

# Pareto-aware Neural Architecture Generation for Diverse Computational Budgets

Yong Guo, Yaofu Chen, Yin Zheng, Qi Chen, Peilin Zhao, Jian Chen, Junzhou Huang, Mingkui Tan\*

**Abstract**—Designing feasible and effective architectures under diverse computational budgets, incurred by different applications/devices, is essential for deploying deep models in real-world applications. To achieve this goal, existing methods often perform an independent architecture search process for each target budget, which is very inefficient yet unnecessary. More critically, these independent search processes cannot share their learned knowledge (*i.e.*, the distribution of good architectures) with each other and thus often result in limited search results. To address these issues, we propose a Pareto-aware Neural Architecture Generator (PNAG) which only needs to be trained once and dynamically produces the Pareto optimal architecture for any given budget via inference. To train our PNAG, we learn the whole Pareto frontier by jointly finding multiple Pareto optimal architectures under diverse budgets. Such a joint search algorithm not only greatly reduces the overall search cost but also improves the search results. Extensive experiments on three hardware platforms (*i.e.*, mobile device, CPU, and GPU) show the superiority of our method over existing methods.

**Index Terms**—Neural Architecture Generation, Pareto Frontier Learning, Architecture Design under Budgets.



## 1 INTRODUCTION

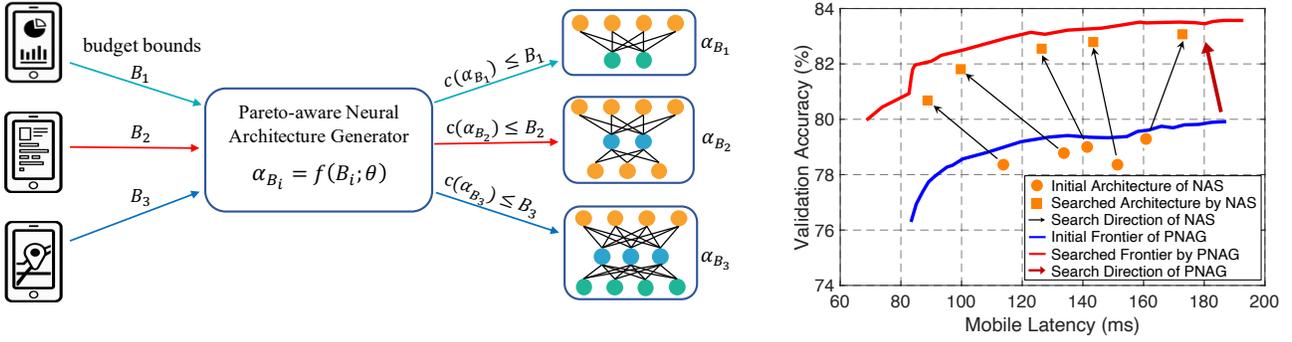
DEEP neural networks (DNNs) [1] have been the workhorse of many challenging tasks, including image classification [2], [3], [4], [5], [6], semantic segmentation [7], [8], [9], [10] and object detection [11], [12], [13], [14]. However, designing effective architectures often relies heavily on human expertise. To alleviate this issue, neural architecture search (NAS) methods have been proposed to automatically design effective architectures [15]. Existing studies show that these automatically searched architectures often outperform the manually designed ones in many computer vision tasks [16], [17], [18], [19], [20], [21], [22], [23].

However, the state-of-the-art deep networks often contain a large number of parameters and come with extremely high computational cost. As a result, it is hard to deploy these models to real-world scenarios with limited computation resources. Regarding this issue, we have to carefully design architectures to fulfill a specific computational budget (*e.g.*, a feasible model should have a latency lower than 100ms on a specified mobile device). More critically, we may have to consider different computational budgets in the real world. For example, a company may simultaneously develop/maintain multiple applications and each of them has a specific budget of latency.

In order to design feasible architectures, most methods [24], [25] only considers a single computational budget and incorporates architecture’s computational cost into the objective function of NAS. When we consider diverse budgets, they have to conduct an independent search process for each budget [24], which is very inefficient yet unnecessary. Unlike these methods, one can also exploit the population-based methods to simultaneously find multiple architectures and then select an appropriate one from them to fulfill a specific budget [26], [27]. However, due to the limited population size, these searched architectures do not necessarily satisfy the required budget. More critically, all these searched architectures are fixed after search and cannot be easily adapted for a slightly changed budget. Thus, how to design effective architectures under diverse computational budgets in an efficient and flexible way still remains an open question.

In this paper, we propose a Pareto-aware Neural Architecture Generator (PNAG) which only needs to be trained once and then dynamically produces Pareto optimal architectures for diverse budgets via *inference* (as shown in Fig. 1(a)). Note that the Pareto optimal architectures under different budgets should lie on a distribution, *i.e.*, the Pareto frontier over model performance and computational cost [28]. We propose to jointly learn the whole Pareto frontier (*i.e.*, improving the blue curve to the red curve in Fig. 1(b)) instead of finding a single Pareto optimal architecture. During training, we randomly sample budgets from a predefined distribution and maximize the expected reward of the searched architectures to approximate the ground-truth Pareto frontier. It is worth noting that learning the Pareto frontier is able to share the learned knowledge across different budgets and greatly improve the search results in practice (see results in Table 5). Furthermore, when evaluating architectures under diverse budgets, we design an architecture evaluator that learns a Pareto dominance rule to determine which architecture is a relatively better one in pairwise comparisons. Unlike existing methods, we highlight that our PNAG designs architectures through a generation process instead of search, which is very efficient (see results in Table 6) and practically useful in real-world model deployment.

- Yong Guo is with the School of Software Engineering, South China University of Technology. Yong Guo is also with Max Planck Institute for Informatics. E-mail: guo.yong@mail.scut.edu.cn
- Yaofu Chen is with the School of Software Engineering, South China University of Technology. E-mail: sechenyaofu@mail.scut.edu.cn
- Yin Zheng is with the Weixin Group, Tencent, China. E-mail: yzheng3xg@gmail.com
- Qi Chen is with the School of Computer Science, the University of Adelaide. E-mail: qi.chen04@adelaide.edu.au
- Peilin Zhao and Junzhou Huang are with Tencent AI Lab, Tencent, China. E-mail: {masonzhao, joe.huang}@tencent.com
- Mingkui Tan and Jian Chen are with the School of Software Engineering, South China University of Technology. Mingkui Tan is also with the Key Laboratory of Big Data and Intelligent Robot (South China University of Technology), Ministry of Education. Mingkui Tan is also with the Pazhou Laboratory, Guangzhou, China. E-mail: {mingkuitan, ellsachen}@scut.edu.cn
- \* Corresponding author



(a) An illustration of generating feasible architectures for diverse budgets using PNAG. (b) Comparisons between PNAG and conventional NAS methods.

Fig. 1. We show an illustration of how to apply PNAG to generate feasible architectures for diverse computational budgets and the comparisons between PNAG and conventional NAS methods. (a) PNAG takes an arbitrary budget as input and flexibly generates architectures. (b) PNAG learns the whole Pareto frontier rather than finding discrete architectures. Here, the accuracy is measured on the constructed validation set.

We summarize the contributions of our paper as follows.

- Instead of designing architectures for a single budget, we propose a Pareto-aware Neural Architecture Generator (PNAG) which is only trained once and flexibly generates effective architectures for arbitrary budget via inference (see Fig. 1(a)). In this way, our architecture generation process becomes very efficient and practically useful in real-world applications.
- To train our PNAG, we propose to explicitly learn the Pareto frontier by maximizing the expected reward of the searched architectures over diverse budgets. Interestingly, learning the Pareto frontier shares the learned knowledge across the search processes under diverse budgets and greatly improves the search results (see results in Table 5).
- Since an architecture should have different rewards/scores under different budgets, we propose an architecture evaluator to adaptively evaluate architectures for any given budget. To train the evaluator, we propose to learn a Pareto dominance rule which determines whether an architecture is better than the other in pairwise comparisons.
- We measure the latencies on three hardware platforms and take them as the computational budgets to generate feasible architectures. Extensive experiments show that the architectures produced by PNAG consistently outperform the architectures searched by existing methods across different budgets and platforms.

## 2 RELATED WORK

In this section, we provide a brief overview of existing work on neural architecture search, architecture design under resource constraints, as well as Pareto frontier learning.

### 2.1 Neural Architecture Search (NAS)

Unlike manually designing architectures with expert knowledge, NAS seeks to automatically design more effective architectures [29], [30], [31], [32], [33]. Existing NAS methods can be roughly divided into three categories, namely, reinforcement-learning-based methods, evolutionary approaches, and gradient-based methods. Specifically, reinforcement-learning-based methods [15], [34], [35], [36], [37] learn a controller to produce architectures. Evolutionary approaches [38], [39], [40], [41], [42], [43] search for promising architectures by gradually evolving a

population. Gradient-based methods [44], [45], [46], [47], [48], [49] relax the search space to be continuous and optimize architectures by gradient descent. Besides designing effective search algorithms, many efforts have also been made to improve the accuracy of architecture evaluation [50], [51], [52]. Unlike these methods that find a single architecture, one can design different architectures by training an architecture generator. Specifically, RandWire [53] designs stochastic network generators to generate randomly wired architectures. NAGO [54] is the first work to learn an architecture generator and proposes a hierarchical and graph-based search space to reduce the optimization difficulty. However, these generated architectures tend to perform very similarly (*i.e.*, low diversity) in terms of both model performance and computational cost [53], [54]. Thus, these architectures may not satisfy an arbitrary required budget. In other words, they still have to learn a generator for a required budget to produce feasible architectures.

### 2.2 Architecture Design under Resource Constraints

Many efforts have been made in designing architectures under a resource constraint [55], [56], [57], [58], [59], [60]. Specifically, PONAS [56] builds an accuracy table to find architectures satisfying a single budget constraint. TuNAS [58] proposes a reward function to restrict the latency of the searched architecture, which omits additional hyper-parameter tuning. Related to our work, SGNAS [61] proposes an architecture generator which generates architectures for specific budget constraints. Nevertheless, SGNAS optimizes a regression loss *w.r.t.* budget constraint and the resultant architecture does not necessarily have lower cost than the target budget, *i.e.*, violating the budget. More critically, SGNAS considers a fixed hyper-parameter  $\lambda$  to balance the regression loss and a classification loss. Due to the large diversity among architectures, their accuracy and computational cost may vary significantly across different budgets, also leading to suboptimal search results (See Table 1).

### 2.3 Pareto Frontier Learning

Given multiple objectives, Pareto frontier learning aims to find a set of Pareto optimal solutions over them. Most methods exploit evolutionary algorithms [62], [63] to solve this problem. Inspired by them, many efforts have been made to simultaneously find a set of Pareto optimal architectures over accuracy and computational cost [64], [65]. Recently, NSGANetV1 [66] presents

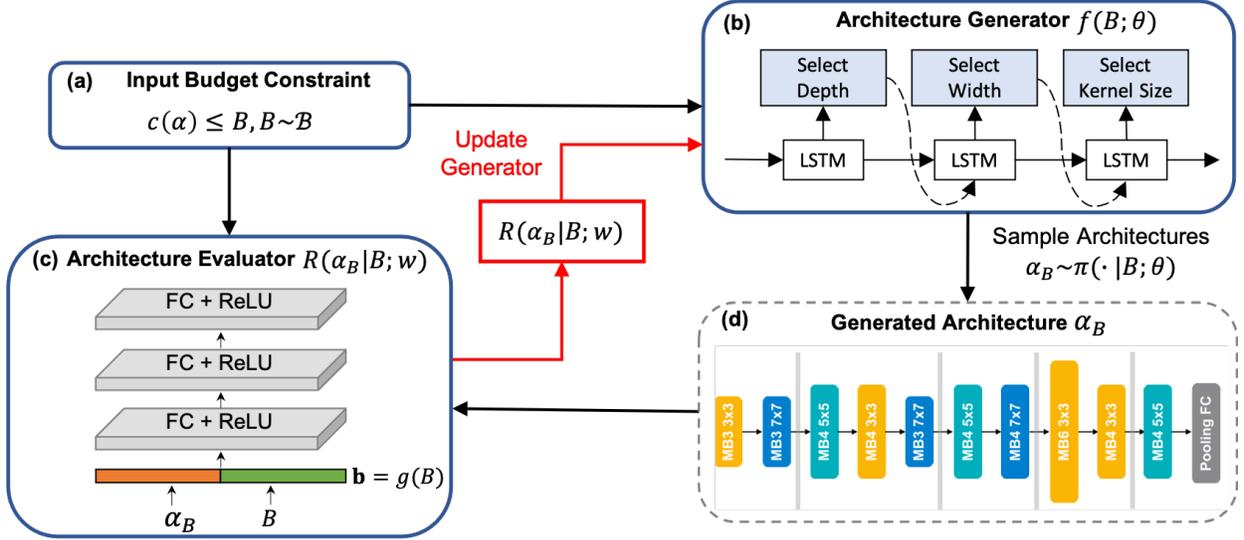


Fig. 2. Overview of the proposed PNAG. Our PNAG mainly consists of two modules: an architecture generator  $f(\cdot; \theta)$  and an architecture evaluator  $R(\cdot; w)$ . Specifically, we build the generator model based on an LSTM network, which takes a budget constraint  $B$  as input and produces a promising architecture  $\alpha_B$  that satisfies the budget constraint, *i.e.*,  $c(\alpha)$ . To optimize the generator model, we design the evaluator using three fully connected (FC) layers to estimate the performance of the generated architectures  $\alpha_B$ . The orange and green boxes in (c) denote the embeddings of architecture  $\alpha_B$  and the budget w.r.t.  $B$ , respectively.

an evolutionary approach to find a set of trade-off architectures over multiple objectives in a single run. NSGANetV2 [27] further presents two surrogates (at the architecture and weights level) to produce task-specific models under multiple competing objectives. Given a target budget, these methods may manually select an appropriate architecture from a set of searched architectures. However, given limited population size, the selected architectures do not necessarily satisfy a required budget. More critically, all the searched architectures are fixed after search and cannot be easily adapted for a slightly changed budget. Thus, how to learn the Pareto frontier and use it to generate architectures for arbitrary budget in a flexible way still remains unexplored.

### 3 PARETO-AWARE ARCHITECTURE GENERATION

In this paper, we focus on the architecture generation problem and intend to generate effective architectures for diverse computational budgets via *inference* instead of search/training. Note that the optimal architectures under different budgets lie on the Pareto frontier over model performance and computational cost [28]. Thus, we develop a Pareto-aware Neural Architecture Generator (PNAG) to explicitly learn the whole Pareto frontier. To locate the best architecture from the frontier for a given budget, we build our PNAG as a conditional model which takes the budget as input and directly produces a feasible architecture. In Section 3.1, we depict our architecture generator model and present a novel learning algorithm to learn the Pareto frontier. In Section 3.2, we propose an architecture evaluator, as well as its training algorithm, to adaptively evaluate architectures under different budgets. Algorithm 1 shows the whole training process of PNAG.

#### 3.1 Learning the Architecture Generator $f(B; \theta)$

We seek to build an architecture generator model to dynamically and flexibly produce effective architectures for any given computational budget. Let  $B$  be a budget (*e.g.*, latency or MAdds) which can be considered as a random variable drawn from some

distribution  $\mathcal{B}$ , namely  $B \sim \mathcal{B}$ . Let  $\Omega$  be an architecture search space. For any architecture  $\alpha \in \Omega$ , we use  $c(\alpha)$  and  $\text{Acc}(\alpha)$  to measure the cost and validation accuracy of  $\alpha$ , respectively.

Since an architecture can be represented as a sequence of tokens (each token denotes a setting of a layer, *e.g.*, width or kernel size) [15], [34], we cast the architecture generation problem as a sequential decision problem and build the architecture generator  $f(B; \theta)$  using an LSTM network. As shown in Fig. 2, the generator takes a budget  $B$  as input and generates architectures  $\alpha_B = f(B; \theta)$  (satisfying the constraint  $c(\alpha_B) \leq B$ ) by sequentially predicting the token sequences, *i.e.*, the depth, width, and kernel size of each layer. Here,  $\theta$  denotes the learnable parameters. Note that the optimal architecture under a specific budget should lie on the Pareto frontier over model performance and computational cost. To make the generator generalize to arbitrary budget, we seek to learn the Pareto frontier rather than finding discrete architectures. In the following, we first illustrate our training method in Section 3.1.1 and then discuss how to represent a budget with arbitrary value in Section 3.1.2.

##### 3.1.1 Training Method of $f(B; \theta)$

To illustrate the training objective of our method, we first revisit the NAS problem with a single budget and then generalize it to the problem with diverse budgets.

**NAS under a single budget.** Since it is non-trivial to directly find the optimal architecture [15], by contrast, one can first learn a policy  $\pi(\cdot; \theta)$  and then conduct sampling from it to find promising architectures, *i.e.*,  $\alpha \sim \pi(\cdot; \theta)$ . Given a budget  $B$ , the optimization problem becomes

$$\max_{\theta} \mathbb{E}_{\alpha \sim \pi(\cdot; \theta)} [R(\alpha|B; w)], \text{ s.t. } c(\alpha) \leq B. \quad (1)$$

Here,  $\pi(\cdot; \theta)$  is the learned policy parameterized by  $\theta$ , and  $R(\alpha|B; w)$  is the reward function parameterized by  $w$  that measures the joint performance of both the accuracy and latency of  $\alpha$ .  $\mathbb{E}_{\alpha \sim \pi(\cdot; \theta)} [\cdot]$  is the expectation over the searched architectures.

**Algorithm 1** Training method of PNAG.

---

**Require:** Search space  $\Omega$ , latency distribution  $\mathcal{B}$ , learning rate  $\eta$ , training data set  $\mathcal{D}$ , parameters  $M$ ,  $N$ , and  $K$ .

- 1: Initialize model parameters  $\theta$  for the generator and  $w$  for the architecture evaluator.  
*// Collect the architectures with accuracy and latency*
- 2: Train a supernet  $S$  on  $\mathcal{D}$ .
- 3: Randomly sample architectures  $\{\beta_i\}_{i=1}^M$  from  $\Omega$ .
- 4: Construct tuples  $\{(\beta_i, c(\beta_i), \text{Acc}(\beta_i))\}_{i=1}^M$  using  $S$ .  
*// Learn the architecture evaluator*
- 5: **while** not convergent **do**
- 6:     Sample a set of latencies  $\{B_k\}_{k=1}^K$  from  $\mathcal{B}$ .
- 7:     Update the architecture evaluator by:
- 8:      $w \leftarrow w - \eta \nabla_w L(w)$ .
- 9: **end while**  
*// Learn the architecture generator*
- 10: **while** not convergent **do**
- 11:     Sample a set of latencies  $\{B_k\}_{k=1}^K$  from  $\mathcal{B}$ .
- 12:     Obtain  $\{\alpha_{B_k}^{(i)}\}_{i=1}^N$  from  $\pi(\cdot|B_k; \theta)$  for each  $B_k$ .
- 13:     Update the generator via policy gradient by:
- 14:      $\theta \leftarrow \theta + \eta \nabla_{\theta} J(\theta)$ .
- 15: **end while**

---

**NAS under diverse budgets.** Problem (1) only focuses on one specific budget constraint. In fact, we seek to learn the Pareto frontier over the whole range of budgets (*e.g.*, latency). However, this problem is hard to solve since there may exist infinite Pareto optimal architectures with different computational cost. To address this, one can learn an approximated Pareto frontier by finding a set of uniformly distributed Pareto optimal points [67]. Here, we evenly sample  $K$  budgets from the range of latency and maximize the expected reward over them. Thus, the problem becomes

$$\begin{aligned} \max_{\theta} \mathbb{E}_{B \sim \mathcal{B}} \left[ \mathbb{E}_{\alpha_B \sim \pi(\cdot|B; \theta)} [R(\alpha_B|B; w)] \right], \\ \text{s.t. } c(\alpha_B) \leq B, B \sim \mathcal{B}, \end{aligned} \quad (2)$$

where  $\mathbb{E}_{B \sim \mathcal{B}} [\cdot]$  denotes the expectation over the distribution of budget. Unlike Eqn. (1),  $\pi(\cdot|B; \theta)$  is the learned policy conditioned on the budget of  $B$ . In practice, we use policy gradient to learn the architecture generator. To encourage exploration, we follow [34], [68], [69] to introduce an entropy regularization. Please refer to the supplementary materials for more details.

**Advantages over existing NAS methods.** Our PNAG exhibits two advantages over existing NAS methods. *First*, our PNAG is able to share the learned knowledge across the search processes under different budgets, which greatly improves the search results (see Table 5). The main reason is that, once we find a good architecture for one budget, we may easily obtain a competitive architecture for a larger/smaller budget by slightly modifying some components (model width or kernel size). *Second*, given a well-trained PNAG, we can directly use it to generate feasible architectures for any required budget via inference, which is very efficient and practically useful (see Table 6).

**3.1.2 Vector Representation of Budget Bounds**

To learn the architecture generator, we still have to consider how to represent the budget bound  $B$  as the inputs of PNAG. As mentioned before, our PNAG considers  $K$  discrete budgets during training. To represent different budgets, we use an embedding vector [34] to represent different budgets (See details in Section 3.1.2). Following [34], we build a learnable embedding vector  $\mathbf{b} = g(B)$  for each sampled budget  $B$ . We incorporate these learnable embedding vectors into the parameters of the

architecture generator and train them jointly. In this way, we are able to automatically learn the vectors of these budgets and encourage PNAG to produce feasible architectures.

As mentioned before, we only sample a set of discrete budgets to train PNAG. To accommodate all the budgets belonging to a continuous space, we propose an embedding interpolation method to represent a budget with any possible value. Specifically, we perform a linear interpolation between the embedding of two adjacent discrete budgets to represent the considered budgets. For a target budget  $B$  between two sampled budgets  $B_1 < B < B_2$ , the linear interpolation of the budget vector  $\mathbf{b}$  can be computed by

$$\mathbf{b} = g(B) = \xi g(B_1) + (1 - \xi)g(B_2), \text{ where } \xi = \frac{B_2 - B}{B_2 - B_1},$$

Here,  $\xi \in [0, 1]$  denotes the weight of  $B_1$  in interpolation.

**3.2 Learning the Architecture Evaluator  $R(\cdot|B; w)$**

Given diverse budgets, an architecture should have different rewards/scores regarding whether it satisfies the corresponding budget constraint. However, it is non-trivial to manually design a reward function for each budget. Instead, we propose to learn an architecture evaluator to automatically predict the score. To this end, we build an evaluator with three fully connected layers. Given any architecture  $\beta$  and a budget  $B$ , we seek to predict the performance  $R(\beta|B; w)$  of  $\beta$  under the budget  $B$ . Since we have no ground-truth labels for training, following [70], [71], [72], we learn the evaluator via pairwise architecture comparisons.

**3.2.1 Training Method of  $R(\cdot|B; w)$**

To obtain a promising evaluator, we train the architecture evaluator using a pairwise ranking loss, which has been widely used in ranking problems [70], [71], [72]. Specifically, we collect  $M$  architectures with accuracy and latency, and record them as a set of triplets  $\{(\beta_i, c(\beta_i), \text{Acc}(\beta_i))\}_{i=1}^M$ . Thus, given  $M$  architectures, we have  $M(M-1)$  architecture pairs  $\{(\beta_i, \beta_j)\}$  in total after omitting the pairs with themselves. Assuming that we have  $K$  budgets, the pairwise ranking loss becomes

$$\begin{aligned} L(w) = \frac{1}{KM(M-1)} \sum_{k=1}^K \sum_{i=1}^M \sum_{j=1, j \neq i}^M \phi \left( d(\beta_i, \beta_j, B_k) \right. \\ \left. \cdot [R(\beta_i|B_k; w) - R(\beta_j|B_k; w)] \right), \end{aligned} \quad (3)$$

where  $d(\beta_1, \beta_2, B_k)$  denotes a function to indicate whether  $\beta_i$  is better than  $\beta_j$  under the budget  $B_k$ , as will be discussed in Section 3.2.2.  $\phi(z) = \max(0, 1 - z)$  is a hinge loss function and we use it to enforce the predicted ranking results  $R(\beta_i|B_k; w) - R(\beta_j|B_k; w)$  to be consistent with the results of  $d(\beta_i, \beta_j, B_k)$  obtained by a comparison rule based on Pareto dominance.

**3.2.2 Pareto Dominance Rule**

To compare the performance between two architectures, we need to define a reasonable function  $d(\beta_1, \beta_2, B)$  in Eqn. (3). To this end, we define a Pareto dominance to guide the design of this function. Specifically, Pareto dominance requires that the quality of an architecture should depend on both the satisfaction of budget and accuracy. That means, given a specific budget  $B$ , a good architecture should be the one with the cost lower than or equal to  $B$  and with high accuracy. In this sense, we use Pareto dominance to compare two architectures and judge which one is dominative.

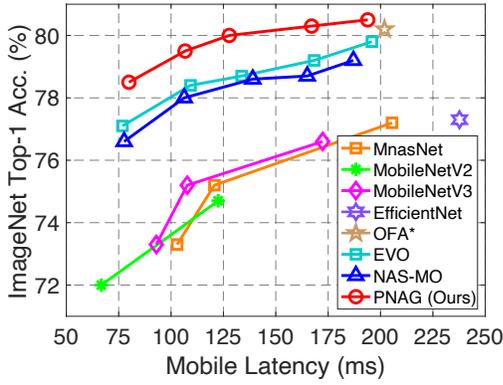


Fig. 3. Comparisons of the architectures obtained by different methods on a mobile device (Qualcomm Snapdragon 821).

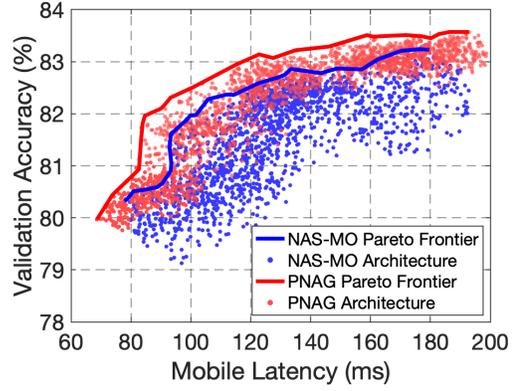


Fig. 4. Comparisons of the Pareto frontiers of the generated architectures between NAS-MO and PNAG. Here, we report the accuracy evaluated on the constructed validation set.

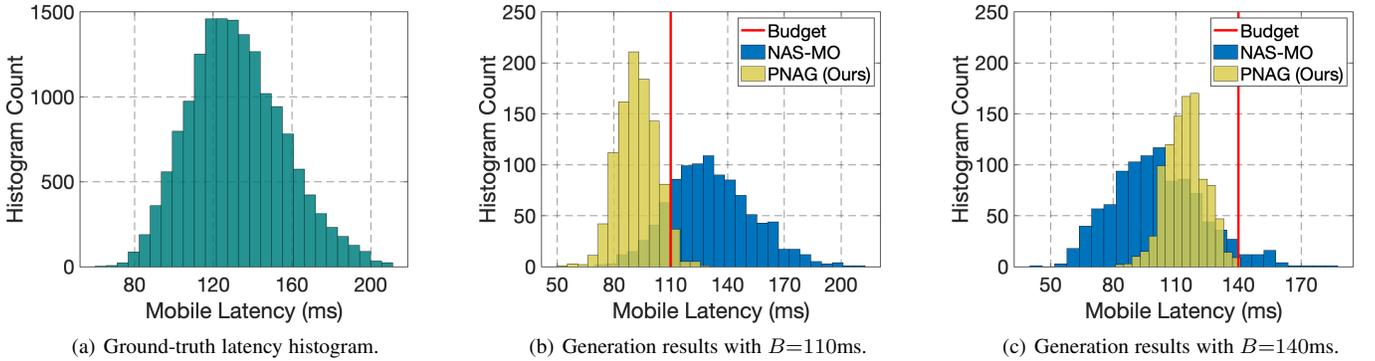


Fig. 5. Latency histograms of sampled architectures on mobile devices. (a) Ground-truth latency histogram of 16,000 architectures that are uniformly sampled from the search space. (b) The latency histogram of 1,000 architectures sampled by different methods given  $B=110$ ms. (c) The latency histogram of 1,000 architectures sampled by different methods given  $B=140$ ms.

Given any two architectures  $\beta_1, \beta_2$ , if both of them satisfy the budget constraints (namely  $c(\beta_1) \leq B$  and  $c(\beta_2) \leq B$ ), then  $\beta_1$  dominates  $\beta_2$  if  $\text{Acc}(\beta_1) \geq \text{Acc}(\beta_2)$ . Moreover, when at least one of  $\beta_1, \beta_2$  violates the budget constraint, clearly we have that  $\beta_1$  dominates  $\beta_2$  if  $c(\beta_1) \leq c(\beta_2)$ . Formally, we define the Pareto dominance function  $d(\beta_1, \beta_2, B)$  to reflect the above rules:

$$d(\beta_1, \beta_2, B) = \begin{cases} 1, & \text{if } (c(\beta_1) \leq B \wedge c(\beta_2) \leq B) \\ & \wedge (\text{Acc}(\beta_1) \geq \text{Acc}(\beta_2)); \\ -1, & \text{else if } (c(\beta_1) \leq B \wedge c(\beta_2) \leq B) \\ & \wedge (\text{Acc}(\beta_1) < \text{Acc}(\beta_2)); \\ 1, & \text{else if } c(\beta_1) \leq c(\beta_2); \\ -1, & \text{otherwise.} \end{cases} \quad (4)$$

Based on Eqn. (4), we have  $d(\beta_1, \beta_2, B) = -d(\beta_2, \beta_1, B)$  if  $\beta_1 \neq \beta_2$ , making it a symmetric metric w.r.t.  $\beta_1$  and  $\beta_2$ .

**Remark 1** The accuracy constraint  $\text{Acc}(\beta_1) \geq \text{Acc}(\beta_2)$  plays an important role in the proposed Pareto dominance function  $d(\beta_1, \beta_2, B)$ . Without the accuracy constraint, we may easily find the architectures with very low computation cost and poor performance (See results in Table 4).

## 4 EXPERIMENTS

We apply the proposed PNAG to produce architectures under diverse latency budgets evaluated on three hardware platforms, including a mobile device (equipped with a Qualcomm Snapdragon 821 processor), a CPU processor (Intel Core i5-7400), and a GPU card (NVIDIA TITAN X). For convenience, we use “Architecture- $B$ ” to represent the generated architecture that satisfies the latency budget w.r.t.  $B$ , e.g., PNAG-80. Our code and all the pretrained models are available at <https://github.com/guoyongcs/PNAG>.

### 4.1 Implementation Details

**Search space.** Following [55], we use MobileNetV3 [73] as the backbone to build the search space [55], [56]. We divide a network into several units. To find promising architectures, we allow each unit to have 1) any numbers of layers (*i.e.*, depth) chosen from  $\{2, 3, 4\}$ , 2) any width expansion ratios in each layer (*i.e.*, width) chosen from  $\{3, 4, 6\}$ , and 3) any kernel sizes chosen from  $\{3, 5, 7\}$ . We build the model with 5 units. Thus, there are  $3 \times 3$  combinations of widths and kernel sizes for each layer.

**Training the supernet.** To accelerate the training of supernet, we follow [76] to randomly choose 100 classes from original 1000 classes in ImageNet for training and train the supernet with progressive shrinking strategy [55] for 90 epochs. We treat 80% of these data as the training set to train the supernet and the rest

TABLE 1

Comparisons with state-of-the-art architectures on mobile devices. \* denotes the best architecture reported in the original paper. “-” denotes the results that are not reported. All the models are evaluated on  $224 \times 224$  images of ImageNet.

Architecture	Latency (ms)	Test Accuracy (%)		#Params (M)	#MAdds (M)	Search Cost (GPU Days)
		Top-1	Top-5			
MobileNetV3-Small (1.0 $\times$ ) [73]	39.8	67.4	-	2.4	56	-
MobileNetV3-Large (0.75 $\times$ ) [73]	93.0	73.3	-	4.0	155	-
MobileNetV2 (1.0 $\times$ ) [74]	90.3	72.0	-	3.4	300	-
FBNetV2 [75]	-	76.3	92.9	-	321	30.0
MnasNet-A1 (0.5 $\times$ ) [24]	37.5	68.9	88.4	2.1	105	-
SGNAS-B [61]	-	76.8	-	-	326	0.3
EVO-80	76.8	77.1	93.3	6.1	350	0.7
NAS-MO-80	77.6	76.6	93.2	7.9	340	0.7
PNAG-80	79.9	<b>78.3</b>	<b>94.0</b>	7.3	349	0.7
FBNet-A [76]	91.7	73.0	-	4.3	249	9.0
SGNAS-A [61]	-	77.1	-	-	373	0.3
ProxylessNAS-Mobile [77]	97.3	74.6	-	4.1	319	8.3
ProxylessNAS-CPU [77]	98.5	75.3	-	4.4	438	8.3
MobileNetV3-Large (1.0 $\times$ ) [73]	107.7	75.2	-	5.4	219	-
EVO-110	109.3	78.4	94.0	10.2	482	0.7
NAS-MO-110	106.3	78.0	93.8	8.4	478	0.7
PNAG-110	106.8	<b>79.4</b>	<b>94.5</b>	9.9	451	0.7
RandWire [53]	-	74.7	92.2	5.6	583	-
ProxylessNAS-GPU [77]	123.3	75.1	-	7.1	463	8.3
MnasNet-A1 (1.0 $\times$ ) [24]	120.7	75.2	92.5	3.4	300	$\sim$ 3792
FBNet-C [76]	135.2	74.9	-	5.5	375	9.0
EVO-140	133.7	78.7	94.1	9.1	488	0.7
NAS-MO-140	139.0	78.6	94.0	9.5	486	0.7
PNAG-140	127.8	<b>79.8</b>	<b>94.7</b>	9.2	492	0.7
NSGANetV1 [66]	-	76.2	93.0	5.0	585	27
PONAS-C [56]	145.1	75.2	-	5.6	376	8.8
P-DARTS [45]	168.7	75.6	92.6	4.9	577	3.8
BigNAS-L [78]	-	79.5	-	6.4	586	1.5
EVO-170	168.3	79.2	94.4	10.7	661	0.7
NAS-MO-170	165.0	78.7	94.4	8.5	584	0.7
PNAG-170	167.1	<b>80.3</b>	<b>95.0</b>	10.0	606	0.7
NSGANetV2 [27]	-	79.1	94.5	8.0	400	1
NAGO [54]	-	76.8	93.4	5.7	-	20.0
PC-DARTS [46]	194.1	75.8	92.7	5.3	597	0.1
EfficientNet B0 [79]	237.7	77.3	93.5	5.3	390	-
Cream-L [80]	-	80.0	94.7	9.7	604	12
OFA* [55]	201.9	80.2	95.1	9.1	743	51.7
EVO-200	195.9	79.8	94.5	11.0	783	0.7
NAS-MO-200	187.4	79.2	94.4	9.1	630	0.7
PNAG-200	193.9	<b>80.5</b>	<b>95.2</b>	10.4	724	0.7

20% as the validation set to measure the validation accuracy of candidate architectures (we report such validation accuracy in Figs. 1(b) and 4). We consider the original ImageNet validation set as the test data and report the test accuracy of candidate architectures on them in all the other tables and figures. Based on a NVIDIA V100 GPU, the training process of the supernet takes around *15 GPU hours* (*i.e.*, 0.6 GPU days).

**Training the architecture evaluator.** We collect  $M=16,000$  architectures by uniformly sampling architectures from the search space  $\Omega$  (See Fig. 5(a)) following [55] and obtain the latency ranges on three hardware devices. We deploy these architectures to different devices and measure the latency over a batch of images. Specifically, we measure the latency on mobile and CPU devices with a batch size of 1. Since the inference on GPU is too fast to obtain the accurate latency, we measure the latency with a batch size of 64 on NVIDIA TITAN X. We compute the accuracy  $\text{Acc}(\cdot)$  on our validation set (*i.e.*, 20% samples of 100 selected classes in

ImageNet). We train the architecture evaluator for 250 epochs. The learning rate is initialized to 0.1 and decreased to  $1 \times 10^{-3}$  with a cosine annealing. Following [55], we train two predictors to predict the latency and validation accuracy, respectively. We set the dimension of the embedding vector of budgets to 64. We emphasize that training the architecture evaluator is very efficient and only takes less than *0.2 GPU hours*.

**Training the architecture generator.** We train the model for 120k iterations using an Adam optimizer with a learning rate of  $3 \times 10^{-4}$ . Following ENAS [34], we sample  $N=1$  architecture at each iteration and find it works well in practice. We select  $K=10$  latency budgets by evenly dividing the range. We add an entropy regularization term to the reward weighted by  $1 \times 10^{-3}$ . Note that training the architecture generator approximately takes *2 GPU hours*. When evaluating the searched architectures, following [40], [55], we first obtain the parameters from the OFA full network and then finetune them for 75 epochs to obtain the final performance.

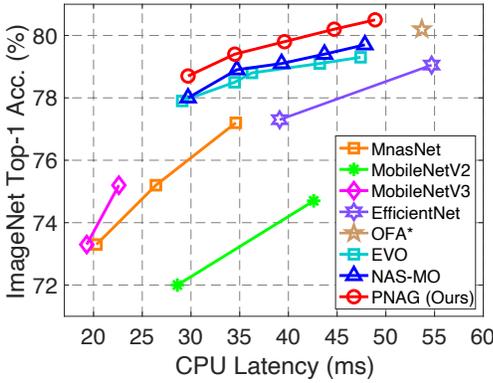


Fig. 6. Comparisons of the architectures obtained by different methods on a Core i5-7400 CPU.

### 4.2 Compared Methods

To investigate the effectiveness of the proposed method, we compare our PNAG with two variants: 1) **EVO** uses the evolutionary search method [39] to perform architecture search. 2) **NAS-MO** conducts architecture search based by exploiting the multi-objective reward [24]. We also compare our method with several state-of-the-art methods, including ENAS [34], DARTS [44], P-DARTS [45], PC-DARTS [46], MNasNet [24], MobileNetV2 [74], MobileNetV3 [73], FBNet [76], FBNetV2 [75], Proxyless-NAS [77], EfficientNet [79], OFA [39], Cream [80], PONAS [56], BigNAS [78], and NSGANetV2 [27].

### 4.3 Architecture Search for Mobile Devices

In this experiment, we train our PNAG to produce feasible architectures for the latency budgets based on a mobile device (Qualcomm Snapdragon 821 processor). Based on the proposed budget interpolation method in Section 3.1.2, our PNAG is flexible to generate feasible architectures for any arbitrary budget. To evaluate our method, for simplicity, we manually choose 5 latency budgets {80ms, 110ms, 140ms, 170ms, 200ms} and reports the results under each of them. The other budgets are also possible.

We compare our PNAG with state-of-the-art methods given different latency budgets evaluated on the considered mobile device. In Fig. 3, we compare the architectures searched by different methods in terms of both accuracy and latency. We draw the following conclusions. *First*, our PNAG (red line) consistently generates better architectures than the considered variants EVO and NAS-MO under diverse budgets. *Second*, our best architecture (the rightmost point of the red line) yields a better trade-off between accuracy and latency than a strong baseline OFA\*, *i.e.*, the best architecture reported in [55]. For convenience, we put more detailed comparison results in Table 1. Given diverse latency budgets, our PNAG greatly outperforms the compared NAS methods in terms of the accuracy of the generated/searched architectures. Specifically, our PNAG-200 yields the best accuracy of 80.5, which is better than the best reported results in OFA [55], namely OFA\*. We also highlight that, besides the superior performance, the total training cost of our PNAG is about 0.7 GPU days<sup>1</sup>, which is much more efficient than most state-of-the-art NAS methods, such as [55], [78], [80].

1. We report the training cost of each component of PNAG in Section 4.1

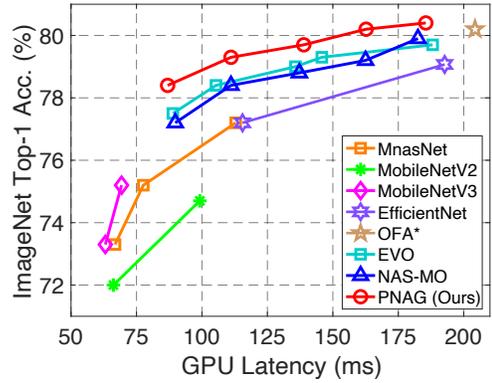


Fig. 7. Comparisons of the architectures obtained by different methods on a NVIDIA TITAN X GPU.

Moreover, we compare the learned/searched frontiers of different methods and show the comparisons of Pareto frontiers in Fig. 4. We plot all the architectures produced by different methods to form the Pareto frontier. Specifically, we use the architectures searched by multiple independent runs under different budgets for NAS-MO. For PNAG, we use linear interpolation to generate architectures that satisfy different budgets. From Fig. 4, our PNAG finds a better frontier than NAS-MO due to the shared knowledge across the search process under different budgets. We also visualize the latency histograms of the architectures evaluated on mobile devices in Fig. 5(b) and Fig. 5(c). Given latency budgets of 110ms and 140ms, NAS-MO is prone to produce a large number of architectures that cannot satisfy the target budgets. These results show that it is hard to design the multi-objective reward to obtain the preferred architectures. Instead, PNAG uses the Pareto dominance reward to encourage the architectures to satisfy the desired budget constraints. In this sense, most architectures generated by our PNAG are able to fulfill the target budgets. We put more visual results of latency histograms *w.r.t.* other latency budgets in the supplementary.

### 4.4 Architecture Search for CPU Devices

We further exploit our PNAG to generate architectures under the latency budgets evaluated on a CPU device (Core i5-7400). Similar to the experiments for mobile devices, we evaluate our PNAG under 5 latency budgets, *i.e.*, {30ms, 35ms, 40ms, 45ms, 50ms}.

As shown in Fig. 6, our PNAG yields a large performance improvement over the considered two variants, *i.e.*, EVO and NAS-MO, under diverse budgets. Moreover, our PNAG also outperforms popular NAS based (MnasNet, OFA\*) and manually designed architectures (MobileNetV2, MobileNetV3, and EfficientNet). As for the quantitative comparisons, in Table 2, our PNAG consistently yields the best results across all the considered latency budgets. To be specific, given a small latency budget  $B=35ms$ , our PNAG-35 yields better accuracy than the compared NAS methods with much lower search cost. Given a relatively large budget  $B=50ms$ , our PNAG-50 yields the same accuracy (80.5%) as the best result on mobile devices (*i.e.*, PNAG-200). This indicates that our PNAG generalizes well across the latency budgets based on different hardware. Overall, these results demonstrate that our PNAG is able to generate very competitive architectures while satisfying diverse latency budgets.

TABLE 2

Comparisons with state-of-the-art architectures on Intel Core i5-7400 CPU. \* denotes the best architecture reported in the original paper. “-” denotes the results that are not reported. All the models are evaluated on  $224 \times 224$  images of ImageNet.

Architecture	Latency (ms)	Test Accuracy (%)		#Params (M)	#MAdds (M)	Search Cost (GPU Days)
		Top-1	Top-5			
MobileNetV2 (1.0 $\times$ ) [74]	28.6	72.0	-	3.4	300	-
MobileNetV3-Large (1.0 $\times$ ) [73]	22.6	75.2	-	5.4	219	-
FBNet-C [76]	25.7	74.9	-	5.5	375	9.0
SGNAS-B [61]	-	76.8	-	-	326	0.3
EVO-30	29.1	77.9	93.8	7.9	385	0.7
NAS-MO-30	29.7	77.5	93.7	6.6	353	0.7
PNAG-30 (Ours)	29.7	<b>78.3</b>	<b>94.1</b>	7.6	335	0.7
ProxylessNAS-CPU [77]	34.6	75.3	-	4.4	438	8.3
MnasNet-A1 (1.4 $\times$ ) [24]	34.6	77.2	93.5	6.1	592	$\sim$ 3792
EVO-35	34.5	78.5	94.3	8.2	354	0.7
NAS-MO-35	34.7	78.3	94.0	7.9	478	0.7
PNAG-35 (Ours)	34.5	<b>79.4</b>	<b>94.5</b>	8.4	431	0.7
ResNet-18 [4]	38.6	69.8	90.1	11.7	1814	-
EfficientNet B0 [79]	39.1	77.3	93.5	5.3	390	-
EVO-40	36.3	78.8	94.6	8.4	388	0.7
NAS-MO-40	39.3	78.6	94.3	8.3	491	0.7
PNAG-40 (Ours)	39.6	<b>79.8</b>	<b>94.9</b>	9.4	502	0.7
MobileNetV2 (1.4 $\times$ ) [74]	42.6	74.7	-	6.9	585	-
EVO-45	43.2	79.1	94.6	9.1	481	51.7
NAS-MO-45	43.7	78.8	94.4	9.3	626	0.7
PNAG-45 (Ours)	44.7	<b>80.2</b>	<b>95.0</b>	10.4	620	0.7
PONAS-C [56]	52.2	75.2	-	5.6	376	8.8
OFA* [55]	53.7	80.2	95.1	9.1	743	51.7
EVO-50	47.4	79.3	94.7	9.1	511	0.7
NAS-MO-50	46.7	78.9	94.4	9.1	632	0.7
PNAG-50 (Ours)	48.9	<b>80.5</b>	<b>95.1</b>	10.5	682	0.7

TABLE 3

Comparisons with state-of-the-art architectures on NVIDIA TITAN X GPU. \* denotes the best architecture reported in the original paper. “-” denotes the results that are not reported. All the models are evaluated on  $224 \times 224$  images of ImageNet.

Architecture	Latency (ms)	Test Accuracy (%)		#Params (M)	#MAdds (M)	Search Cost (GPU Days)
		Top-1	Top-5			
ProxylessNAS-GPU [77]	84.7	75.1	-	7.1	463	8.3
MobileNetV2 (1.0 $\times$ ) [74]	71.6	72.0	-	3.4	300	-
NAGO [54]	-	76.8	93.4	5.7	-	20.0
EVO-90	88.9	77.3	93.1	5.9	332	0.7
NAS-MO-90	89.8	75.4	92.4	4.9	266	0.7
PNAG-90 (Ours)	86.9	<b>78.3</b>	<b>94.0</b>	5.7	310	0.7
MnasNet-A1 (1.4 $\times$ ) [24]	112.9	77.2	93.5	6.1	592	$\sim$ 3792
EfficientNet B0 [79]	115.5	77.3	93.5	5.3	390	-
ENAS [34]	110.8	73.8	91.7	5.6	607	0.5
EVO-115	105.4	78.4	94.1	8.4	388	51.7
NAS-MO-115	111.2	78.1	94.0	8.8	431	0.7
PNAG-115 (Ours)	111.2	<b>79.3</b>	<b>94.6</b>	8.9	411	0.7
EVO-140	135.7	78.9	94.4	9.1	481	0.7
NAS-MO-140	137.2	78.4	94.1	8.8	470	0.7
PNAG-140 (Ours)	138.9	<b>79.7</b>	<b>94.9</b>	9.7	510	0.7
ResNet-50 [4]	159.8	76.2	92.9	25.6	4087	-
EVO-165	164.1	79.1	94.5	10.7	597	51.7
NAS-MO-165	162.6	78.8	94.4	10.5	583	0.7
PNAG-165 (Ours)	162.7	<b>80.3</b>	<b>95.0</b>	10.5	582	0.7
NASNet-A [16]	162.3	74.0	91.6	5.3	564	$\sim$ 3
PONAS [56]	182.4	75.2	-	5.6	376	8.8
EfficientNet B1 [79]	192.7	79.2	94.5	7.8	700	-
OFA* [55]	204.3	80.2	95.1	9.1	743	51.7
EVO-190	188.1	79.5	94.8	11.3	687	0.7
NAS-MO-190	183.2	78.8	94.5	10.7	652	0.7
PNAG-190 (Ours)	185.5	<b>80.4</b>	<b>95.0</b>	10.4	640	0.7

TABLE 4  
Comparisons of different reward functions based on PNAG. We report the latency on mobile devices.

Reward	$B_1=80\text{ms}$		$B_2=110\text{ms}$		$B_3=140\text{ms}$		$B_4=170\text{ms}$		$B_5=200\text{ms}$	
	Acc. (%)	Lat. (ms)	Acc.	Lat. (ms)						
Multi-objective Reward [24]	77.0	77.6	78.5	106.3	78.9	139.0	79.3	165.1	79.5	187.3
Multi-objective Absolute Reward [58]	78.1	76.8	78.9	109.2	79.2	130.1	79.5	163.6	79.9	197.5
Pareto Dominance Reward (w/o acc. constraint)	73.8	74.4	73.6	64.9	74.3	66.5	73.9	70.0	74.0	70.8
Pareto Dominance Reward (Ours)	<b>78.4</b>	79.9	<b>79.5</b>	106.8	<b>79.8</b>	127.8	<b>80.3</b>	167.1	<b>80.5</b>	193.9

TABLE 5  
Effect of different search strategies on the performance of PNAG. We report the accuracy on ImageNet.

Search Strategy	$B_1=80\text{ms}$	$B_2=110\text{ms}$	$B_3=140\text{ms}$	$B_4=170\text{ms}$	$B_5=200\text{ms}$
Repeated Independent Search	76.7	78.6	79.1	79.4	79.7
Pareto Frontier Search	<b>78.4</b>	<b>79.5</b>	<b>79.8</b>	<b>80.3</b>	<b>80.5</b>

TABLE 6  
Comparisons of the time cost for architecture generation/design among different methods.

Method	PNAG	PC-DARTS	ENAS	DARTS
Time Cost	$\leq 5$ s	2 hours	12 hours	4 days

TABLE 7  
Effect of  $K$  on the generation performance of PNAG. We compare the generated architectures using different values of  $K$  with the target latency  $B=140\text{ms}$  on ImageNet.

$K$	1	2	5	10	30
Top-1 Acc. (%)	78.5	79.1	79.4	<b>79.8</b>	<b>79.8</b>

#### 4.5 Architecture Search for GPU Devices

Besides the mobile and CPU devices, we also consider GPUs and adopt the latency on them as the computational budget. Since the inference speed on GPU is much faster than mobile processor and CPU, we measure the latency of deep models on a NVIDIA TITAN X GPU with a batch size of 64. In this experiments, we compare different architecture design/search methods under the budgets of {90ms, 115ms, 140ms, 165ms, 190ms}.

As shown in Fig. 7, similar to the results on mobile and CPU devices, our PNAG outperforms existing methods and the constructed variants by a large margin. We also reported the detailed comparisons in terms of accuracy and computational cost in Table 3. Again, compared with both the hand-crafted methods (e.g., MobileNetV2 [74] and EfficientNet [79]) and NAS methods (e.g., ENAS [34] and MnasNet [24]), our PNAG consistently produces better architectures under diverse budgets. These results further emphasize the generalization ability of our PNAG to the latency budgets evaluated on different hardware devices.

### 5 FURTHER EXPERIMENTS

In this section, we conduct ablation studies on our method. Then we compare the architecture generation cost of our proposed method among different methods and discuss the impact of the number of considered budgets  $K$ .

#### 5.1 Effect of the Pareto Dominance Reward

We investigate the effectiveness of the Pareto frontier learning strategy and the Pareto dominance reward. From Table 4 and Table 5, the Pareto frontier learning strategy tends to find better architectures than the independent search process due to the shared knowledge across the search processes under different budgets. Compared with two existing multi-objective rewards [24], [58], the Pareto dominance reward encourages the generator to produce architectures that satisfy the considered budget constraints. Moreover, if we do not consider accuracy constraint in the Pareto

dominance reward, the generated architectures have low latency and poor accuracy. With both the Pareto frontier learning strategy and the Pareto dominance reward, our method yields the best results under all budgets.

#### 5.2 Comparisons of Architecture Generation Cost

In this part, we compare the architecture generation cost of different methods for 5 different budgets and show the comparison results in Table 6. Given an arbitrary target budget, existing NAS methods need to perform an independent search to find feasible architectures. By contrast, since PNAG directly learns the whole Pareto frontier, we are able to generate promising architectures based on a learned generator model via *inference*. Thus, the architecture generation cost of PNAG is much less than other existing methods (See results in Table 6). In this sense, we are able to greatly accelerate the architecture design process in real-world scenarios. These results demonstrate the efficiency of our PNAG in generating architectures.

#### 5.3 Effect of $K$ on the Generation Performance

We investigate the effect of  $K$  on the generation performance of PNAG based on mobile device. Note that we evenly select  $K$  budgets from the range of latency. To investigate the effect of  $K$ , we consider several candidate values of  $K \in \{2, 5, 10, 30\}$ . We show the Top-1 accuracies of the architectures generated by PNAG with different  $K$  on ImageNet in Table 7. Since a small number of selected budgets  $K$  cannot accurately approximate the ground-truth Pareto frontier or provide enough shared knowledge between different search processes, our method yields poor results with  $K = 2$ . When we increase  $K$  larger than 5, we are able to greatly improve the performance of the generated architectures. From Table 7, our method yields the best result when  $K \geq 10$  and we use this setting in the experiments.

## 6 CONCLUSION

In this paper, we focus on designing effective and feasible architectures via an architecture generation process. To this end, we have proposed a novel Pareto-aware Neural Architecture Generator (PNAG) which only needs to be trained once and dynamically generates promising architectures satisfying any given budget via inference. Unlike existing methods, we seek to learn the whole Pareto frontier instead of finding a single or several discrete Pareto optimal architectures. Based on the learned Pareto frontier, our PNAG consistently outperforms existing NAS methods across diverse budgets. Extensive experiments on three hardware platforms (*i.e.*, mobile devices, CPU, and GPU) demonstrate the effectiveness of the proposed method.

## REFERENCES

- [1] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel, "Backpropagation applied to handwritten zip code recognition," *Neural Computation*, vol. 1, no. 4, pp. 541–551, 1989.
- [2] 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.
- [3] R. K. Srivastava, K. Greff, and J. Schmidhuber, "Training very deep networks," in *Advances in Neural Information Processing Systems*, 2015, pp. 2377–2385.
- [4] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [5] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," in *International Conference on Learning Representations*, 2021.
- [6] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *IEEE International Conference on Computer Vision*, 2021, pp. 9992–10002.
- [7] E. Shelhamer, J. Long, and T. Darrell, "Fully convolutional networks for semantic segmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 4, pp. 640–651, 2017.
- [8] L. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *European Conference on Computer Vision*, V. Ferrari, M. Hebert, C. Sminchisescu, and Y. Weiss, Eds., vol. 11211, 2018, pp. 833–851.
- [9] E. Xie, W. Wang, Z. Yu, A. Anandkumar, J. M. Alvarez, and P. Luo, "Segformer: Simple and efficient design for semantic segmentation with transformers," *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [10] J. Wang, K. Sun, T. Cheng, B. Jiang, C. Deng, Y. Zhao, D. Liu, Y. Mu, M. Tan, X. Wang, W. Liu, and B. Xiao, "Deep high-resolution representation learning for visual recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 43, no. 10, pp. 3349–3364, 2021.
- [11] J. Redmon and A. Farhadi, "Yolov3: An incremental improvement," *arXiv preprint arXiv:1804.02767*, 2018.
- [12] Z.-Q. Zhao, P. Zheng, S.-t. Xu, and X. Wu, "Object detection with deep learning: A review," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 30, no. 11, pp. 3212–3232, 2019.
- [13] Z. Tian, C. Shen, H. Chen, and T. He, "Fcos: Fully convolutional one-stage object detection," in *IEEE International Conference on Computer Vision*, 2019, pp. 9627–9636.
- [14] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *European Conference on Computer Vision*, A. Vedaldi, H. Bischof, T. Brox, and J. Frahm, Eds., vol. 12346, 2020, pp. 213–229.
- [15] B. Zoph and Q. V. Le, "Neural architecture search with reinforcement learning," in *International Conference on Learning Representations*, 2017.
- [16] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le, "Learning transferable architectures for scalable image recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8697–8710.
- [17] C. Li, J. Peng, L. Yuan, G. Wang, X. Liang, L. Lin, and X. Chang, "Block-wisely supervised neural architecture search with knowledge distillation," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1989–1998.
- [18] M. Tan, R. Pang, and Q. V. Le, "Efficientdet: Scalable and efficient object detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10778–10787.
- [19] X. Dai, A. Wan, P. Zhang, B. Wu, Z. He, Z. Wei, K. Chen, Y. Tian, M. Yu, P. Vajda, and J. E. Gonzalez, "Fbnetv3: Joint architecture-recipe search using predictor pretraining," in *IEEE Conference on Computer Vision and Pattern Recognition*, June 2021, pp. 16276–16285.
- [20] C. White, A. Zela, R. Ru, Y. Liu, and F. Hutter, "How powerful are performance predictors in neural architecture search?" *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [21] W. Chen, X. Gong, and Z. Wang, "Neural architecture search on imagenet in four gpu hours: A theoretically inspired perspective," *arXiv preprint arXiv:2102.11535*, 2021.
- [22] Z. Yan, X. Dai, P. Zhang, Y. Tian, B. Wu, and M. Feiszli, "Fp-nas: Fast probabilistic neural architecture search," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 15139–15148.
- [23] Y. Guo, Y. Chen, Y. Zheng, P. Zhao, J. Chen, J. Huang, and M. Tan, "Breaking the curse of space explosion: Towards efficient nas with curriculum search," in *International Conference on Machine Learning*, 2020.
- [24] M. Tan, B. Chen, R. Pang, V. Vasudevan, M. Sandler, A. Howard, and Q. V. Le, "Mnasnet: Platform-aware neural architecture search for mobile," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2820–2828.
- [25] 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.
- [26] 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, pp. 419–427.
- [27] Z. Lu, K. Deb, E. Goodman, W. Banzhaf, and V. N. Boddeti, "Nsganetv2: Evolutionary multi-objective surrogate-assisted neural architecture search," in *European Conference on Computer Vision*. Springer, 2020, pp. 35–51.
- [28] I. Y. Kim and O. L. De Weck, "Adaptive weighted-sum method for bi-objective optimization: Pareto front generation," *Structural and Multidisciplinary Optimization*, vol. 29, no. 2, pp. 149–158, 2005.
- [29] C. He, H. Ye, L. Shen, and T. Zhang, "Milenas: Efficient neural architecture search via mixed-level reformulation," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11990–11999.
- [30] G. Li, G. Qian, I. C. Delgadillo, M. Müller, A. K. Thabet, and B. Ghanem, "SGAS: sequential greedy architecture search," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1617–1627.
- [31] T.-J. Yang, Y.-L. Liao, and V. Sze, "Netadaptv2: Efficient neural architecture search with fast super-network training and architecture optimization," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 2402–2411.
- [32] X. Zhang, Z. Huang, N. Wang, S. Xiang, and C. Pan, "You only search once: Single shot neural architecture search via direct sparse optimization," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 43, no. 9, pp. 2891–2904, 2021.
- [33] X. Zheng, R. Ji, Y. Chen, Q. Wang, B. Zhang, J. Chen, Q. Ye, F. Huang, and Y. Tian, "MIGO-NAS: towards fast and generalizable neural architecture search," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 43, no. 9, pp. 2936–2952, 2021.
- [34] H. Pham, M. Guan, B. Zoph, Q. Le, and J. Dean, "Efficient neural architecture search via parameter sharing," in *International Conference on Machine Learning*, 2018, pp. 4095–4104.
- [35] R. Pasunuru and M. Bansal, "Continual and multi-task architecture search," in *Proceedings of Conference of the Association for Computational Linguistics*, A. Korhonen, D. R. Traum, and L. Màrquez, Eds., 2019, pp. 1911–1922.
- [36] Y. Tian, Q. Wang, Z. Huang, W. Li, D. Dai, M. Yang, J. Wang, and O. Fink, "Off-policy reinforcement learning for efficient and effective GAN architecture search," in *European Conference on Computer Vision*, A. Vedaldi, H. Bischof, T. Brox, and J. Frahm, Eds., vol. 12352, 2020, pp. 175–192.
- [37] A. Vahdat, A. Mallya, M. Liu, and J. Kautz, "UNAS: differentiable architecture search meets reinforcement learning," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11263–11272.

- [38] 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 *International Conference on Machine Learning*, 2017, pp. 2902–2911.
- [39] E. Real, A. Aggarwal, Y. Huang, and Q. V. Le, "Regularized evolution for image classifier architecture search," in *AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 4780–4789.
- [40] Z. Lu, G. Sreeksumar, E. Goodman, W. Banzhaf, K. Deb, and V. N. Boddeti, "Neural architecture transfer," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 43, no. 9, pp. 2971–2989, 2021.
- [41] M. Lin, P. Wang, Z. Sun, H. Chen, X. Sun, Q. Qian, H. Li, and R. Jin, "Zen-nas: A zero-shot nas for high-performance deep image recognition," in *IEEE International Conference on Computer Vision*, 2021.
- [42] M. Chen, H. Peng, J. Fu, and H. Ling, "Autoformer: Searching transformers for visual recognition," in *IEEE International Conference on Computer Vision*, 2021, pp. 12 250–12 260.
- [43] Y. Liu, Y. Sun, B. Xue, M. Zhang, G. G. Yen, and K. C. Tan, "A survey on evolutionary neural architecture search," *IEEE transactions on neural networks and learning systems*, 2021.
- [44] H. Liu, K. Simonyan, and Y. Yang, "Darts: Differentiable architecture search," in *International Conference on Learning Representations*, 2019.
- [45] X. Chen, L. Xie, J. Wu, and Q. Tian, "Progressive differentiable architecture search: Bridging the depth gap between search and evaluation," in *IEEE International Conference on Computer Vision*, 2019, pp. 1294–1303.
- [46] Y. Xu, L. Xie, X. Zhang, X. Chen, G.-J. Qi, Q. Tian, and H. Xiong, "Pc-darts: Partial channel connections for memory-efficient differentiable architecture search," in *International Conference on Learning Representations*, 2020.
- [47] X. Chu, X. Wang, B. Zhang, S. Lu, X. Wei, and J. Yan, "DARTS-: robustly stepping out of performance collapse without indicators," in *International Conference on Learning Representations*, 2021.
- [48] X. Chen, R. Wang, M. Cheng, X. Tang, and C. Hsieh, "Dnras: Dirichlet neural architecture search," in *International Conference on Learning Representations*, 2021.
- [49] Y. Guo, J. Wang, Q. Chen, J. Cao, Z. Deng, Y. Xu, J. Chen, and M. Tan, "Towards lightweight super-resolution with dual regression learning," *arXiv preprint arXiv:2207.07929*, 2022.
- [50] Y. Zhao, L. Wang, Y. Tian, R. Fonseca, and T. Guo, "Few-shot neural architecture search," in *International Conference on Machine Learning*. PMLR, 2021, pp. 12 707–12 718.
- [51] Y. Xu, Y. Wang, K. Han, Y. Tang, S. Jui, C. Xu, and C. Xu, "Renas: Relativistic evaluation of neural architecture search," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 4411–4420.
- [52] X. Chu, B. Zhang, and R. Xu, "Fairnas: Rethinking evaluation fairness of weight sharing neural architecture search," in *IEEE International Conference on Computer Vision*. IEEE, 2021, pp. 12 219–12 228.
- [53] S. Xie, A. Kirillov, R. Girshick, and K. He, "Exploring randomly wired neural networks for image recognition," in *IEEE International Conference on Computer Vision*, 2019, pp. 1284–1293.
- [54] B. Ru, P. Esperanca, and F. Carlucci, "Neural architecture generator optimization," in *Advances in Neural Information Processing Systems*, 2020.
- [55] H. Cai, C. Gan, and S. Han, "Once for all: Train one network and specialize it for efficient deployment," in *International Conference on Learning Representations*, 2020.
- [56] S.-Y. Huang and W.-T. Chu, "Ponas: Progressive one-shot neural architecture search for very efficient deployment," *arXiv preprint arXiv:2003.05112*, 2020.
- [57] J. H. M. Thomas Elsken and F. Hutter, "Efficient multi-objective neural architecture search via lamarckian evolution," in *International Conference on Learning Representations*, 2019.
- [58] G. Bender, H. Liu, B. Chen, G. Chu, S. Cheng, P. Kindermans, and Q. V. Le, "Can weight sharing outperform random architecture search? an investigation with tunas," in *IEEE International Conference on Computer Vision*, 2020, pp. 14 311–14 320.
- [59] 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," in *European Conference on Computer Vision*, 2020, pp. 544–560.
- [60] C. Li, Z. Yu, Y. Fu, Y. Zhang, Y. Zhao, H. You, Q. Yu, Y. Wang, and Y. Lin, "Hw-nas-bench: Hardware-aware neural architecture search benchmark," *arXiv preprint arXiv:2103.10584*, 2021.
- [61] S.-Y. Huang and W.-T. Chu, "Searching by generating: Flexible and efficient one-shot nas with architecture generator," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 983–992.
- [62] 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.
- [63] M. Kim, T. Hiroyasu, M. Miki, and S. Watanabe, "SPEA2+: improving the performance of the strength pareto evolutionary algorithm 2," in *Parallel Problem Solving from Nature*, vol. 3242. Springer, 2004, pp. 742–751.
- [64] A.-C. Cheng, J.-D. Dong, C.-H. Hsu, S.-H. Chang, M. Sun, S.-C. Chang, J.-Y. Pan, Y.-T. Chen, W. Wei, and D.-C. Juan, "Searching toward pareto-optimal device-aware neural architectures," in *Proceedings of the International Conference on Computer-Aided Design*, 2018, pp. 1–7.
- [65] J.-D. Dong, A.-C. Cheng, D.-C. Juan, W. Wei, and M. Sun, "Dpp-net: Device-aware progressive search for pareto-optimal neural architectures," in *European Conference on Computer Vision*, 2018, pp. 517–531.
- [66] 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*, 2020.
- [67] C. Grosan and A. Abraham, "Generating uniformly distributed pareto optimal points for constrained and unconstrained multicriteria optimization," *International Conference on Informatics and Systems*, pp. 27–29, 2008.
- [68] 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*, 2019.
- [69] Y. Guo, Y. Zheng, M. Tan, Q. Chen, Z. Li, J. Chen, P. Zhao, and J. Huang, "Towards accurate and compact architectures via neural architecture transformer," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [70] Y. Freund, R. Iyer, R. E. Schapire, and Y. Singer, "An efficient boosting algorithm for combining preferences," *Journal of Machine Learning Research*, vol. 4, no. Nov, pp. 933–969, 2003.
- [71] C. Burges, T. Shaked, E. Renshaw, A. Lazier, M. Deeds, N. Hamilton, and G. Hullender, "Learning to rank using gradient descent," in *International Conference on Machine Learning*, 2005, pp. 89–96.
- [72] W. Chen, T.-Y. Liu, Y. Lan, Z.-M. Ma, and H. Li, "Ranking measures and loss functions in learning to rank," in *Advances in Neural Information Processing Systems*, 2009, pp. 315–323.
- [73] A. Howard, M. Sandler, G. Chu, L.-C. Chen, B. Chen, M. Tan, W. Wang, Y. Zhu, R. Pang, V. Vasudevan *et al.*, "Searching for mobilenetv3," in *IEEE International Conference on Computer Vision*, 2019, pp. 1314–1324.
- [74] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4510–4520.
- [75] A. Wan, X. Dai, P. Zhang, Z. He, Y. Tian, S. Xie, B. Wu, M. Yu, T. Xu, K. Chen *et al.*, "Fbnetv2: Differentiable neural architecture search for spatial and channel dimensions," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 12 965–12 974.
- [76] 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 *IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 10 734–10 742.
- [77] H. Cai, L. Zhu, and S. Han, "ProxylessNAS: Direct neural architecture search on target task and hardware," in *International Conference on Learning Representations*, 2019.
- [78] J. Yu, P. Jin, H. Liu, G. Bender, P.-J. Kindermans, M. Tan, T. Huang, X. Song, R. Pang, and Q. Le, "Bignas: Scaling up neural architecture search with big single-stage models," in *European Conference on Computer Vision*. Springer, 2020, pp. 702–717.
- [79] M. Tan and Q. V. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," in *International Conference on Machine Learning*, 2019, pp. 6105–6114.
- [80] H. Peng, H. Du, H. Yu, Q. Li, J. Liao, and J. Fu, "Cream of the crop: Distilling prioritized paths for one-shot neural architecture search," in *Advances in Neural Information Processing Systems*, 2020.