e-prints*, p. arXiv:2007.07743, Jul. 2020.
- [145] X. Chu, B. Zhang, H. Ma, R. Xu, J. Li, and Q. Li, "Fast, accurate and lightweight super-resolution with neural architecture search," *CoRR*, vol. abs/1901.07261, 2019. [Online]. Available: <http://arxiv.org/abs/1901.07261>
- [146] G. Bender, H. Liu, B. Chen, G. Chu, S. Cheng, P. J. Kindermans, and Q. V. Le, "Can weight sharing outperform random architecture search? an investigation with tunas," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 14 311–14 320.
- [147] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le, "Learning transferable architectures for scalable image recognition," 2018.
- [148] E. Real, S. Moore, A. Selle, S. Saxena, Y. L. Suematsu, J. Tan, Q. V. Le, and A. Kurakin, "Large-scale evolution of image classifiers," in *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, 2017, p. 2902–2911.
- [149] N. Nayman, A. Noy, T. Ridnik, I. Friedman, R. Jin, and L. Zelnik-Manor, "XNAS: neural architecture search with expert advice," in *Advances in Neural Information Processing System 32 (NeurIPS)*, 2019, pp. 1975–1985. [Online]. Available: <http://papers.nips.cc/paper/8472-xnas-neural-architecture-search-with-expert-advice>
- [150] H. Liu, K. Simonyan, and Y. Yang, "DARTS: differentiable architecture search," in *7th International Conference on Learning Representations, ICLR*, 2019.
- [151] M. Courbariaux, Y. Bengio, and J.-P. David, "Binaryconnect: Training deep neural networks with binary weights during propagations," in *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2*, 2015, p. 3123–3131.
- [152] N. Cesa-Bianchi and G. Lugosi, *Prediction, Learning, and Games*. Cambridge University Press, 2006.
- [153] J. Kivinen and M. K. Warmuth, "Exponentiated gradient versus gradient descent for linear predictors," *information and computation*, vol. 132, no. 1, pp. 1–63, 1997.
- [154] L. Li and A. Talwalkar, "Random search and reproducibility for neural architecture search," in *Uncertainty in Artificial Intelligence*. PMLR, 2020, pp. 367–377.
- [155] H. Pham, M. Y. Guan, B. Zoph, Q. V. Le, and J. Dean, "Efficient neural architecture search via parameter sharing," in *Proceedings of the 35th International Conference on Machine Learning, ICML*, 2018, pp. 4092–4101. [Online]. Available: <http://proceedings.mlr.press/v80/pham18a.html>
- [156] E. Jang, S. Gu, and B. Poole, "Categorical reparameterization with gumbel-softmax," in *5th International Conference on Learning Representations, ICLR*, 2017. [Online]. Available: <https://openreview.net/forum?id=rkE3y85ee>
- [157] B. Wu, Y. Wang, P. Zhang, Y. Tian, P. Vajda, and K. Keutzer, "Mixed precision quantization of convnets via differentiable neural architecture search," *CoRR*, vol. abs/1812.00090, 2018. [Online]. Available: <http://arxiv.org/abs/1812.00090>
- [158] H. Mikami, H. Suganuma, P. U.-Chupala, Y. Tanaka, and Y. Kageyama, "Imagenet/resnet-50 training in 224 seconds," *CoRR*, vol. abs/1811.05233, 2018. [Online]. Available: <http://arxiv.org/abs/1811.05233>
- [159] B. Deng, J. Yan, and D. Lin, "Peephole: Predicting network performance before training," *CoRR*, vol. abs/1712.03351, 2017. [Online]. Available: <http://arxiv.org/abs/1712.03351>
- [160] C. Liu, B. Zoph, M. Neumann, J. Shlens, W. Hua, L.-J. Li, L. Fei-Fei, A. Yuille, J. Huang, and K. Murphy, "Progressive neural architecture search," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 19–34.
- [161] B. Baker, O. Gupta, R. Raskar, and N. Naik, "Accelerating neural architecture search using performance prediction," in *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Workshop Track Proceedings*. OpenReview.net, 2018. [Online]. Available: <https://openreview.net/forum?id=HJqk3N1vG>
- [162] Y. Tang, Y. Wang, Y. Xu, H. Chen, B. Shi, C. Xu, C. Xu, Q. Tian, and C. Xu, "A semi-supervised assessor of neural architectures," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1810–1819.
- [163] H. Bouzidi, H. Ouarnoughi, S. Niar, and A. A. E. Cadi, "Performance prediction for convolutional neural networks in edge devices," 2020.
- [164] T. Yang, A. G. Howard, B. Chen, X. Zhang, A. Go, M. Sandler, V. Sze, and H. Adam, "Netadapt: Platform-aware neural network adaptation for mobile applications," in *ECCV 15th European Conference Computer Vision*, V. Ferrari, M. Hebert, C. Sminchisescu, and Y. Weiss, Eds., vol. 11214, 2018, pp. 289–304. [Online]. Available: [https://doi.org/10.1007/978-3-030-01249-6\\_18](https://doi.org/10.1007/978-3-030-01249-6_18)
- [165] C. Gong, Z. Jiang, D. Wang, Y. Lin, Q. Liu, and D. Z. Pan, "Mixed precision neural architecture search for energy efficient deep learning," in *ICCAD*, 2019, pp. 1–7.
- [166] H. Kwon, P. Chatarasi, M. Pellauer, A. Parashar, V. Sarkar, and T. Krishna, "Understanding reuse, performance, and hardware cost of dnn dataflow: A data-centric approach," in *Proceedings of the 52nd Annual IEEE/ACM International Symposium on Microarchitecture*, 2019, p. 754–768.
- [167] A. Anderson, J. Su, R. Dahyot, and D. Gregg, "Performance-oriented neural architecture search," in *International Conference on High Performance Computing Simulation (HPCS)*, 2019, pp. 177–184.

- [168] C. Ying, A. Klein, E. Christiansen, E. Real, K. Murphy, and F. Hutter, “Nas-bench-101: Towards reproducible neural architecture search,” in *Proceedings of the 36th International Conference on Machine Learning, ICML*, ser. Proceedings of Machine Learning Research, K. Chaudhuri and R. Salakhutdinov, Eds., vol. 97. PMLR, 2019, pp. 7105–7114.
- [169] X. Zhang, S. Liu, R. Zhang, C. Liu, D. Huang, S. Zhou, J. Guo, Y. Kang, Q. Guo, Z. Du *et al.*, “Adaptive precision training: Quantify back propagation in neural networks with fixed-point numbers,” *arXiv preprint arXiv:1911.00361*, 2019.
- [170] U. Asif, J. Tang, and S. Harrer, “Ensemble knowledge distillation for learning improved and efficient networks,” in *ECAI - 24th European Conference on Artificial Intelligence*, G. D. Giacomo, A. Catalá, B. Dilkina, M. Milano, S. Barro, A. Bugarin, and J. Lang, Eds., 2020. [Online]. Available: <https://doi.org/10.3233/FAIA200188>
- [171] G. Boukli Hacene, “Processing and learning deep neural networks on chip,” Theses, Ecole nationale supérieure Mines-Télécom Atlantique, Oct. 2019. [Online]. Available: <https://tel.archives-ouvertes.fr/tel-02438921>
- [172] I. Hubara, M. Courbariaux, D. Soudry, R. El-Yaniv, and Y. Bengio, “Binarized neural networks,” in *Advances in Neural Information Processing Systems* 29, D. D. Lee, M. Sugiyama, U. von Luxburg, I. Guyon, and R. Garnett, Eds., 2016, pp. 4107–4115. [Online]. Available: <http://papers.nips.cc/paper/6573-binarized-neural-networks>
- [173] H. Yu, Q. Han, J. Li, J. Shi, G. Cheng, and B. Fan, “Search what you want: Barrier panelty NAS for mixed precision quantization,” in *ECCV - 16th European Conference Computer Vision*, A. Vedaldi, H. Bischof, T. Brox, and J. Frahm, Eds., vol. 12354, 2020, pp. 1–16. [Online]. Available: [https://doi.org/10.1007/978-3-030-58545-7\\_1](https://doi.org/10.1007/978-3-030-58545-7_1)
- [174] V. Camus, L. Mei, C. Enz, and M. Verhelst, “Review and benchmarking of precision-scalable multiply-accumulate unit architectures for embedded neural-network processing,” *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, vol. 9, no. 4, pp. 697–711, 2019.
- [175] Anonymous, “Uniform-precision neural network quantization via neural channel expansion,” in *Submitted to International Conference on Learning Representations*, 2021, under review. [Online]. Available: <https://openreview.net/forum?id=oGq4d9TbyIA>
- [176] X. Dong and Y. Yang, “Network pruning via transformable architecture search,” in *Advances in Neural Information Processing Systems*, 2019, pp. 760–771.
- [177] M. Lin, R. Ji, Y. Zhang, B. Zhang, Y. Wu, and Y. Tian, “Channel pruning via automatic structure search,” in *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI*, C. Bessiere, Ed. ijcai.org, 2020, pp. 673–679. [Online]. Available: <https://doi.org/10.24963/ijcai.2020/94>
- [178] X. Li, Y. Zhou, Z. Pan, and J. Feng, “Partial order pruning: For best speed/accuracy trade-off in neural architecture search,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [179] D. V. Vargas and S. Kotyan, “Evolving robust neural architectures to defend from adversarial attacks,” *CoRR*, vol. abs/1906.11667, 2019. [Online]. Available: <http://arxiv.org/abs/1906.11667>
- [180] E. D. Cubuk, B. Zoph, S. S. Schoenholz, and Q. V. Le, “Intriguing properties of adversarial examples,” in *6th International Conference on Learning Representations, ICLR*. OpenReview.net, 2018.
- [181] M. Guo, Y. Yang, R. Xu, Z. Liu, and D. Lin, “When NAS meets robustness: In search of robust architectures against adversarial attacks,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*. IEEE, 2020, pp. 628–637. [Online]. Available: <https://doi.org/10.1109/CVPR42600.2020.00071>
- [182] T. T. Le, W. Fu, and J. H. Moore, “Scaling tree-based automated machine learning to biomedical big data with a feature set selector,” *Bioinformatics*, vol. 36, no. 1, pp. 250–256, 2020.
- [183] X. Dong and Y. Yang, “Nas-bench-201: Extending the scope of reproducible neural architecture search,” in *International Conference on Learning Representations*, 2020. [Online]. Available: <https://openreview.net/forum?id=HJxyZkBKDr>
- [184] A. Zela, J. Siems, and F. Hutter, “Nas-bench-1shot1: Benchmarking and dissecting one-shot neural architecture search,” in *8th International Conference on Learning Representations, ICLR*. OpenReview.net, 2020. [Online]. Available: <https://openreview.net/forum?id=SJx9ngStPH>
- [185] X. Dong, L. Liu, K. Musial, and B. Gabrys, “Nats-bench: Benchmarking NAS algorithms for architecture topology and size,” *CoRR*, vol. abs/2009.00437, 2020. [Online]. Available: <https://arxiv.org/abs/2009.00437>
- [186] Anonymous, “{NAS}-bench-301 and the case for surrogate benchmarks for neural architecture search,” in *Submitted to International Conference on Learning Representations*, 2021, under review. [Online]. Available: <https://openreview.net/forum?id=1flmvXGGJaa>
- [187] C. Li, Z. Yu, Y. Fu, Y. Zhang, Y. Zhao, H. You, Q. Yu, Y. Wang, C. Hao, and Y. Lin, “{HW}-{nas}-{bench}: Hardware-aware neural architecture search benchmark,” in *International Conference on Learning Representations*, 2021. [Online]. Available: [https://openreview.net/forum?id=\\_0kaDkv3dVf](https://openreview.net/forum?id=_0kaDkv3dVf)
- [188] S. Merity, N. S. Keskar, and R. Socher, “Regularizing and optimizing LSTM language models,” in *6th International Conference on Learning Representations, ICLR*, 2018. [Online]. Available: <https://openreview.net/forum?id=SyyGPP0TZ>
- [189] M. P. Marcus, B. Santorini, and M. A. Marcinkiewicz, “Building a large annotated corpus of english: The penn treebank,” *Comput. Linguistics*, vol. 19, no. 2, pp. 313–330, 1993.
- [190] S. Merity, C. Xiong, J. Bradbury, and R. Socher, “Pointer sentinel mixture models,” in *5th International Conference on Learning Representations, ICLR*. OpenReview.net, 2017. [Online]. Available: <https://openreview.net/forum?id=Byj72udxe>
- [191] J. Mellor, J. Turner, A. J. Storkey, and E. J. Crowley, “Neural architecture search without training,” *CoRR*, vol. abs/2006.04647, 2020. [Online]. Available: <https://arxiv.org/abs/2006.04647>
- [192] G. J. van Wyk and A. S. Bosman, “Evolutionary neural architecture search for image restoration,” in *International Joint Conference on Neural Networks (IJCNN)*, 2019, pp. 1–8.
- [193] H. Zhang, Y. Li, H. Chen, and C. Shen, “IR-NAS: neural architecture search for image restoration,” *CoRR*, vol. abs/1909.08228, 2019. [Online]. Available: <http://arxiv.org/abs/1909.08228>
- [194] X. Gong, S. Chang, Y. Jiang, and Z. Wang, “Autogan: Neural architecture search for generative adversarial networks,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 3224–3234.
- [195] K. Ho, A. Gilbert, H. Jin, and J. P. Collomosse, “Neural architecture search for deep image prior,” *CoRR*, vol. abs/2001.04776, 2020. [Online]. Available: <https://arxiv.org/abs/2001.04776>
- [196] Z. Zhu, C. Liu, D. Yang, A. Yuille, and D. Xu, “V-nas: Neural architecture search for volumetric medical image segmentation,” in *2019 International Conference on 3D Vision (3DV)*. IEEE, 2019, pp. 240–248.
- [197] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, “Imagenet large scale visual recognition challenge,” *International journal of computer vision*, vol. 115, no. 3, pp. 211–252, 2015.
- [198] Z. Lu, G. Sreekumar, E. D. Goodman, W. Banzhaf, K. Deb, and V. N. Boddeti, “Neural architecture transfer,” *CoRR*, vol. abs/2005.05859, 2020. [Online]. Available: <https://arxiv.org/abs/2005.05859>
- [199] M. G. D. Nascimento, T. W. Costain, and V. A. Prisacariu, “Finding non-uniform quantization schemes using multi-task gaussian processes,” in *ECCV - 16th European Conference Computer Vision*, A. Vedaldi, H. Bischof, T. Brox, and J. Frahm, Eds., vol. 12362, 2020, pp. 383–398.
- [200] A. L. Maas, R. E. Daly, P. T. Pham, D. Huang, A. Y. Ng, and C. Potts, “Learning word vectors for sentiment analysis,” in *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, 2011, pp. 142–150. [Online]. Available: <http://www.aclweb.org/anthology/P11-1015>
- [201] S. Wang, A. Pathania, and T. Mitra, “Neural network inference on mobile socs,” *IEEE Design Test*, vol. 37, no. 5, pp. 50–57, 2020.
- [202] 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*, pp. 1–1, 2020.
- [203] A. Sebastian, I. Boybat, M. Dazzi, I. Giannopoulos, V. Jonnalagadda, V. Joshi, G. Karunaratne, B. Kersting, R. Khaddam-Aljameh, S. R. Nandakumar, A. Petropoulos, C. Piveteau, T. Antonakopoulos, B. Rajendran, M. L. Gallo, and E. Eleftheriou, “Computational memory-based inference and training of deep neural networks,” in *2019 Symposium on VLSI Technology*, 2019, pp. T168–T169.



**Hadjer Benmeziiane** received her Master and Engineering degrees in Computer Science at the Higher National School of Computer Science, Algiers, Algeria in 2020. She is currently pursuing her PhD degree in Computer Science at LAMIH/CNRS, Université Polytechnique Hauts-de-France, Valenciennes, France.



**Dr. Naigang Wang** received his PhD in materials science and engineering from University of Florida in 2010. Since then, he has been affiliated with IBM T. J. Watson Research Center as a research staff member. His current research focus is on hardware-friendly deep learning algorithms for artificial intelligence accelerators.



**Dr. Kaoutar El Maghraoui** received her PhD in computer science from Rensselaer Polytechnic Institute in 2007. Since then, she has been affiliated with the IBM T.J. Watson Research Center. Kaoutar currently holds the position of a principal research scientist at the IBM Research AI organization where she is focusing on innovations at the intersection of systems and artificial intelligence. She leads the End-Use experimental AI testbed of the IBM Research AI Hardware Center, a global research hub focusing on enabling next-generation chips and systems for

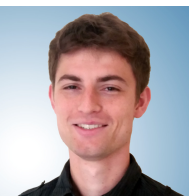
AI workloads. Her research interests include AI platforms, distributed deep learning, neural architecture search, distributed systems, and AI hardware acceleration.



**Dr. Hamza Ouarnoughi** received his PhD in computer science from the University of Western Brittany in 2017. Since then, he has been associate professor at LAMIH/CNRS, Université Polytechnique Hauts-de-France and INSA Hauts-de-France, Valenciennes, France. His research interests include computer architecture and systems, cloud computing, storage systems, and artificial intelligence.



**Pr. Smail Niar** received his PhD in computer science from the University of Lille in 1990. Since then, he has been professor at LAMIH/CNRS, Université Polytechnique Hauts-de-France and INSA Hauts-de-France, Valenciennes, France. He is member of the European Network of Excellence on “High Performance and Embedded Architectures and Compilation” (HIPEAC), EuroMicro society and IEEE senior member. His research interests include heterogeneous multi-core architectures, autonomous driving and reliability for intelligent transportation systems.



**Dr. Martin Wistuba** is a researcher at IBM Research, where he develops tools to automate deep learning. Previously, he received his Ph.D. in Machine Learning from the University of Hildesheim. His research interest includes AutoML, in particular the idea of meta-knowledge transfer to speed up Bayesian optimization and Neural Architecture Search.