

DCP-NAS: Discrepant Child-Parent Neural Architecture Search for 1-bit CNNs

Yanjing Li[†] · Sheng Xu[†] · Xianbin Cao* · Li'an Zhuo · Baochang Zhang* · Tian Wang · Guodong Guo

Received: date / Accepted: date

Abstract Neural architecture search (NAS) proves to be among the effective approaches for many tasks by generating an application-adaptive neural architecture, which is still challenged by high computational cost and memory consumption. At the same time, 1-bit convolutional neural networks (CNNs) with binary weights and activations show their potential for resource-limited embedded devices. One natural approach is to use 1-bit CNNs to reduce the computation and memory cost of NAS by taking advantage of the strengths of each in a unified framework, while searching the 1-bit CNNs is more challenging due to the more complicated processes involved. In this paper, we introduce *Discrepant Child-Parent Neural Architecture Search (DCP-NAS)* to efficiently search 1-bit CNNs, based on a new framework of searching the 1-bit model (Child) under the supervision of a real-valued model (Parent). Particularly, we first utilize a Parent model to calculate a tangent direction, based on which the tangent propagation method is introduced to search the optimized 1-bit Child. We further observe a coupling relationship between the weights and architecture parameters existing in such differentiable frameworks. To address the issue, we propose a decoupled optimization method to search an optimized architecture. Extensive experiments demonstrate that our DCP-NAS achieves much better results than prior arts on both CIFAR-10 and ImageNet datasets. In particular, the backbones achieved by our DCP-NAS

Yanjing Li, Sheng Xu, Xianbin Cao, Li'an Zhuo and Tian Wang
Beihang University, Beijing
E-mail: {yanjingli, shengxu, xbcao, lianzhuo, bczhang, wangtian}@buaa.edu.cn
Baochang Zhang
Beihang University, Beijing; Zhongguancun Laboratory, Beijing; Nanyang Institute of Technology, Nanyang
E-mail: bczhang@buaa.edu.cn
Guodong Guo
UNIUBI Research, Universal Ubiquitous Co., Hangzhou, Zhejiang, CHINA.
E-mail: guodong.guo@mail.wvu.edu
[†] Co-first authors with equal contribution.
* Corresponding authors.

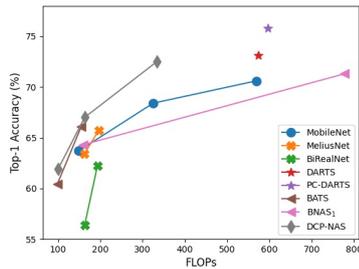


Fig. 1 Comparing DCP-NAS with lightweight CNNs and BNAS on Computational cost vs. ImageNet Accuracy. DCP-NAS achieves the best performance-cost trade-off.

Network Design	Model	32-bit	1-bit	Gap
Hand-crafted	ResNet-18	69.6	65.9	-3.7
	ResNet-34	73.3	69.0	-4.3
	ReActNet-A	72.4	69.4	-3.0
Direct BNAS	BNAS ₁	67.1	64.3	-2.9
	BNAS ₂	64.2	59.8	-4.4
	BNAS ₂ v2	69.5	66.0	-3.5
Auxiliary BNAS	CP-NAS	68.0	66.5	-1.5
	DCP-NAS	68.5	72.4	+4.0

Table 1 Accuracy of both hand-crafted and search-based architectures with real-valued and their counterpart 1-bit CNNs on ImageNet. Note that BNAS₁ [9, 10], BNAS₂ [31] and BNAS₂ v2 [32] are different works.

achieve strong generalization performance on person re-identification and object detection.

Keywords Binary neural network · Neural architecture search · Tangent propagation

1 Introduction

Neural architecture search (NAS) has attracted great attention with a remarkable performance in many computer vision tasks such as classification [81, 56] and detection [12, 57]. NAS attempts to design network architectures automatically and replace conventional hand-crafted counterparts at the expense of searching from a huge search space and a high computational cost in training. Early efforts include sharing weights between searched and newly generated networks [6]. ProxylessNAS [8] introduces a differentiable latency loss into NAS to search architectures on the target task instead of adopting the conventional proxy-based framework. PC-DARTS [67] takes advantage of partial channel connections to improve the memory efficiency by sampling a small part of the supernet to reduce the redundancy of the network search space. EfficientNet [56] introduces a new scaling method that uniformly scales all dimensions of depth, width and resolution using a simple yet highly effective compound coefficient to obtain efficient networks. Besides these efficient NAS methods, developing effective NAS for efficient 1-bit convolutional neural networks (1-bit CNNs) [41, 66, 63, 64, 62] has drawn increasing attention, which retains the advantages of 1-bit CNNs on memory saving and computational cost reduction [9, 3]. 1-bit CNNs directly compress real-valued weights and activations (32-bit) of CNNs into a single bit and directly decrease the memory consumption by $32\times$ and the computation cost by up to $58\times$ [53].

Comparatively speaking, 1-bit CNNs based on hand-crafted architectures have been extensively researched. Binarized filters have been used in conven-

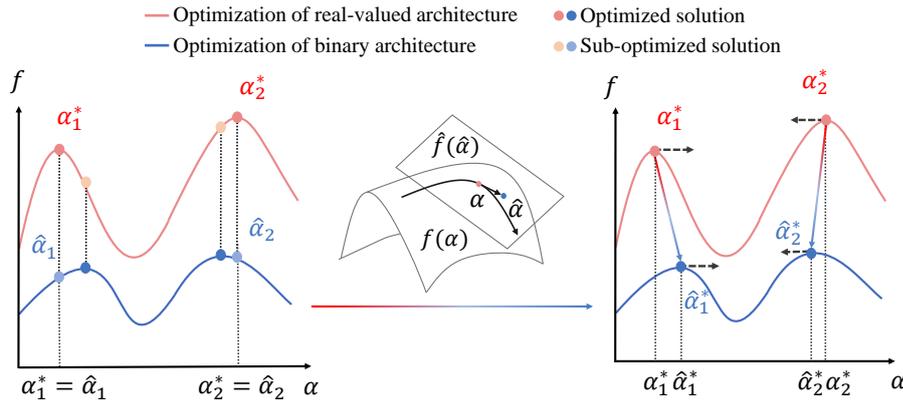


Fig. 2 Motivation for DCP-NAS. We first show directly binarizing real-valued architecture to 1-bit is sub-optimal. Thus we use tangent propagation (middle) to find an optimized 1-bit neural architecture along the tangent direction, leading to a better-performed 1-bit neural architecture.

tional CNNs to compress deep models [53, 14, 13, 30], which is widely considered as one of the most efficient ways to perform computing on embedded devices with low computational cost. In [30], the XNOR network is presented where both the weights and inputs attached to the convolution are approximated with binarized values. This results in an efficient implementation of convolutional operations by reconstructing the unbinarized filters with a single scaling factor. In [21], a projection convolutional neural network (PCNN) is proposed to implement 1-bit CNNs based on a simple back propagation algorithm. [72] proposes Bayesian optimized 1-bit CNNs, taking the advantages of Bayesian learning to significantly improve the performance of extreme 1-bit CNNs. 1-bit CNNs show the advantages on computational cost reduction and memory saving, however, there still exists a significant gap (about 3%~4%) between 1-bit hand-crafted models and real-valued counterparts, as shown in Tab. 1.

Binary neural architecture search (BNAS) is a simple and promising way to search the 1-bit network architectures. Previous BNAS methods are to directly search binary architecture [9, 31], *i.e.*, direct BNAS. BNAS₁ [9] attempts to search binary architectures from well-designed binary search spaces. Likewise, BNAS₂ [31] utilizes the diversity in the early search to learn better performing 1-bit neural architectures. However, such search strategy is only effective if an exhaustive search process is exploited, as revealed in [3]. Likewise, we observe a performance gap about 3%~4% compared with the real-valued counterparts from such direct BNAS, as depicted in Tab. 1.

In this paper, we first introduce a Child-Parent framework to efficiently search a 1-bit CNN in a unified framework. The real-valued models converge much faster than 1-bit models as revealed in [45], which motivate us to use the tangent direction of Parent supernet (real-valued model) as an indicator of the optimization direction for Child supernet (1-bit model).

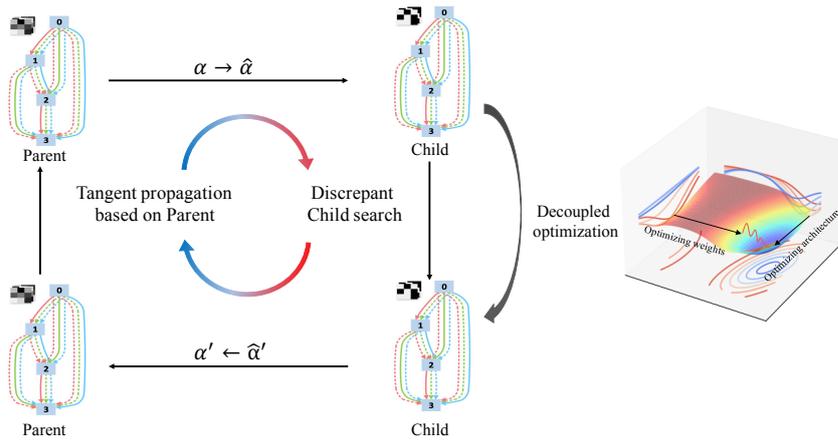


Fig. 3 The main framework of the proposed DCP-NAS, where α and $\hat{\alpha}$ denote real-valued and binary architecture, respectively. In one single round, we first conduct the real-valued NAS and generate the corresponding tangent direction. Then we learn a discrepant binary architecture via tangent propagation. In this process, real-valued and 1-bit CNNs inherit architectures from their counterparts in turn.

We assume that all the possible 1-bit neural architectures can be learnt from the tangent space of the Parent model, based on which we introduce a *Discrepant Child-Parent Neural Architecture Search* (DCP-NAS) method to produce an optimized 1-bit CNN. Specifically, as shown in Fig. 2, we utilize the Parent model to find a tangent direction to learn the 1-bit Child through tangent propagation, rather than directly binarizing the Parent to Child. Since the tangent direction is based on the second order information, we further accelerate the search process by Generalized Gauss-Newton matrix (GGN), leading to an efficient search process. Moreover, a coupling relationship between the weights and architecture parameters exists in such DARTS-based [43] methods, leading to an asynchronous convergence and insufficient training process. To overcome this obstacle, we propose a decoupled optimization for training the Child-Parent model, leading to an effective and optimized search process. The overall framework of our DCP-NAS is shown in Fig. 3. This paper is an extended version of our conference paper [80], where we make the following new contributions:

- 1) We propose a new Discrepant Child-Parent model (DCP-NAS) to guide the binary architecture search, which further explore the discrepancy between real-valued and binary architectures.
- 2) We utilize the Parent model to generate a tangent direction for Child, which learns the 1-bit neural architecture discrepancy through tangent propagation. The GGN method is further introduced to reduce the computation cost in the DCP-NAS optimization.

3) We introduce a new decoupled optimization to address the asynchronous convergence in such differentiable NAS process, which further improve the performance of our DCP-NAS.

4) We add a group of comprehensive experiments including new application on person re-identification and object detection tasks, which validate the acceleration, efficiency, and practicability of the proposed method.

2 Related Work

Extensive work has been reported on compressing and accelerating neural networks through quantization [53, 37, 61, 37, 63, 64], low-rank approximation [71, 48], network pruning [35, 22, 24], and more recently, architecture search [42, 82]. Amongst them, quantization and architecture search are most relevant to our work and are reviewed in this section.

2.1 1-bit CNNs

BNN [14] first binarized the weights and activations of convolutional layers in CNNs. BinaryNet based on BinaryConnect [13] is introduced to train DCNNs with binary weights, where activations are triggered at the run-time while parameters are computed at the training time [29]. In [1], the straight-through estimator (STE) is introduced to optimize BNNs. In the forward propagation, the sign function is directly applied to obtain binary parameters. In [53], XNOR-Net is presented where both the weights and inputs of the convolution kernels are approximated with binary values, to improve the efficiency of the convolution operations. Instead of directly binarizing the weights into ± 1 , XNOR-Net reconstruct the binary weights with the mean absolute value (MAV) of each layer as the optimal value for α_l . XNOR++ [5] further designs an efficient estimator for layer-wise feature maps of the BNNs. Based on XNOR-Net, Bi-Real Net [47] designs a new variant of residual structure to preserve the real activations before the sign function and a magnitude-aware gradient scheme *w.r.t.* the weight to update the binarized weight parameters. Instead of using a scale factor to estimate the real-valued kernel weights from the binary weights, HQRQ [38] adopts a high-order binarization method to achieve a more accurate approximation while retaining the advantage of the binarized weights. Inspired by approximation of the real-valued weights helping the performance of BNNs, ABC-Net [41] introduces multiple binary weights and activations for better estimation of the real-valued parameters to alleviate the degradation in prediction accuracy. [21] introduces a quantization method based on a discrete back-propagation algorithm to learn a better BNNs. Furthermore, BONNs [72] extends a Bayesian method to the prior distributions of real-valued weights, features, and filters for constructing a BNNs in a comprehensive end-to-end manner, which further improves the performance. ReActNet [46] replaces the conventional PReLU [25] and the sign

function of the BNNs with RPreLU and RSign with a learnable threshold, thus improving the performance of BNNs. LWS-Det [66] effectively propose the angular and amplitude minimization technology to narrow the gap between real-valued kernels and 1-bit kernels. [68] explores the effect of dead weights, which refer to a group of weights, and proposes the rectified clamp unit (ReCU) to revive the dead weights for updating.

2.2 Neural Architecture Search

Thanks to the rapid development of deep learning, significant gains in performance have been realized in a wide range of computer vision tasks, most of which are manually designed network architectures [34, 55, 26, 28]. Recently, the new approach called neural architecture search (NAS) has been attracting increased attention. The goal is to find automatic ways of designing neural architectures to replace conventional hand-crafted ones. Existing NAS approaches need to explore a very large search space and can be roughly divided into three types of approaches: evolution based, reinforcement-learning-based and one-shot-based.

In order to implement the architecture search within a short period of time, researchers try to reduce the cost of evaluating each searched candidate. Early efforts include sharing weights between searched and newly generated networks [6]. Later, this method was generalized into a more elegant framework named one-shot architecture search [2, 7, 43, 51, 60, 75]. In these approaches, an over-parameterized network or super network covering all candidate operations is trained only once, and the final architecture is obtained by sampling from this super network. For example, [2] trained the over-parameterized network using a HyperNet [23], and [51] proposed to share parameters among Child models to avoid retraining each candidate from scratch. DARTS [43] introduces a differentiable framework and thus combines the search and evaluation stages into one. Despite its simplicity, researchers have found some of its drawbacks and proposed a few improved approaches over DARTS [60, 11]. PDARTS [11] presents an efficient algorithm which allows the depth of searched architectures to grow gradually during the training procedure, with a significantly reduced search time. ProxylessNAS [8] adopted the differentiable framework and proposed to search architectures on the target task instead of adopting the conventional proxy-based framework. IDARTS [69] focuses on decoupled optimization of NAS. In particular, IDARTS [69] decouple and backtrack the edge and operation parameters, in which the constraint $R(\cdot)$ of the edge parameter is utilized for controlling backtracking.

Binary neural architecture search replaces real-valued weights and activations with binarized ones, which consumes much less memory and computation resource to search 1-bit CNNs and provides a more promising way to efficiently find network architectures. These methods can be categorized into two classes, *direct binary architecture search* and *auxiliary binary architecture search*. Direct binary architecture search yields binary architectures directly from

well-designed binary search spaces. As one of the pioneer arts of this field, BNAS_1 [9] effectively reduces the search time by channel sampling and search space pruning in the early training stages for a differentiable NAS. Another pioneer art BNAS_2 [31] utilizes the diversity in the early search to learn better performing binary architectures. BMES [52] learns an efficient binary MobileNet [27] architecture via the evolution-based search. However, the accuracy of direct binary architecture search can be improved by auxiliary binary architecture search [3]. BATS [3] designs a novel search space specially tailored to the 1-bit CNNs and incorporates it into the DARTS framework. NASB [79] explores an optimal architecture and connection for a convolutional group with introducing a three stage training scheme into BinaryNAS. EBN [4] propose Expert Binary Convolution to select data-specific expert binary and design a principled 1-bit CNNs search mechanism that obtain a set of network architectures of preferred properties. BARS [73] design a joint search strategy for both macro-level and micro-level search space tailored for 1-bit CNNs and analyze the information bottlenecks relating to both the topology and layout architecture design choices.

Different from the aforementioned methods, our work is driven by the performance discrepancy between the 1-bit neural architecture and its real-valued counterpart. We introduce tangent propagation to explore the accuracy discrepancy, and further accelerate the search process by applying the GGN to the Hessian matrix in the optimization. Furthermore, we introduce a novel decoupled optimization to address the asynchronous convergence in such a differentiable NAS process, leading to better performed 1-bit CNNs. The overall framework leads to a novel and effective BNAS process.

3 Discrepant Child-Parent NAS

In this section, we first give the preliminaries of our method. Then, we describe our Child-Parent framework and the search space for binary NAS. Then we present the proposed architecture discrepancy-aware binary neural architecture search strategy together with tangent propagation. Finally, the decoupled optimization for solving the coupling relationship between weights and architecture parameters is introduced for effectively searching optimized binary neural architecture. For clarity, Tab. 2 describes the main notations used in the following sections.

3.1 Search Space

We search for computation cells as the building blocks of the final architecture. Following [81, 82, 43, 54], we construct the network with a pre-defined number of cells and each cell is a fully-connected directed acyclic graph (DAG) \mathcal{G} with N nodes. For simplicity, we assume that each cell only takes the outputs of the two previous cells as input and each input node has pre-defined convolutional

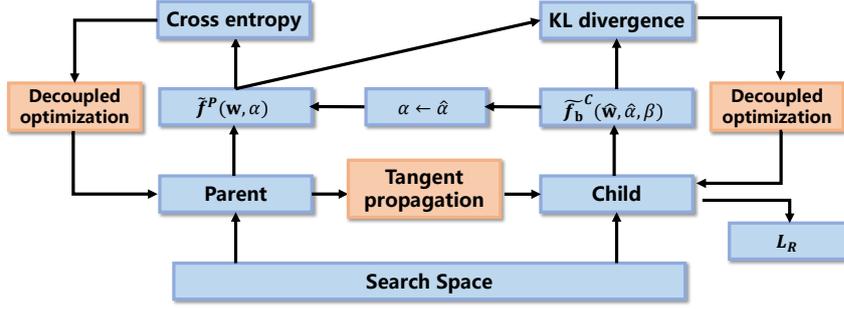


Fig. 4 The main framework of Discrepant Child-Parent model. We show the key novelty of DCP-NAS, *i.e.*, tangent propagation and decoupled optimization, in orange.

operations for preprocessing. Each node j is obtained by

$$\begin{aligned} \mathbf{a}^{(j)} &= \sum_{i < j} o^{(i,j)}(\mathbf{a}^{(i)}) \\ o^{(i,j)}(\mathbf{a}^i) &= \mathbf{w}^{(i,j)} \otimes \mathbf{a}^i, \end{aligned} \quad (1)$$

where i is the dependent nodes of j with the constraints $i < j$ to avoid cycles in a cell, and \mathbf{a}^j is the output of the node j node. $\mathbf{w}^{(i,j)}$ denotes the weights of the convolution operation between i -th and j -th nodes and \otimes denotes the convolution operation. Each node is a specific tensor like a feature map, and each directed edge (i, j) denotes an operation $o^{(i,j)}(\cdot)$, which is sampled from following $M = 8$ operations:

- no connection (zero)
- skip connection (identity)
- 3×3 dilated convolution with rate 2
- 5×5 dilated convolution with rate 2
- 3×3 max pooling
- 3×3 average pooling
- 3×3 depth-wise separable convolution
- 5×5 depth-wise separable convolution

We replace the depth wise separable convolution with a binarized form, *i.e.*, binary weights and activations. Note that, skip connection is identity mapping in NAS, instead of additional shortcut. The optimization of BNNs is more challenging than that of conventional CNNs [21, 53], since binarization brings additional burdens to NAS. Following [43], to reduce the undesired fluctuation in the performance evaluation, we normalize the architecture parameter of M operations for each edge to obtain the final architecture indicator as

$$\hat{o}_m^{(i,j)}(\mathbf{a}^{(j)}) = \frac{\exp\{\alpha_m^{(i,j)}\}}{\sum_{m'} \exp\{\alpha_{m'}^{(i,j)}\}} o_m^{(i,j)}(\mathbf{a}^{(j)}). \quad (2)$$

Table 2 A brief description of the main notations used in our paper.

$\mathbf{w}, \hat{\mathbf{w}}$: real-valued weight	$\mathbf{a}, \hat{\mathbf{a}}$: real-valued activation
$\mathbf{b}^{\hat{\mathbf{w}}}$: binary weight	$\mathbf{b}^{\hat{\mathbf{a}}}$: binary activation
β : scale factor	$\alpha, \hat{\alpha}$: real-valued and 1-bit architecture parameter
$f(\cdot)$: real-valued NAS objective	$\tilde{f}(\cdot)$: relaxed real-valued NAS objective
$f_b(\cdot)$: BNAS objective	$\tilde{f}_b(\cdot)$: relaxed BNAS objective
$p_n(\cdot), \hat{p}_n(\cdot)$: output logits	i, j : node index
$o^{(i,j)}(\cdot)$: operation	\mathcal{L}_{NAS} : searching objective of NAS
$\mathcal{L}_{\text{CP-NAS}}$: searching objective of CP-NAS	$\mathcal{L}_{\text{DCP-NAS}}$: searching objective of DCP-NAS
$G(\cdot)$: objective of Child model training	$\mathbf{D}(\cdot)$: objective of tangent propagation
$\mathcal{L}_R(\cdot)$: reconstruction error	$\mathbf{H}(\cdot)$: Hessian matrix
$\mathbf{I}(\cdot)$: Fisher information matrix	$R(\cdot)$: norm constraint in decoupled optimization

3.2 Preliminary

3.2.1 Neural architecture search (NAS)

Given a conventional CNN model, we denote $\mathbf{w} \in \mathcal{W}$ and $\mathcal{W} = \mathbb{R}_{C_{out} \times C_{in} \times K \times K}$ and $\mathbf{a}_{in} \in \mathbb{R}_{C_{in} \times W \times H}$ as its weights and feature maps in the specific layer. C_{out} and C_{in} represent the number of output and input channels of the specific layer. (W, H) is the feature maps' width and height, and K is the kernel size. We then have

$$\mathbf{a}_{out} = \mathbf{a}_{in} \otimes \mathbf{w}, \quad (3)$$

where \otimes is the convolution operation. We omit the batch normalization (BN) and activation layers for simplicity. Based on this, a normal NAS problem is given as

$$\max_{\mathbf{w} \in \mathcal{W}, \alpha \in \mathcal{A}} f(\mathbf{w}, \alpha), \quad (4)$$

where $f: \mathcal{W} \times \mathcal{A} \rightarrow \mathbb{R}$ is a differentiable objective function *w.r.t.* the network weight $\mathbf{w} \in \mathcal{W}$ and the architecture space $\mathcal{A} \in \mathbb{R}_{M \times E}$, where E and M denote the number of edges and operators, respectively. Considering that minimizing $f(\mathbf{w}, \alpha)$ is a black-box optimization, we relax the objective function into $\tilde{f}(\mathbf{w}, \alpha)$ as the objective of NAS

$$\begin{aligned} \min_{\mathbf{w} \in \mathcal{W}, \alpha \in \mathcal{A}} \mathcal{L}_{\text{NAS}} &= -\tilde{f}(\mathbf{w}, \alpha) \\ &= -\sum_{n=1}^N p_n(\mathcal{X}) \log(p_n(\mathbf{w}, \alpha)), \end{aligned} \quad (5)$$

where N denotes the number of classes and \mathcal{X} is the input data. $\tilde{f}(\mathbf{w}, \alpha)$ represents the performance of a specific architecture with real-valued weights, where $p_n(\mathcal{X})$ and $p_n(\mathbf{w}, \alpha)$ denote the true distribution and the distribution of network prediction, respectively.

3.2.2 Binary neural architecture search (BNAS)

The 1-bit model aims to quantize $\hat{\mathbf{w}}$ and $\hat{\mathbf{a}}_{in}$ into $\mathbf{b}^{\hat{\mathbf{w}}} \in \{-1, +1\}_{C_{out} \times C_{in} \times K \times K}$ and $\mathbf{b}^{\hat{\mathbf{a}}_{in}} \in \{-1, +1\}_{C_{in} \times H \times W}$ using the efficient XNOR and Bit-count operations to replace real-valued operations. More specifically, all 1-bit CNNs are trained using the fake quantization strategy, *i.e.*, the weights and activations are represented by $\{-1, +1\}$ in FP32 format using $\text{sign}(\cdot)$ during training, as mentioned in [20]. Efficient XNOR and Bit-count operations are used in hardware deployment. Following [13, 29], the forward process of the 1-bit CNN is

$$\hat{\mathbf{a}}_{out} = \beta \circ \mathbf{b}^{\hat{\mathbf{a}}_{in}} \odot \mathbf{b}^{\hat{\mathbf{w}}}, \quad (6)$$

where \odot is the XNOR, and bit-count operations and \circ denotes the channel-wise multiplication. $\beta = [\beta_1, \dots, \beta_{C_{out}}] \in \mathbb{R}_{C_{out}}^+$ is the vector consisting of channel-wise scale factors. $\mathbf{b} = \text{sign}(\cdot)$ denotes the binarized variable using the sign function, which returns 1 if the input is greater than zero, and -1 otherwise. It then enters several non-linear layers, *e.g.*, BN layer, non-linear activation layer, and max-pooling layer. We omit these for simplification. Then, the output $\hat{\mathbf{a}}_{out}$ is binarized to $\mathbf{b}^{\hat{\mathbf{a}}_{out}}$ via the sign function. The fundamental objective of BNNs is calculating $\hat{\mathbf{w}}$. We want it to be as close as possible before and after binarization, such that the binarization effect is minimized. Then, we define the reconstruction error following [21] as

$$\mathcal{L}_R(\hat{\mathbf{w}}, \beta) = \|\hat{\mathbf{w}} - \beta \circ \mathbf{b}^{\hat{\mathbf{w}}}\|_2^2. \quad (7)$$

Based on the derivation above, the vanilla direct BNAS [9, 31] can be defined as

$$\max_{\hat{\mathbf{w}} \in \mathcal{W}, \hat{\alpha} \in \mathcal{A}, \beta \in \mathbb{R}^+} f_{\mathbf{b}}(\hat{\mathbf{w}}, \hat{\alpha}, \beta), \quad (8)$$

where $\mathbf{b}^{\hat{\mathbf{w}}} = \text{sign}(\hat{\mathbf{w}})$ is used for inference and $\hat{\alpha}$ is neural architecture with binary weights. Prior direct BNAS [9] learning the BNAS from such objective as

$$\max_{\hat{\mathbf{w}} \in \mathcal{W}, \hat{\alpha} \in \mathcal{A}, \beta \in \mathbb{R}^+} \tilde{f}_{\mathbf{b}}(\hat{\mathbf{w}}, \hat{\alpha}, \beta) = \sum_{n=1}^N \hat{p}_n(\hat{\mathbf{w}}, \hat{\alpha}, \beta) \log(\hat{p}_n(\mathcal{X})), \quad (9)$$

where we use similar denotations as Eq. 5. Eq. 9 means that vanilla direct BNAS only focus on the binary search space under the supervision of cross entropy loss, which is less effective due to the search process is not exhaustive [3].

3.3 Child-Parent Framework for Network Binarization

Network binarization, which calculates neural networks with binary weights and activations to fit the real-valued network, can significantly compress the CNNs. Prior work [72] usually investigates the binarization problem by exploring the real-valued model to guide the optimization of 1-bit CNNs. Based on the investigation, we reformulate NAS-based network binarization as

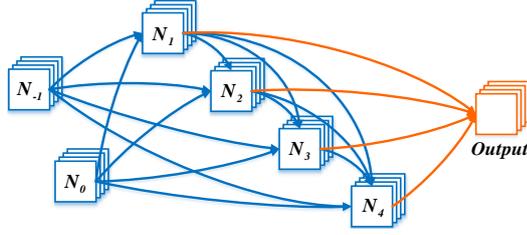


Fig. 5 The cell architecture for DCP-NAS. One cell includes 2 input nodes, 4 intermediate nodes and 14 edges (blue).

a Child-Parent model as shown in Fig. 4. The 1-bit CNNs and the real-valued counterpart are the Child and Parent model, respectively.

Conventional NAS is inefficient due to the complicated reward computation in network training where the evaluation of a structure is usually done after the network training converges. There are also some methods to perform the evaluation of a cell during the training of the network. [76] points out that the best choice in early stages is not necessarily the final optimal one, however, the worst operation in the early stages usually has a bad performance in the end. This phenomenon will become more and more significant as the training goes. Based on this observation, we propose a simple yet effective operation removing process, which is the key task of the proposed CP-model. In a searching loop, we first sample the operation without replacement for each edge from the search space, and then train the sampled model generated by the corresponding supernet. Second, we validate our model on the validation set. Until all the operations are selected, we remove the operation in each edge with the worst performance, *i.e.*, in terms of the accuracy of both models on the validation dataset.

Intuitively, the representation difference between the Children and Parent, and how much Children can independently handle their problems, are two main aspects that should be considered to define a reasonable performance evaluation measure. Based on these analysis, we introduce the Child-Parent framework for binary NAS, which defines the objective as

$$\begin{aligned}
 \hat{\mathbf{w}}^*, \hat{\alpha}^*, \beta^* &= \underset{\hat{\mathbf{w}} \in \hat{\mathcal{W}}, \alpha \in \mathcal{A}, \beta \in \mathbb{R}^+}{\operatorname{argmin}} \quad \mathcal{L}_{\text{CP-NAS}}(\tilde{f}^P(\mathbf{w}, \alpha), \tilde{f}_{\mathbf{b}}^C(\hat{\mathbf{w}}, \hat{\alpha}, \beta)) \\
 &= \underset{\hat{\mathbf{w}} \in \hat{\mathcal{W}}, \alpha \in \mathcal{A}, \beta \in \mathbb{R}^+}{\operatorname{argmin}} \quad \tilde{f}^P(\mathbf{w}, \alpha) - \tilde{f}_{\mathbf{b}}^C(\hat{\mathbf{w}}, \hat{\alpha}, \beta),
 \end{aligned} \tag{10}$$

where $\tilde{f}^P(\mathbf{w}, \alpha)$ denotes the performance of the real-valued Parent model as pre-defined in Eq. 5. $\tilde{f}_{\mathbf{b}}^C$ is further defined as $\tilde{f}_{\mathbf{b}}^C(\hat{\mathbf{w}}, \alpha, \beta) = \sum_{n=1}^N \hat{p}_n(\hat{\mathbf{w}}, \alpha, \beta) \log(\hat{p}_n(\mathcal{X}))$ following Eq. 9. As shown in Eq. 10, we propose the \mathcal{L} for estimating the performance of the candidate architectures with binary weights and activations, which take both the real-valued architectures and 1-bit architectures into consideration.

3.4 Tangent propagation for DCP-NAS

In this section, we first propose the generation of the tangent direction based on the Parent model, and then present the tangent propagation for effectively searching the optimized architecture in the binary NAS. As shown in Fig. 3, the DCP-NAS novelty introduce tangent propagation and decoupled optimization, thus leading to a effective discrepancy-based search framework. The main motivation of DCP-NAS is to “fine-tune” the Child model architecture based on the real-valued Parent rather than directly binarizing the Parent. Thus we first take advantage of the Parent model to generate the tangent direction from the architecture gradient of model as

$$\frac{\partial \tilde{f}^C(\mathbf{w}, \alpha)}{\partial \alpha} = \sum_{n=1}^N \frac{\partial p_n(\mathbf{w}, \alpha)}{\partial \alpha}, \quad (11)$$

where $\tilde{f}(\mathbf{w}, \alpha)$ is pre-defined in Eq. 5.

Then we conduct the second step, *i.e.*, tangent propagation for Child model. For each epoch of binary NAS in our DCP-NAS, we inherit weights from the real-valued architecture $\hat{\alpha} \leftarrow \alpha$ and enforce the 1-bit CNNs to learn similar distributions as real-valued networks

$$\begin{aligned} \max_{\hat{\mathbf{w}} \in \mathcal{W}, \hat{\alpha} \in \mathcal{A}, \beta \in \mathbb{R}^+} G(\hat{\mathbf{w}}, \hat{\alpha}, \beta) &= \tilde{f}^P(\mathbf{w}, \alpha) \log \frac{\tilde{f}^P(\mathbf{w}, \alpha)}{\tilde{f}_{\mathbf{b}}^C(\hat{\mathbf{w}}, \hat{\alpha}, \beta)} \\ &= \sum_{n=1}^N p_n(\mathbf{w}, \alpha) \log \left(\frac{\hat{p}_n(\hat{\mathbf{w}}, \hat{\alpha}, \beta)}{p_n(\mathbf{w}, \alpha)} \right), \end{aligned} \quad (12)$$

where the KL divergence is employed to supervise the binary search process. $G(\hat{\mathbf{w}}, \hat{\alpha}, \beta)$ calculates the output logits similarity between real-valued network $p(\cdot)$ and 1-bit CNNs $\hat{p}(\cdot)$, where the output of teacher is already given.

To further optimize the binary architecture, we constrain the gradient of binary NAS using the tangent direction as

$$\min_{\hat{\alpha} \in \mathcal{A}} \mathbf{D}(\hat{\alpha}) = \left\| \frac{\partial \tilde{f}^P(\mathbf{w}, \alpha)}{\partial \alpha} - \frac{\partial G(\hat{\mathbf{w}}, \hat{\alpha}, \beta)}{\partial \hat{\alpha}} \right\|_2^2. \quad (13)$$

We use Eqs. 12 - 13 to jointly learn the DCP-NAS and rewrite the objective function in Eq. 10 as

$$\begin{aligned} \mathcal{L}_{\text{DCP-NAS}}(\tilde{f}^P(\mathbf{w}, \alpha), \tilde{f}_{\mathbf{b}}^C(\hat{\mathbf{w}}, \hat{\alpha}, \beta)) \\ = -G(\hat{\mathbf{w}}, \hat{\alpha}, \beta) + \lambda \mathbf{D}(\hat{\alpha}) + \mu \mathcal{L}_R(\hat{\mathbf{w}}, \beta). \end{aligned} \quad (14)$$

Then we optimize the binary architecture $\hat{\alpha}$ along the tangent direction of the real-valued model, which inherits from the real-valued one. Note that, when we set $\lambda = 0$, the Eq. 14 is equivalent to the objective of original CP-NAS [80]. As revealed in [45], the real-valued weights generally converge faster than the binary ones. Motivated by this observation, the tangent direction of Parent supernet can be used to approximate the optimization direction of the

more slowly converged Child supernet. To conclude, in Eq. 13, we improve the optimization of Child architecture based on tangent direction of the Parent architecture, which lead Child supernet to be more efficiently trained.

Considering the binary weights are learned by KL divergence, we optimize our DCP-NAS as

$$\begin{aligned} & \nabla_{\hat{\alpha}} \mathcal{L}_{\text{DCP-NAS}}(\tilde{f}^P(\mathbf{w}, \alpha), \tilde{f}_{\mathbf{b}}^C(\hat{\mathbf{w}}, \hat{\alpha}, \beta)) \\ &= -\frac{\partial G(\hat{\mathbf{w}}, \hat{\alpha}, \beta)}{\partial \hat{\alpha}} + \lambda \frac{\partial \mathbf{D}(\hat{\alpha})}{\partial \hat{\alpha}} \\ &= -\frac{\partial G(\hat{\mathbf{w}}, \hat{\alpha}, \beta)}{\partial \hat{\alpha}} + 2\lambda \left(\frac{\partial G(\hat{\mathbf{w}}, \hat{\alpha}, \beta)}{\partial \hat{\alpha}} - \frac{\partial \tilde{f}^P(\mathbf{w}, \alpha)}{\partial \alpha} \right) \frac{\partial^2 G(\hat{\mathbf{w}}, \hat{\alpha}, \beta)}{\partial \hat{\alpha}^2} \quad (15) \\ &= -\frac{\partial G(\hat{\mathbf{w}}, \hat{\alpha}, \beta)}{\partial \hat{\alpha}} + 2\lambda \left(\frac{\partial G(\hat{\mathbf{w}}, \hat{\alpha}, \beta)}{\partial \hat{\alpha}} - \frac{\partial \tilde{f}^P(\mathbf{w}, \alpha)}{\partial \alpha} \right) \mathbf{H}_G(\hat{\alpha}), \end{aligned}$$

$$\nabla_{\hat{\mathbf{w}}} \mathcal{L}_{\text{DCP-NAS}}(\tilde{f}^P(\mathbf{w}, \alpha), \tilde{f}_{\mathbf{b}}^C(\hat{\mathbf{w}}, \hat{\alpha}, \beta)) = -\frac{\partial G(\hat{\mathbf{w}}, \hat{\alpha}, \beta)}{\partial \mathbf{b}^{\hat{\mathbf{w}}}} \frac{\partial \mathbf{b}^{\hat{\mathbf{w}}}}{\partial \hat{\mathbf{w}}}, \quad (16)$$

where

$$\frac{\partial \mathbf{b}^{\hat{\mathbf{w}}}}{\partial \hat{\mathbf{w}}} = \mathbf{1}_{|\hat{\mathbf{w}}| \leq 1}. \quad (17)$$

and λ is a hyper-parameter, $\mathbf{H}_{\tilde{f}_{\mathbf{b}}^C}(\hat{\alpha}) = \frac{\partial^2 \tilde{f}_{\mathbf{b}}^C(\hat{\mathbf{w}}, \hat{\alpha})}{\partial \hat{\alpha}^2}$ denotes the Hessian matrix. The process of DCP-NAS is outlined in Fig. 3.

We minimize the difference of the gradient (tangent direction) $\frac{\partial G(\hat{\mathbf{w}}, \hat{\alpha}, \beta)}{\partial \hat{\alpha}}$ and the gradient (tangent direction) $\frac{\partial \tilde{f}^P(\mathbf{w}, \alpha)}{\partial \alpha}$, aiming to search the architectures (both real-valued NAS and binary NAS) in the same direction to generate a better 1-bit architecture. Note that α inherits from $\hat{\alpha}$ at the beginning of each real-valued NAS iteration, indicating that we only utilize \mathbf{w} and α of real-valued NAS for heuristic optimization direction for 1-bit NAS instead of searching a better architecture for real-valued networks. Since a better tangent direction $\frac{\partial \tilde{f}^P(\mathbf{w}, \alpha)}{\partial \alpha}$ is achieved, DCP-NAS can have a more suitable $\hat{\alpha}$ for 1-bit CNNs. We note that α is different from $\hat{\alpha}$, which is not an optimized architecture for real-valued weights but an optimized architecture for binary weights.

The expression above contains an expensive matrix gradient computation in its second term. Thus we introduce a first-order approximation of the Hessian matrix for accelerating the searching efficiency in Section 3.5.

3.5 Generalized Gauss-Newton matrix (GN) for Hessian Matrix

Since the Hessian matrix is computationally expensive, this section mainly tries to accelerate the aforementioned Hessian matrix calculation by deriving a second-order expansion based on Eq. 16.

Below we proved that the Hessian matrix of loss function is directly related to the expectation of the covariance of the gradient. Considering the loss

function as the negative logarithm of the likelihood, let \mathcal{X} be a set of input data of network and $p(\mathcal{X}; \hat{\mathbf{w}}, \hat{\alpha})$ be the predicted distribution over \mathcal{X} under the parameters of the network are $\hat{\mathbf{w}}$ and $\hat{\alpha}$, *i.e.*, output logits of the head layer.

With omitting $\hat{\mathbf{w}}$ for simplicity, the Fisher information of the probability distribution set $P = \{p_n(\mathcal{X}; \hat{\alpha}), n \in N\}$ can be described by a matrix whose value at the i -th row and j -th column

$$I_{i,j}(\hat{\alpha}) = \mathbb{E}_{\mathcal{X}} \left[\frac{\partial \log p_n(\mathcal{X}; \hat{\alpha})}{\partial \hat{\alpha}_i} \frac{\partial \log p_n(\mathcal{X}; \hat{\alpha})}{\partial \hat{\alpha}_j} \right]. \quad (18)$$

Remembering that N denotes the number of classes as described in Eq. 5. It is then trivial to prove that the Fisher information of the probability distribution set P approaches a scaled version of the Hessian of log-likelihood as

$$I_{i,j}(\hat{\alpha}) = -\mathbb{E}_{\mathcal{X}} \left[\frac{\partial^2 \log p_n(\mathcal{X}; \hat{\alpha})}{\partial \hat{\alpha}_i \partial \hat{\alpha}_j} \right]. \quad (19)$$

Let $H_{i,j}$ denotes the second-order partial derivatives $\frac{\partial^2}{\partial \hat{\alpha}_i \partial \hat{\alpha}_j}$. Noting that the first derivatives of log likelihood is

$$\frac{\partial \log p_n(\mathcal{X}; \hat{\alpha})}{\partial \hat{\alpha}_i} = \frac{\partial p_n(\mathcal{X}; \hat{\alpha})}{p_n(\mathcal{X}; \hat{\alpha}) \partial \hat{\alpha}_i}, \quad (20)$$

the second derivatives is

$$\mathbf{H}_{i,j} \log p_n(\mathcal{X}; \hat{\alpha}) = \frac{\mathbf{H}_{i,j} p_n(\mathcal{X}; \hat{\alpha})}{p_n(\mathcal{X}; \hat{\alpha})} - \frac{\partial p_n(\mathcal{X}; \hat{\alpha})}{p_n(\mathcal{X}; \hat{\alpha}) \partial \hat{\alpha}_i} \frac{\partial p_n(\mathcal{X}; \hat{\alpha})}{p_n(\mathcal{X}; \hat{\alpha}) \partial \hat{\alpha}_j}. \quad (21)$$

Considering that

$$\begin{aligned} \mathbb{E}_{\mathcal{X}} \left(\frac{\mathbf{H}_{i,j} p_n(\mathcal{X}; \hat{\alpha})}{p_n(\mathcal{X}; \hat{\alpha})} \right) &= \int \frac{\mathbf{H}_{i,j} p_n(\mathcal{X}; \hat{\alpha})}{p_n(\mathcal{X}; \hat{\alpha})} p_n(\mathcal{X}; \hat{\alpha}) d\mathcal{X} \\ &= \mathbf{H}_{i,j} \int p_n(\mathcal{X}; \hat{\alpha}) d\mathcal{X} = 0, \end{aligned} \quad (22)$$

we take the expectation of the second derivative and then obtain

$$\begin{aligned} \mathbb{E}_{\mathcal{X}} (\mathbf{H}_{i,j} \log p_n(\mathcal{X}; \hat{\alpha})) &= -\mathbb{E}_{\mathcal{X}} \left\{ \frac{\partial p_n(\mathcal{X}; \hat{\alpha})}{p_n(\mathcal{X}; \hat{\alpha}) \partial \hat{\alpha}_i} \frac{\partial p_n(\mathcal{X}; \hat{\alpha})}{p_n(\mathcal{X}; \hat{\alpha}) \partial \hat{\alpha}_j} \right\} \\ &= -\mathbb{E}_{\mathcal{X}} \left\{ \frac{\partial p_n(\mathcal{X}; \hat{\alpha})}{\partial \hat{\alpha}_i} \frac{\partial p_n(\mathcal{X}; \hat{\alpha})}{\partial \hat{\alpha}_j} \right\}. \end{aligned} \quad (23)$$

Thus an equivalent substitution for the Hessian matrix $\mathbf{H}_{\hat{\mathbf{f}}_b}(\hat{\alpha})$ in Eq. 14 is the product of two first-order derivatives. This concludes the proof that we can use the covariance of gradients to represent the Hessian matrix for efficient computation.

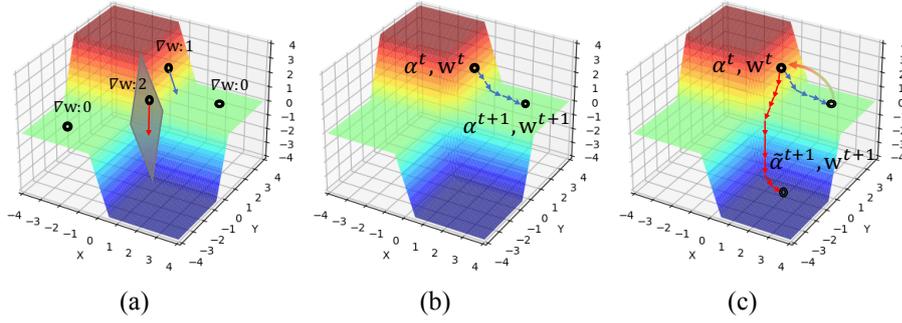


Fig. 6 The loss landscape illustration of supernet. (a) The gradient of current weights with different α , (b) The vanilla α^{t+1} with backpropagation, (c) $\tilde{\alpha}^{t+1}$ with the decoupled optimization.

3.6 Decoupled optimization for training the DCP-NAS

As one of the critical steps to achieve a powerful NAS strategy, the optimization of α plays an essential role. However, existing NAS and BinaryNAS methods neglect the intrinsic coupling relationship of the weights and architecture parameters, resulting in a sub-optimal architecture caused by an insufficient searching process. In this section, we first describe the coupling relationship between the weights and architecture parameters existing in the DCP-NAS. Then we present the decoupled optimization during backpropagation of the sampled supernet for fully and effectively optimizing such two coupling parameters.

3.6.1 Coupled models for DCP-NAS

Combing Eq. 1 and Eq. 2, we first show how parameters in DCP-NAS are formulated in a coupling relationship as $\mathbf{a}^{(j)} = \sum_{i < j} \text{softmax}(\alpha_m^{(i,j)})(\mathbf{w}^{(i,j)} \otimes \mathbf{a}^{(i)})$, where $\mathbf{w}^{(i,j)} = [[\mathbf{w}_m]] \in \mathbb{R}_{M \times 1}$, $\mathbf{w}_m \in \mathbb{R}_{C_{out} \times C_{in} \times K_m \times K_m}$ denote the weights of all candidate operations between i -th and j -th nodes and K_m denotes the kernel size of m -th operation. Specifically, for pooling and identity operations, K_m equals to the downsample size and the feature map size, \mathbf{w}_m equals to $\mathbf{1}/(K_m \times K_m)$ and $\mathbf{1}$, respectively. For each intermediate node, its output $\mathbf{a}^{(j)}$ are jointly determined by $\alpha_m^{(i,j)}$ and $\mathbf{w}_m^{(i,j)}$, while $\mathbf{a}^{(i)}$ is independent with both $\alpha_m^{(i,j)}$ and $\mathbf{w}_m^{(i,j)}$. As shown in Fig 6 (a) and (b), with different α , the gradient of corresponding \mathbf{w} can be various and sometimes hard to be optimized, which are possibly trapped into a local minima. However, with decouple the α and \mathbf{w} , the supernet has the chance to jump out of a local minima and be optimized with a better convergence.

Based on the deviation and analysis above, we propose our objective for optimizing the neural architecture search process

$$\arg \min_{\alpha, \mathbf{w}} \mathcal{L}(\mathbf{w}, \alpha) = \begin{cases} \mathcal{L}_{\text{NAS}} + \text{reg}(\mathbf{w}), & \text{for Parent model} \\ \mathcal{L}_{\text{DCP-NAS}} + \text{reg}(\mathbf{w}), & \text{for Child model} \end{cases} \quad (24)$$

where $\alpha \in \mathbb{R}_{E \times M}$, $\mathbf{w} \in \mathbb{R}_{M \times 1}$ and $\text{reg}(\cdot)$ denotes the regularization item. Following [43, 67], the weights \mathbf{w} and the architectural parameters α are optimized sequentially, in which \mathbf{w} and α are updated independently. However, it is improper to optimize \mathbf{w} and α independently due to their coupling relationship. We consider the searching and training process of differentiable Child-Parent neural architecture search as a coupling optimization problem and solve the problem using a new backtracking method. The details will be shown in Section 3.6.2.

3.6.2 Decoupled optimization for Child-Parent model

We reconsider the coupling relation between \mathbf{w} and α from a new perspective. The derivative calculation process of \mathbf{w} should consider its coupling parameters α . Based on the chain rule [50] and its notations, we have

$$\begin{aligned} \tilde{\alpha}^{t+1} &= \alpha^t + \eta_1 \left(-\frac{\partial \mathcal{L}(\alpha^t, \mathbf{w}^t)}{\partial \alpha^t} + \eta_2 \text{Tr} \left[\left(\frac{\partial \mathcal{L}(\alpha^t, \mathbf{w}^t)}{\partial \mathbf{w}^t} \right)^T \frac{\partial \mathbf{w}^t}{\partial \alpha^t} \right] \right) \\ &= \alpha^{t+1} + \eta_1 \eta_2 \text{Tr} \left[\left(\frac{\partial \mathcal{L}(\alpha^t, \mathbf{w}^t)}{\partial \mathbf{w}^t} \right)^T \frac{\partial \mathbf{w}^t}{\partial \alpha^t} \right], \end{aligned} \quad (25)$$

where η_1 represents the learning rate, η_2 represents the coefficient of backtracking, $\tilde{\alpha}^{t+1}$ denotes the value after backtracking from the vanilla α^{t+1} , while vanilla α^{t+1} is calculated from the backpropagation rule and corresponding optimizer in the neural network. $\text{Tr}(\cdot)$ represents the trace of a matrix. However, the item $\frac{\partial \mathbf{w}^t}{\partial \alpha^t}$ of Eq. 25 is undefined and unsolvable based on the normal backpropagation process. To address this problem, we propose a decoupled optimization method as following. In the following, we omit the superscript \cdot^t and define $\tilde{\mathcal{L}}$ as

$$\tilde{\mathcal{L}} = \left(\frac{\partial \mathcal{L}(\alpha, \mathbf{w})}{\partial \mathbf{w}} \right)^T / \alpha, \quad (26)$$

which considers the coupling optimization problem as in Eq. 24. Note that $R(\cdot)$ is only considered when backtracking. Thus we have

$$\frac{\partial \mathcal{L}(\alpha, \mathbf{w})}{\partial \mathbf{w}} = \text{Tr} \left[\alpha \tilde{\mathcal{L}} \frac{\partial \mathbf{w}}{\partial \alpha} \right]. \quad (27)$$

For simplifying the derivation, we rewrite $\tilde{\mathcal{L}}$ as $[\tilde{g}_1, \tilde{g}_e, \dots, \tilde{g}_E]$, where each \tilde{g}_e is a column vector. Assuming that \mathbf{w}_m and $\alpha_{i,j}$ are independent when $m \neq j$,

$\alpha_{i,j}$ denotes an specific element in matrix α , we have

$$\left(\frac{\partial \mathbf{w}}{\partial \alpha}\right)_m = \begin{bmatrix} 0 & \dots & \frac{\partial \mathbf{w}_m}{\partial \alpha_{1,m}} & \dots & 0 \\ \cdot & & \cdot & & \cdot \\ 0 & \dots & \frac{\partial \mathbf{w}_m}{\partial \alpha_{e,m}} & \dots & 0 \\ \cdot & & \cdot & & \cdot \\ 0 & \dots & \frac{\partial \mathbf{w}_m}{\partial \alpha_{E,m}} & \dots & 0 \end{bmatrix}_{E \times M} \quad (28)$$

and with rewritten α as a column vector $[\alpha_1, \alpha_e, \dots, \alpha_E]^T$ with each α_e is a row vector, we have

$$\alpha \tilde{\mathcal{L}} = \begin{bmatrix} \alpha_1 \tilde{g}_1 & \dots & \alpha_1 \tilde{g}_e & \dots & \alpha_1 \tilde{g}_E \\ \cdot & & \cdot & & \cdot \\ \alpha_e \tilde{g}_1 & \dots & \alpha_e \tilde{g}_e & \dots & \alpha_e \tilde{g}_E \\ \cdot & & \cdot & & \cdot \\ \alpha_E \tilde{g}_1 & \dots & \alpha_E \tilde{g}_e & \dots & \alpha_E \tilde{g}_E \end{bmatrix}_{E \times E} \quad (29)$$

Combing Eq. 28 and Eq. 29, the matrix in the trace item of Eq. 26 can be written as

$$\alpha \tilde{\mathcal{L}} \left(\frac{\partial \mathbf{w}}{\partial \alpha}\right)_m = \begin{bmatrix} 0 & \dots & \alpha_1 \sum_{e'=1}^E \tilde{g}_{e'} \frac{\partial \mathbf{w}_m}{\partial \alpha_{e',m}} & \dots & 0 \\ \cdot & & \cdot & & \cdot \\ 0 & \dots & \alpha_e \sum_{e'=1}^E \tilde{g}_{e'} \frac{\partial \mathbf{w}_m}{\partial \alpha_{e',m}} & \dots & 0 \\ \cdot & & \cdot & & \cdot \\ 0 & \dots & \alpha_E \sum_{e'=1}^E \tilde{g}_{e'} \frac{\partial \mathbf{w}_m}{\partial \alpha_{e',m}} & \dots & 0 \end{bmatrix}_{E \times M} \quad (30)$$

Thus the whole matrix $\alpha \tilde{\mathcal{L}} \left(\frac{\partial \mathbf{w}}{\partial \alpha}\right)$ is with the size of $E \times M \times M$. After the above derivation, we compute the e -th component of the trace item in Eq. 26 as

$$Tr[\alpha \tilde{\mathcal{L}} \left(\frac{\partial \mathbf{w}}{\partial \alpha}\right)]_e = \alpha_e \sum_{m=1}^M \sum_{e'=1}^E \tilde{g}_{e'} \frac{\mathbf{w}_m}{\partial \alpha_{e',m}} \quad (31)$$

Noting that in the vanilla propagation process, $\alpha^{t+1} = \alpha^t - \eta_1 \frac{\partial \mathcal{L}(\alpha^t)}{\partial \alpha^t}$, thus combining Eq. 31 we have

$$\begin{aligned} \tilde{\alpha}^{t+1} &= \alpha^{t+1} - \eta \begin{bmatrix} \sum_{m=1}^M \sum_{e'=1}^E \tilde{g}_{e'} \frac{\partial \mathbf{w}_m}{\partial \alpha_{e',m}} \\ \sum_{m=1}^M \sum_{e'=1}^E \tilde{g}_{e'} \frac{\partial \mathbf{w}_m}{\partial \alpha_{e',m}} \\ \sum_{m=1}^M \sum_{e'=1}^E \tilde{g}_{e'} \frac{\partial \mathbf{w}_m}{\partial \alpha_{e',m}} \end{bmatrix} \otimes \begin{bmatrix} \alpha_1 \\ \alpha_e \\ \alpha_E \end{bmatrix} \\ &= \alpha^{t+1} + \eta \psi^t \otimes \alpha^t, \end{aligned} \quad (32)$$

where \otimes represents the Hadamard product and $\eta = \eta_1 \eta_2$. We take $\psi^t = -[\sum_{m=1}^M \sum_{e'=1}^E \tilde{g}_{e'} \frac{\partial \mathbf{w}_m}{\partial \alpha_{e',m}}, \dots, \sum_{m=1}^M \sum_{e'=1}^E \tilde{g}_{e'} \frac{\partial \mathbf{w}_m}{\partial \alpha_{e',m}}]^T$. Note that, $\frac{\partial \mathbf{w}}{\partial \alpha}$ is unsolvable and has no explicit form in NAS, which causes a unsolvable ψ^t . Thus

we introduce a learnable parameter $\tilde{\psi}^t$ for approximating ψ^t , which back-propagation process is calculated as

$$\tilde{\psi}^{t+1} = |\tilde{\psi}^t - \eta_{\psi} \frac{\partial \mathcal{L}}{\partial \tilde{\psi}^t}|. \quad (33)$$

Eq. 32 shows our method is actually based on a projection function to solve the coupling problem of the optimization by the learnable parameter $\tilde{\psi}^t$. In this method, we consider the influence of α^t and backtrack the optimized state at the $(t+1)$ -th step to form $\tilde{\alpha}^{t+1}$. However, where and when the backtracking should be applied is the key point in the optimization, thus we define the update rule as

$$\tilde{\alpha}_m^{t+1} = \begin{cases} P(\alpha_{:,m}^{t+1}, \alpha_{:,m}^t), & \text{if } \text{ranking}(R(\mathbf{w}_m)) > \tau \\ \alpha_{:,m}^{t+1}, & \text{otherwise} \end{cases} \quad (34)$$

where $P(\alpha_{:,m}^{t+1}, \alpha_{:,m}^t) = \alpha_{:,m}^{t+1} + \eta \tilde{\psi}^t \otimes \alpha_{:,m}^t$ and subscript \cdot_m denotes a specific edge. $R(\mathbf{w}_m)$ denotes the norm constraint of \mathbf{w}_m and is further defined as

$$R(\mathbf{w}_m) = \|\mathbf{w}_m\|_2^2, \quad \forall m = 1, \dots, M, \quad (35)$$

where τ denotes the threshold of deciding or not to backtrack. We further define the threshold as

$$\tau = \lfloor \epsilon \cdot M \rfloor \quad (36)$$

where ϵ denotes a hyper-parameter for controlling the percentage of edges backtracking. With backtracking α , the supernet can learn to jump out of the local minima. The overall process of DCP-NAS is outlined in Algorithm 1. Note that the decoupled optimization can be employed to both Parent and Child model. When applied to Child model, the \mathbf{w} here denotes the reconstructed weights from the binary weights, *i.e.*, $\mathbf{w} = \beta \circ \mathbf{b}^{\tilde{\mathbf{w}}}$.

4 Experiments

We quantitatively demonstrate the effectiveness of our DCP-NAS in this section. We first provide detailed ablation studies of the proposed DCP-NAS in Section 4.2. Then we evaluate the proposed DCP-NAS on the architecture search for image classification on the widely-used CIFAR-10 [33] and ImageNet ILSVRC12 [15] datasets with different search spaces and constraints in Section 4.3. Finally, to further validate the effectiveness and generalizability, we transfer the architectures searched on ImageNet to person re-identification and object detection task from Sections 4.4 to 4.5. In the following experiments, both the kernels and the activations are binarized. The leading performances reported in the following sections verify the superiority of our DCP-NAS.

Algorithm 1: Search process of DCP-NAS

Input: Training data, validation data
Parameter: Searching hyper-graph: $\mathcal{G}, M = 8, e(o_m^{(i,j)}) = 0$ for all edges
Output: Optimized $\hat{\alpha}^*$.

- 1: **while** DCP-NAS **do**
- 2: **while** Training real-valued Parent **do**
- 3: Search a temporary real-valued architecture $p(\mathbf{w}, \alpha)$.
- 4: Decoupled optimization from Eqs. 25 to 35.
- 5: Generate the tangent direction $\frac{\partial \tilde{f}(\mathbf{w}, \alpha)}{\partial \alpha}$ from Eqs. 5 to 11.
- 6: **end while**
- 7: Architecture inheriting $\hat{\alpha} \leftarrow \alpha$.
- 8: **while** Training 1-bit Child **do**
- 9: Calculate the learning objective from Eqs. 10 to 14.
- 10: Tangent propagation from Eqs. 15 to 23 and decoupled optimization from Eqs. 25 to 35.
- 11: Obtain the $\hat{p}(\hat{\mathbf{w}}, \hat{\alpha})$.
- 12: **end while**
- 13: Architecture inheriting $\alpha \leftarrow \hat{\alpha}$.
- 14: **end while**
- 15: **return** Optimized architecture $\hat{\alpha}^*$.

4.1 Datasets and Implementation Details

4.1.1 Datasets

CIFAR-10 [33] is a natural image classification dataset, which is composed of a training set and a test set, with 50,000 and 10,000 32×32 color images, respectively. These images span 10 different classes, including airplanes, automobiles, birds, cats, deer, dogs, frogs, horses, ships, and trucks. ImageNet ILSVRC12 object classification dataset [15] is more diverse and challenging. It contains 1.2 million training images, and 50,000 validation images, across 1000 classes.

The PASCAL VOC dataset contains natural images from 20 different classes. We train our model on the VOC `trainval2007` and VOC `trainval2012` sets, which consist of approximately 16k images. We then evaluate our method on the VOC `test2007` set, which includes 4952 images. Following [17], we use the mean average precision (mAP) as the evaluation criterion.

The COCO dataset consists of images from 80 categories. We conduct experiments on the COCO 2014 [40] object detection track. Models are trained with the combination of 80k images from the COCO `train2014` and 35k images sampled from COCO `val2014`, *i.e.*, COCO `trainval35k`. On the remaining 5k images from the COCO `minival`, we test our method based on the average precision (AP) for $\text{IoU} \in [0.5 : 0.05 : 0.95]$ denoted as $\text{mAP}@[.5, .95]$. We also report AP_{50} , AP_{75} , AP_s , AP_m , and AP_l to further analyze our method.

The Market-1501 [74], DukeMTMC-reID [77] and CUHK03 [36] are among the largest Re-ID benchmark datasets based on images. In Market-

1501, 32,668 labeled bounding boxes of 1,501 identities captured from six different perspectives are contained, which are detected using Deformable Part Model (DPM) [18]. The Market-1501 dataset is split into two parts: a training dataset consisting of 12,936 images with 751 identities and a test dataset containing 19,732 images with 750 identities. For testing, 3,368 hand-painted images with 750 identities are used as probe set to classify the identities of test dataset. DukeMTMC-reID is a subset of the DukeMTMC for Re-ID based on images, with the same format as the Market-1501 dataset. The original dataset is sampled from 85-minute high-resolution videos recorded by eight different cameras. Hand-painted pedestrian bounding boxes are also available for classification during test. CUHK03 offers both hand-labeled and DPM-detected bounding boxes, and we use the latter in this article. CUHK03 originally adopts 20 random train/test splits, which is time-consuming for deep learning.

4.1.2 Search protocol

In the search process, we consider a total of 6 cells with the initial 16 channels in the network, where the reduction cell are inserted in the second and the fourth layers, and the others are normal cells. There are 4 intermediate nodes in each cell. The initial number of operations M is set as 8 and the number of search epochs is set as 105 (same as CP-NAS [80]) for fairness. λ is set as $1e - 4$, and the batch size is set to 512. We use Adam with momentum to optimize the network weights, with an initial learning rate of 0.025 (annealed down to zero following a cosine schedule), a momentum of 0.9, and a weight decay of $5e - 4$. When we search for the architecture directly on ImageNet, we use the same parameters for searching with CIFAR-10 except that the initial learning rate is set to 0.05 and λ is set to $1e - 3$. Due to the efficient guidance of DCP-NAS model, we use 50% of the training set with CIFAR-10 and ImageNet for architecture search and 5% of the training set for evaluation, leading to a faster search. After search, in the architecture evaluation step, our experimental settings are similar to [43, 82, 51].

4.1.3 Training protocol

A larger network of 10 cells (8 normal cells and 2 reduction cells) is trained on CIFAR-10 for 600 epochs with a batch size of 96 and an additional regularization cutout [16]. The initial number of channels is set as 54, 72, 108 for different model sizes, *i.e.*, DCP-NAS-S/M/L in this paper. We use the Adam optimizer with an initial learning rate of 0.0025 (annealed down to zero following a cosine schedule without restart), a momentum of 0.9, a weight decay of $1e - 5$, and a gradient clipping at 5.

When stacking the cells to evaluate on ImageNet, the evaluation stage follows that of DARTS [43], which starts with three convolutional layers with a stride of 2 to reduce the input image resolution from 224×224 to 28×28 . 10 cells (8 normal cells and 2 reduction cells) are stacked after these three layers,

Table 3 Effect of with/without the reconstruction error and tangent direction constraint on the ImageNet dataset. The architecture conducted the experiments is DCP-NAS-L.

Tangent direction ($\mathbf{D}(\hat{\alpha})$)		✗	✓	✗	✓
Reconstruction error ($\mathcal{L}_R(\hat{\mathbf{w}}, \beta)$)		✗	✗	✓	✓
Accuracy	Top-1	66.7	68.3	68.2	72.4
	Top-5	83.3	85.0	85.1	89.2

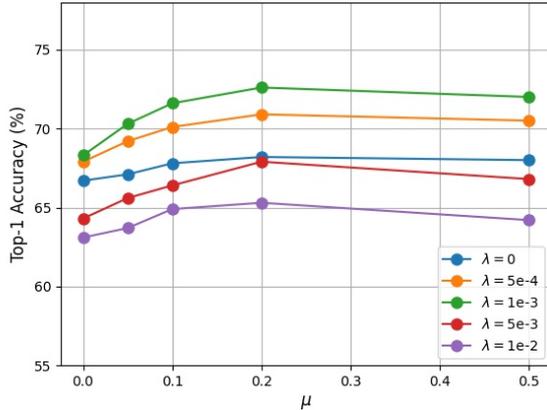


Fig. 7 With different λ and μ , we evaluate the Top-1 accuracies of DCP-NAS-L on ImageNet.

with the initial channel number being 54, 72, 108 for different model sizes, *i.e.*, DCP-NAS-S/M/L. The network is trained from scratch for 512 epochs using a batch size of 512. We use Adam optimizer with a momentum for 0.9, an initial learning rate of 0.001 (decayed down to zero following a cosine schedule), and a weight decay of $3e - 5$. Further enhancements are adopted including label smoothing and an auxiliary loss tower during training. All the experiments and models are implemented in PyTorch [49].

On person re-identification task, we utilize SGD as the optimizer with momentum set as 0.9 and weight decay set as $1e - 4$. On all three datasets, initial learning rates are both set to 0.05 with a Cosine Annealing learning rate decay. The total epochs are 120 and the batch size is set as 64. For object detection task, we implement experiments following [66] on two-stage Faster-RCNN with FPN neck, the lateral connection of the FPN [39] neck is replaced with 3×3 1-bit convolution for improving performance. This adjustment is implemented in all of the Faster-RCNN experiments. For Faster-RCNN, we train the model for 12 epochs at a learning rate of 0.02, which decays by multiplying 0.1 in the 9-th and 11-th epochs.

4.2 Ablation Study

4.2.1 Effectiveness of Tangent Propagation

In this section, we evaluate the effects of the tangent propagation on the performance of DCP-NAS, the hyper-parameter used in this section includes λ , μ . Furthermore, we also discuss the effectiveness of the reconstruction error. The implementation details are given below.

For searching the binary neural architecture in a better form, λ and μ are used to balance the KL divergence $\tilde{f}(\hat{\mathbf{w}}, \hat{\alpha}, \beta)$ for supervising the Child, the reconstruction error for binary weights $\mathcal{L}_R(\hat{\mathbf{w}}, \beta)$ and the tangent direction constraint $\mathbf{D}(\hat{\alpha})$. We evaluate λ and μ on ImageNet dataset with the architecture DCP-NAS-L. To better understand the tangent propagation on the large-scale ImageNet ILSVRC12 dataset, we experiment to examine how the tangent direction constraint affects performance. According to the above experiments, we first set λ to $5e-3$ and μ to 0.2, if they are used. As shown in Tab. 3, both the tangent direction constraint and the reconstruction error can independently improve the accuracy on ImageNet. When applied together, the Top-1 accuracy reaches the highest value of 72.4%. Then we conduct experiments with various value of λ and μ as shown in Tab. 7. We observe that with a fixed value of μ , the Top-1 accuracy increases at the beginning with λ increasing, but decreases when λ is greater than $1e-3$. For when λ becomes larger, DCP-NAS tends to select the binary architecture with similar gradient to its real-valued counterpart. To some extent, the accuracy of 1-bit model is neglected, which leads to a performance drop. Another performance variation phenomenon is that the Top-1 accuracy rise first and then fall with the increase of μ while λ containing fixed values. For too much attention paid to minimizing the distance between 1-bit parameters and their counterparts may introduces a representation ability collapse to the 1-bit models and severely degenerates the performance of DCP-NAS.

To better understand the acceleration rate of applying Generalized Gauss-Newton matrix (GGN) in the searching process, we conduct experiments to examine the searching cost with and without GGN. As shown in Tab. 4, we compare the searching efficiency and the accuracy of the architecture obtained by Random Search (random selection), Real-valued NAS methods, Binarized NAS methods, CP-NAS, DCP-NAS without GGN method, DCP-NAS with GGN applied. In random search, the 1-bit supernet randomly samples and trains an architecture in each epoch, then assign the expectation of all performances to each corresponding edges and operations and returns the architecture with the highest score, which lacks the necessary guidance in the searching process, thus has a poor performance for binary architecture search. Notably, our DCP-NAS without GGN are highly computation consumed for the second order gradient is necessarily computed in the tangent propagation. Note that directly optimizing two supernet is computationally redundant. However, the introduction of GGN for Hessian matrix significantly accelerates the searching process, which reduces the search cost to nearly 10% with

Table 4 We compare the searching efficiency of different search strategies on ImageNet, including previous NAS on both real-valued and 1-bit search space, random search and our DCP-NAS. T.P. and D.O. denote Tangent Propagation and Decoupled Optimization, respectively.

Method		T.P.	GGN	D.O.	Top-1 Acc (%)	Search Cost (GPU days)
Real-valued NAS	PNAS	-	-	-	74.2	225
	DARTS	-	-	-	73.1	4
	PC-DARTS	-	-	-	75.8	3.8
Direct BNAS	BNAS ₁	-	-	-	64.3	2.6
	BNAS ₂ -H	-	-	-	63.5	-
	Random Search	-	-	-	51.3	4.4
Auxiliary BNAS	CP-NAS (Baseline)	-	-	-	66.5	2.8
	DCP-NAS-L	✓	✗	✗	71.4	27.9
	DCP-NAS-L	✓	✓	✗	71.2	2.9
	DCP-NAS-L	✓	✗	✓	72.6	27.9
	DCP-NAS-L	✓	✓	✓	72.4	2.9

Table 5 Comparison results on ImageNet dataset with DCP-NAS of distance calculation method used in constraining the gradient of binary NAS under the tangent direction, *i.e.*, Eq. 13. We use the small model size, *i.e.*, DCP-NAS-S, to evaluate the searched architecture.

Method	Accuracy(%)		Memory (Mbits)	Search Cost (GPU days)
	Top-1	Top-5		
Cosine similarity	62.5	83.9	4.2	2.9
L1-norm	62.7	84.3	4.3	2.9
F-norm	63.0	84.5	4.2	2.9

Table 6 Effect of the decoupled optimization on real-valued NAS using ImageNet dataset. D.O. denotes Decoupled Optimization for simplicity.

Architecture	Accuracy(%)		Memory (Mbits)	Search Cost (GPU days)
	Top-1	Top-5		
DARTS	73.1	91.0	156.8	4
DARTS + D.O.	73.8	92.1	157.2	4
PC-DARTS	75.8	92.7	169.6	3.8
PC-DARTS + D.O.	76.6	93.9	170.3	3.8

negligible accuracy vibration. As shown in Table 4, with GGN utilized, our method reduces the search cost from 29 to 2.9, which is more efficient than DARTS. Also, our DCP-NAS achieves much lower performance gap between the real-valued NAS with less search cost in a clear margin. To further clarify the tangent propagation, we conduct ablative experiments for different architecture discrepancy calculation methods. As shown in Tab. 5, F-norm applied in Eq. 13 achieves best performance, while cosine similarity and L1-norm are not so effective as the F-norm.

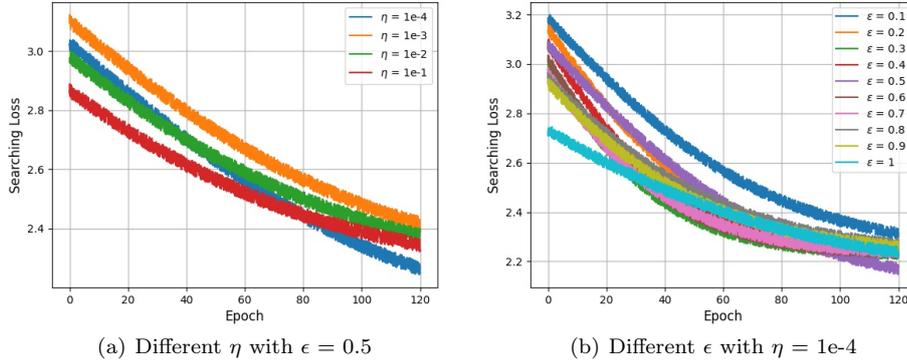


Fig. 8 Effect of hyper-parameters ϵ and η on one-stage and two-stage training using 1-bit ResNet-18.

4.2.2 Effectiveness of decoupled optimization

In this section, we evaluate the effects of the decoupled optimization for solving the coupling relation between \mathbf{w} and α in our DCP-NAS. The implementation details are given below. For architecture parameters α , we use the Adam optimizer with an initial learning rate $\eta = 1e - 4$. We set the hyper-parameter $\epsilon = 0.2$. The learning rate η_ψ is set as $1e - 4$. As shown in Fig. 8 (a), when η is set as $1e-4$, the searching loss converges to a better minimum. In Fig. 8 (b), when ϵ equals to 1, there will be no edges or operations been backtracked through our decoupled optimization, thus the loss in the very beginning of the searching process will be much lower than other conditions. However, the supernet converges to a local minima with higher searching loss at the end of the searching process. Furthermore, when ϵ is set as 0.1, which denotes that over 90 % of the edges and operations will be forced to backtrack during optimization. Too much backtracking degenerates the learning of the binary weights and their corresponding scale factors, thus leading to under fitting of the sampled supernet which affects the performance estimation for a better architecture. We also extend our decoupled optimization into real-valued NAS frameworks shown in Tab. 6. For experiments on DARTS [43] and PC-DARTS [67], our decoupled optimization method further boosts the performance by 0.7% and 0.8% with no extra search cost, which proves the effectiveness of our decoupled optimization for real-valued NAS.

4.2.3 Effectiveness of Search Space

In this section, we evaluate the effectiveness of the pooling layer and the skip connection in the search space. We first remove the pooling layer and skip connection from the search space, respectively. Then we simultaneously remove them for further evaluation. As shown in the 5-th \sim 7-th rows of Tab. 7, the

Table 7 Performance comparison of architectures obtained under different search spaces. The “w/o” denotes without in short.

Search Space	Architecture	Accuracy		Memory (Mbits)
		Top-1	Top-5	
Normal search space	Small	63.0	84.5	4.2
	Medium	69.0	88.4	6.9
	Large	72.4	90.6	6.9
w/o pooling layer	Small	59.1	80.2	5.1
	Medium	64.3	85.3	7.2
	Large	68.9	87.3	7.2
w/o skip connection	Small	60.0	81.3	5.0
	Medium	65.1	85.9	7.1
	Large	69.5	88.6	7.1
w/o pooling layer + w/o skip connection	Small	58.9	80.1	5.3
	Medium	64.0	85.1	7.4
	Large	68.2	88.9	7.4

Table 8 Ablative experiments on ImageNet datasets for binarization proportion in CP-NAS and DCP-NAS. [#] denotes that we “binarize activations of preprocessing operations for 2 input nodes”. [§] denotes that we “binarize activations of the first convolutional layer in the depth-wise separable convolution”.

Architecture	Memory (Mbits)	OPs (M)	Accuracy (%)	
			Top-1	Top-5
DCP-NAS-S [#]	3.5	91	60.7	82.9
DCP-NAS-S [§]	3.9	98	61.2	83.2
DCP-NAS-S	4.2	100	63.0	84.5
DCP-NAS-M [#]	5.8	152	66.9	86.1
DCP-NAS-M [§]	6.2	159	67.2	86.7
DCP-NAS-M	6.9	163	69.0	88.4
CP-NAS-L [#]	11.8	307	63.1	83.0
CP-NAS-L [§]	11.8	316	64.3	84.2
CP-NAS-L	12.5	323	66.5	86.8
DCP-NAS-L [#]	13.5	319	69.5	88.7
DCP-NAS-L [§]	13.8	326	70.1	89.1
DCP-NAS-L	14.4	334	72.4	90.6

performance significantly drops when removing the pooling operations from the search space, primarily due to the insufficient feature fusing. Removing all skip-connection operation (identity mapping) from the search space also degrades the network performance, as well as increasing the memory usage to over $1.3\times$ than the network obtained from the normal search space (described in Sec 3.1). As shown in the bottom of Tab. 7, with removing all pooling and skip-connection operations, the performance of small models drop to only 58.9% Top-1 accuracy. These experiments show the rationality of our search space setup.

Table 9 Ablative experiments on the CIFAR-10 and ImageNet datasets for training scheme based on DCP-NAS-S.

Training scheme on CIFAR-10				Error (%)	
Regularization	Cutout	Gradient Clipping			
✓		✓	6.13		
✓		✗	6.42		
✗		✓	6.49		
✗		✗	6.87		

Training scheme on ImageNet				Accuracy (%)	
Label Smoothing	Auxiliary Loss		Top-1	Top-5	
✓			✓	63.0	84.5
✓			✗	62.9	84.3
✗			✓	62.6	84.2
✗			✗	62.1	83.7

4.2.4 Effectiveness of Binarization Proportion

In this section, we evaluate the effects of the binarization of “activations of the preprocessing operations of the two input nodes” and “activations of the first convolutional layer in the depth-wise separable convolution”. As shown in Tab. 8, with “activations of the preprocessing operations of the two input nodes” binarized, the Top-1 accuracy decreases up to 2.9% compare with the final DCP-NAS models, which is less acceptable compared with the marginal memory saving and OPs decrease.

Furthermore, compared with “activations of the first convolutional layer in the depth-wise separable convolution”, the final DCP-NAS achieves desirable performance improvements with acceptable memory saving and OPs. For example, the DCP-NAS-L surpasses DCP-NAS-L^s by 2.3% Top-1 accuracy, which is significant on ImageNet classification task. Similar phenomenon is also observed on the CP-NAS framework. Therefore, our DCP-NAS achieves significant better performance with acceptable extra memory cost, when saving “activations of the preprocessing operations of the two input nodes” and “activations of the first convolutional layer in the depth-wise separable convolution” as real-valued.

4.2.5 Effectiveness of Training Scheme

In this section, we evaluate the effects of the training scheme from original CP-NAS [80], *i.e.*, additional regularization cutout and gradient clipping used in training DCP-NAS on CIFAR-10, and label smoothing together with auxiliary loss used in training DCP-NAS on ImageNet. As shown in Tab. 9, applying additional regularization cutout and gradient clipping into training on CIFAR-10 brings the best performance, while removing both brings 0.74% performance drop. Likewise, on ImageNet, applying label smoothing and auxiliary loss into training brings the best performance, which surpasses the counterpart with

Table 10 Test accuracies on the CIFAR-10 datasets. Quantization method in DCP-NAS are based on MCN [58]. We calculate the number of parameters for each model, and the numbers refer to the models on CIFAR-10. Err. and Mem. denote Test Error and Memory Usage, respectively. Gradient-based[‡] denotes our DCP-NAS is a discrepancy-aware gradient-based search method.

Architecture	Err. (%)	Mem. (Mbits)	W/A	Search Cost (GPU days)	Search Method
WRN-22	5.04	138.6	32/32	-	Manual
DARTS	2.83	108.8	32/32	4	Gradient-based
PC-DARTS	2.78	112.0	32/32	0.15	Gradient-based
WRN-22	5.69	4.3	1/32	-	Manual
BNAS ₁	3.94	2.6	1/32	0.09	Performance-based
CP-NAS-S	6.50	2.9	1/1	0.1	Child-Parent model
BATS (Small)	6.30	2.8	1/1	0.25	Gradient-based
DCP-NAS-S	6.13	2.9	1/1	0.1	Gradient-based [‡]
WRN-22 (BONN)	8.07	4.3	1/1	-	Manual
BNAS ₁ (Large)	8.29	4.5	1/1	0.09	Performance-based
BNAS ₂ v2-B	6.24	4.5	1/1	0.09	Gradient-based
CP-NAS-M	5.72	4.4	1/1	0.1	Child-Parent model
BATS (Medium)	5.60	5.4	1/1	0.25	Gradient-based
DCP-NAS-M	5.34	4.9	1/1	0.1	Gradient-based [‡]
CP-NAS-L	4.73	10.6	1/1	0.1	Child-Parent model
BATS (Large)	4.50	10.0	1/1	0.25	Gradient-based
DCP-NAS-L	4.22	11.2	1/1	0.1	Gradient-based [‡]

both removed by 0.9% Top-1 accuracy. Thus, we keep the training protocol consistent with previous works [80, 43, 3].

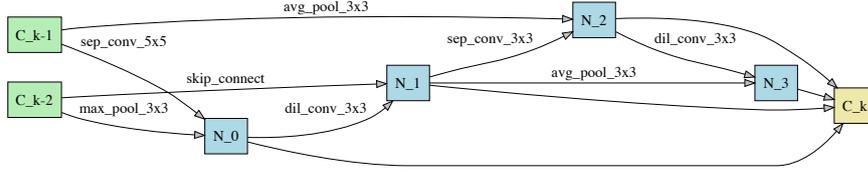
4.3 Results on Image Classification

4.3.1 Results on CIFAR-10 datasets

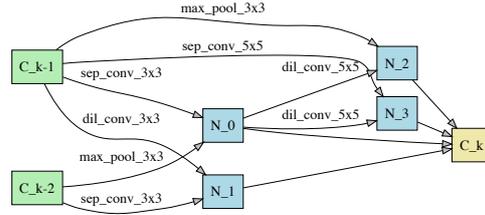
We first evaluate our DCP-NAS on CIFAR-10 and compare results with both manually designed networks [70, 26] and networks searched by NAS at different levels of binarization [43, 67, 9, 80, 3].

The results for different architectures on CIFAR-10 are summarized in Tab. 10. We search for three 1-bit CNNs with different model sizes, which binarize both weights and activations. Note that for the model size, we also consider the number of bits of each parameter. The 1-bit CNNs only need 1 bit to save and compute the weight parameter or the activation, while the real-valued networks need 32.

Compared with manually designed networks, *e.g.*, WRN-22 (BONN) [72], our DCP-NAS achieves comparable or smaller test error (5.34% vs. 8.07%) and similar model size (4.9Mbits vs. 4.3Mbits). Compared with real-valued networks obtained by other NAS methods, our DCP-NAS achieves comparable test errors and significantly more efficient models, with similar or less search time. Compared with state-of-the-art BATS [3], our DCP-NAS achieves 0.17%/0.26%/0.28% performance improvements with different model sizes,



(a) The normal cell



(b) The reduction cell

Fig. 9 The normal cell (a) and the reduction cell (b) searched for ImageNet.

producing a new state-of-the-art. Moreover, our DCP-NAS consumes 60% lower search cost than BATS. The superior results prove our assumption that our DCP-NAS method can learn a well-performed binary architecture from the tangent space.

In terms of search efficiency, compared with the previous work [3], our DCP-NAS is 60.0% faster (tested on our platform - 8 NVIDIA GEFORCE RTX 2080Ti GPUs). We attribute our superior results to the proposed scheme of GGN-based Hessian matrix acceleration.

4.3.2 Results on ImageNet ILSVRC12

To further evaluate the performance of our DCP-NAS, we compare our method with state-of-the-art image classification methods on the ImageNet, including hand-crafted ReActNet [46], BATS [3], and CP-NAS [80]. We also report multi-bit ResNets [26] and NAS methods [8, 43, 67, 9] for reference. We use memory usage and operations (OPs) [47] for the efficiency comparison. All the searched networks are obtained directly by DCP-NAS on ImageNet by stacking the cells. All the training settings of α are the same as other weight parameters. The 1-bit models are trained from scratch. Due to the first convolutional layer in a depth-wise separable convolution with fewer parameters, we do not binarize the activations of the first layers for ImageNet. Tab. 11 shows the test accuracy on ImageNet.

We observe that DCP-NAS significantly outperforms previous binary NAS methods with three different model sizes, by significant margins. With the

Table 11 A performance comparison with SOTAs on ImageNet. W/A denotes the bit length of weights and activations. We report the Top-1 (%) and Top-5 (%) accuracy. Err. and Mem. denote Test Error and Memory Usage (Mbits), respectively. Search cost here is compared through GPU days. ^b denotes models trained with an improved training scheme in BNAS₂ v2 [32]. [†] denotes models trained with an improved training scheme and EBConv in EBN [4]. Gradient-based[‡] denotes our DCP-NAS is a discrepancy-aware gradient-based search method.

Architecture	Accuracy(%)		Mem.	OPs (M)	W/A	Search Cost	GPU Mem. (per-image/ max Mem.)	Search Method
	Top-1	Top-5						
ResNet-18	69.3	89.2	357.4	1819	32/32	-	-	Manual
PNAS	74.2	91.9	163.2	588	32/32	225	-	SMBO
DARTS	73.1	91.0	156.8	574	32/32	4	0.01/5.4	Gradient-based
PC-DARTS	75.8	92.7	169.6	597	32/32	3.8	-	Gradient-based
ResNet-18 (PCNN)	63.5	85.1	11.2	1819	1/32	-	-	Manual
BNAS ₁	71.3	90.3	6.2	778	1/32	2.6	0.01/5.4	Performance-based
BNAS ₂ -D	57.7	79.9	-	148	1/1	-	-	Gradient-based
BNAS ₂ v2-D	64.6	84.9	-	148	1/1	-	-	Gradient-based
BATS	60.4	83.0	-	98	1/1	-	-	Gradient-based
DCP-NAS-S	63.0	84.5	4.2	100	1/1	2.9	0.025/13.1	Gradient-based [‡]
DCP-NAS-S^b	66.9	87.8	4.2	100	1/1	2.9	0.025/13.1	Gradient-based [‡]
DCP-NAS-S[†]	67.2	88.1	4.2	100	1/1	2.9	0.025/13.1	Gradient-based [‡]
BNAS ₂ -G	59.8	81.6	-	193	1/1	-	-	Gradient-based
BNAS ₁	64.3	86.1	6.4	-	1/1	3.2	0.01/5.4	Performance-based
ResNet-18 (ReActNet)	65.9	86.2	11.2	165	1/1	-	-	Manual
BATS (2×-wider)	66.1	87.0	-	155	1/1	-	-	Gradient-based
BARS-D	54.6	79.2	-	129	1/1	-	-	Gradient-based
BARS-E	56.2	80.6	-	183	1/1	-	-	Gradient-based
NASB	60.5	82.2	-	171	1/1	-	-	Gradient-based
NASBV4	65.3	85.9	-	281	1/1	-	-	Gradient-based
EBN	67.5	87.5	-	131	1/1	-	-	Gradient-based
EBN [†]	71.2	90.1	-	131	1/1	-	-	Gradient-based
EBNAS-L	67.8	87.4	-	162	1/1	-	-	Gradient-based
DCP-NAS-M	69.0	88.4	6.9	163	1/1	2.9	0.025/13.1	Gradient-based [‡]
DCP-NAS-M^b	72.3	91.1	6.9	163	1/1	2.9	0.025/13.1	Gradient-based [‡]
DCP-NAS-M[†]	72.7	91.3	6.9	163	1/1	2.9	0.025/13.1	Gradient-based [‡]
BNAS ₂ -H	63.5	83.9	-	656	1/1	-	-	Gradient-based
BARS-F	60.3	81.9	-	293	1/1	-	-	Gradient-based
NASBV5	66.6	87.0	-	352	1/1	-	-	Gradient-based
CP-NAS-L	66.5	86.8	12.5	323	1/1	2.8	0.01/5.4	Child-Parent model
DCP-NAS-L	72.4	90.6	14.4	334	1/1	2.9	0.025/13.1	Gradient-based [‡]
DCP-NAS-L^b	75.0	92.9	14.4	334	1/1	2.9	0.025/13.1	Gradient-based [‡]
DCP-NAS-L[†]	75.2	93.3	14.4	334	1/1	2.9	0.025/13.1	Gradient-based [‡]

small model, our DCP-NAS improves the Top-1 accuracy by 2.6%, 12.8%, 3.7%, 1.5% and 1.2% compared with BATS, BARS-E, NASBV4, EBN and EBNAS-L. Also, our DCP-NAS improves the Top-1 accuracy by 4.7%, and 2.9% compared with BNAS₁ and BATS for the medium model. In large networks, DCP-NAS achieves 72.4% vs. 66.5% (CP-NAS). As shown in Fig. 9,

Table 12 Comparison of mAP(%) and Rank@1(%) for person re-identification on Market-1501, DukeMTMC-reID, and CUHK03 Datasets. We compare with the BiRe-ID using vanilla ResNets as backbone in [65]. Gradient-based[‡] denotes our DCP-NAS is a discrepancy-aware gradient-based search method.

Backbone	Quantization	W/A (bit)	Market-1501		DukeMTMC -reID		CUHK03		Search Method
			mAP	R1	mAP	R1	mAP	R1	
ResNet-18	Real-valued	32/32	64.3	85.1	48.1	72.9	33.0	42.7	Manual
	XNOR-Net		40.1	63.8	24.9	52.1	14.0	26.1	
	Bi-Real Net	1/1	52.7	73.1	36.3	64.9	25.6	36.0	
	BiRe-ID		64.0	84.1	47.1	70.7	31.8	40.9	
DCP-NAS-S		1/1	68.0	84.9	60.2	76.4	48.8	50.0	Gradient-based [‡]
ResNet-34	Real-valued	32/32	68.0	86.7	52.2	73.9	39.2	47.9	Manual
	XNOR-Net		45.9	67.1	27.7	55.7	18.4	29.6	
	Bi-Real Net	1/1	54.1	77.6	40.8	69.1	29.3	38.2	
	BiRe-ID		67.1	85.3	51.0	72.4	37.1	46.4	
DCP-NAS-M		1/1	70.6	86.2	61.9	77.9	51.4	53.6	Gradient-based [‡]
ResNet-50	Real-valued	32/32	71.6	88.8	56.9	75.4	43.8	50.6	Manual
	XNOR-Net		49.1	68.3	31.1	59.4	20.3	31.5	
	Bi-Real Net	1/1	58.9	78.5	45.1	71.0	34.7	42.6	
	BiRe-ID		69.2	86.9	55.0	74.1	41.1	49.0	
DCP-NAS-L		1/1	71.9	87.3	62.4	78.7	54.4	57.6	Gradient-based [‡]

our DCP-NAS obtains more optimized architecture for binarized weights and activations, which increase the representation ability of 1-bit CNNs. It is noteworthy that our DCP-NAS-L trained with improved training scheme and EBN following [4] achieves 75.2% top-1 accuracy on ImageNet.

Note that compared to the human-designed real-valued networks, our DCP-NAS achieves better performance but with higher compression. Furthermore, to obtain a better performance, we do not binarize the activations of the preprocessing operations for the two input nodes. We achieve an accuracy of 72.4%, which surpass the real-valued hand-crafted model, *e.g.*, 69.3% for ResNet-18.

4.4 Results on Person Re-identification

We implement our experiments on three mainstream datasets of person Re-ID task, *i.e.*, Market1501, DukeMTMC-reID, and CUHK03. We report the performance of real-valued models, previous binarization methods XNOR-Net [5] and Bi-Real Net [47], for comparison. We also compare our DCP-NAS with state-of-the-art binarization person reidentification work BiRe-ID [65]. Three mainstream backbones are chosen for comparison: ResNet-18, ResNet-34, and ResNet-50.

Table 13 Comparison of mAP(%) for objects with state-of-the-art 1-bit object detectors in Faster-RCNN detection framework on VOC `test2007`, where the performances of real-valued and 4-bit detectors are reported for reference. Gradient-based[‡] denotes our DCP-NAS is a discrepancy-aware gradient-based search method.

Backbone	Quantization	W/A(bit)	mAP (%)	Mem. (MB)	OPs ($\times 10^9$)	Search Method
ResNet-18	Real-valued	32/32	76.4	112.88	96.40	Manual
	DoReFa-Net	4/4	73.3	21.59	27.15	
	ReActNet	1/1	69.6	16.61	18.49	
	LWS-Det		73.2			
	DCP-NAS-S	1/1	72.8	13.42	15.01	Gradient-based [‡]
DCP-NAS-M	1/1	74.2	13.56	15.10	Gradient-based [‡]	
ResNet-34	Real-valued	32/32	77.8	145.12	118.80	Manual
	DoReFa-Net	4/4	75.6	29.65	32.31	
	ReActNet	1/1	72.3	24.68	21.49	
	LWS-Det		75.8			
	DCP-NAS-L	1/1	76.7	13.79	15.29	Gradient-based [‡]

Market1501: As shown in the column 3 of Tab. 12, on ResNet-18, ResNet-34 and ResNet-50 backbones, our DCP-NAS-S, DCP-NAS-M and DCP-NAS-L both surpass the real-valued counterpart in a clear margin (68.0 % mAP vs. 64.3 % mAP, 70.6 % mAP vs. 68.0 % mAP, 71.9 % mAP vs. 71.6 % mAP). Our DCP-NAS outperforms other 1-bit models (XNOR-Net and Bi-Real Net) by a sizable margin. Moreover, our DCP-NAS outperforms BiRe-ID by 4.0%, 3.5%, 2.7% mAP with different model sizes, respectively.

DukeMTMC-reID: On DukeMTMC-reID dataset, our DCP-NAS also achieves significant performances, as listed in the column 4 of Tab. 12. On ResNet-18, ResNet-34 and ResNet-50 backbones, our DCP-NAS-S, DCP-NAS-M and DCP-NAS-L both outperforms the real-valued models by 12.1%, 9.7%, 5.5% mAP. Our DCP-NAS outperforms other 1-bit models (XNOR-Net and Bi-Real Net) by a sizable margin. Moreover, our DCP-NAS outperforms BiRe-ID by 13.1%, 10.9%, 7.4% mAP with different model sizes, respectively.

CUHK03: As shown in the column 5 of Tab. 12, our DCP-NAS shows significantly representation ability on person ReID task in the CUHK03 dataset. With ResNet-18, ResNet-34 and ResNet-50 backbones, our DCP-NAS-S, DCP-NAS-M and DCP-NAS-L both achieve clear performance advantage than the real-valued models (48.8 % mAP vs. 33.0 % mAP, 51.4 % mAP vs. 39.2 % mAP, 54.4 % mAP vs. 43.8 % mAP). Our DCP-NAS surpasses other 1-bit models (XNOR-Net and Bi-Real Net) by a sizable margin. Moreover, our DCP-NAS outperforms BiRe-ID by 15.8%, 12.2%, 10.6% mAP with different model sizes, respectively.

Table 14 Comparison of mAP@[.5, .95](%), AP with different IoU threshold and AP for objects in various sizes with state-of-the-art 1-bit object detectors in Faster-RCNN framework on COCO `minival`, where the performances of real-valued and 4-bit detectors are reported for reference. Gradient-based[‡] denotes our DCP-NAS is a discrepancy-aware gradient-based search method.

Backbone	Quantization	W/A(bit)	mAP @[.5, .95]	AP ₅₀	AP ₇₅	AP _s	AP _m	AP _l	Search Method
	Real-valued	32/32	32.2	53.8	34.0	18.0	34.7	41.9	
ResNet-18	FQN	4/4	28.1	48.4	29.3	14.5	30.4	38.1	Manual
	ReActNet		21.1	38.5	20.5	9.7	23.5	32.1	
	LWS-Det	1/1	26.9	44.9	27.7	12.9	28.7	38.3	
	DCP-NAS-S	1/1	25.4	42.9	24.6	14.6	26.3	33.8	Gradient-based [‡]
	DCP-NAS-M	1/1	27.3	45.2	27.4	17.4	28.5	35.6	Gradient-based [‡]
	DCP-NAS-L	1/1	29.7	47.8	29.9	19.5	30.7	37.5	Gradient-based [‡]

4.5 Results on Object Detection

4.5.1 Results on the PASCAL VOC dataset

In this section, we compare the proposed DCP-NAS with other state-of-the-art 1-bit detectors based on Faster-RCNN, including ReActNet [46], and LWS-Det [66], on the same framework for the task of object detection on the PASCAL VOC datasets. We also report the detection performance of the multi-bit quantized network DoReFa-Net [78]. Tab. 13 illustrates the comparison of the mAP across different quantization methods and detection frameworks. Our DCP-NAS efficiently generalizes to the downstream task and significantly accelerates the computation and saves the storage on various detectors by binarizing the activations and weights to 1-bit. The results for 1-bit Faster-RCNN on VOC `test2007` are summarized in Tab. 13. Compared with other 1-bit methods, we observe significant performance improvements with our DCP-NAS over other state-of-the-arts. Compared with the ResNet-18 backbone, our DCP-NAS-S outperforms ReActNet 3.2% mAP with the same bit-width activations and weights, but higher acceleration rate (15.01 GOPs vs. 18.49 GOPs). Our DCP-NAS-M achieves surpasses other low-bit detectors including DoReFa-Net, ReActNet and LWS-Det by 0.9%, 4.6% and 1.0%. The lines 8 to 12 in Tab. 13 illustrate that our DCP-NAS-L achieves a performance far more close to the real-valued counterpart compared with other prior low-bit detectors with ResNet-34 backbone. Quantitatively speaking, our DCP-NAS surpasses DoReFa-Net, ReActNet and LWS-Det by 1.1%, 4.4%, and 0.9% mAP, respectively. Moreover, DCP-NAS even achieves comparable performance as real-valued (76.7% vs. 77.8%) with high compression rate and speed acceleration, which demonstrates the superiority of our DCP-NAS.

In short, we achieved a new state-of-the-art performance compared to other 1-bit CNNs on various detection frameworks with various backbones on PASCAL VOC. We also surpass the performance to real-valued models, as demonstrated in extensive experiments, clearly validating the superiority of our DCP-NAS.

Table 15 Comparing DCP-NAS-M with real-valued ResNet-18 and 1-bit ResNet-18 (ReActNet) on hardware (single thread).

Model	W/A	Memory (MBits)	Memory Saving	Latency (ms)	Acceleration
ResNet-18	32/32	357.4	-	255.7	-
ResNet-18 (ReActNet)	1/1	11.2	31.9×	67.1	3.8×
DCP-NAS-M	1/1	6.9	51.8×	101.3	2.5×

4.5.2 Results on the COCO dataset

The COCO dataset is much more challenging for object detection than PASCAL VOC due to its diversity and scale. We compare the proposed DCP-NAS with state-of-the-art 1-bit CNNs, including XNOR-Net [53], Bi-Real-Net [44], and BiDet [59], on COCO. We also report the detection performance of the 4-bit quantized DoReFa-Net [78]. Tab. 14 shows mAP and AP with different IoU thresholds and AP of objects with different scales. Limited by the page width, we do not show the memory usage and FLOPs in Tab. 14. We conduct experiments on Faster-RCNN framework. Compared with the state-of-the-art 1-bit methods, our DCP-NAS outperforms other methods by significant margins. Our DCP-NAS-S improves the mAP@[.5, .95] by 4.3% compared with state-of-the-art ReActNet and our DCP-NAS-M surpasses the state-of-the-art 1-bit detectors ReActNet and LWS-Det by 6.2% and 0.4% on mAP@[.5, .95], respectively. Compared with the 4-bits detectors FQN, our DCP-NAS-L improves the mAP@[.5, .95] by 1.6% with binary weights and activations which can significantly compress the model parameters and accelerate the inference process. Our DCP-NAS-L also achieves much closer performance on mAP@[.5, .95] compared with real-valued ResNet-18 based Faster-RCNN (29.7% vs. 32.2%), but with extremely compressed detectors’ weights and activations. Similarly, on other APs with different IoU thresholds, our DCP-NAS outperforms other methods obviously.

To conclude, compared with the baseline methods of network quantization, our method achieves the best performance in terms of the AP with different IoU thresholds and the AP for objects in different sizes on COCO, demonstrating DCP-NAS’s superiority and universality in different application settings.

4.6 Deployment Efficiency

We implement the 1-bit models achieved by our DCP-NAS-M on ODROID C4, which has a 2.016 GHz 64-bit quad-core ARM Cortex-A55. With evaluating its real speed in practice, the efficiency of our DCP-NAS is proved when deployed into real-world mobile devices. We leverage the SIMD instruction SSSL on ARM NEON to make the inference framework BOLT [19] compatible with DCP-NAS. We compare DCP-NAS-M to the real-valued ResNet-18

and binarized ResNet-18 (ReActNet) in Tab. 15. We can see that DCP-NAS-M's inference speed is substantially faster with the highly efficient BOLT framework. For example, the acceleration rate achieves about $2.5\times$ compared with ResNet-18, which is slightly lower than the acceleration rate of ReActNet's $3.8\times$. However, our DCP-NAS-M has higher memory saving ratio, compared with ReActNet. The deployment efficiency experiment reveals the effectiveness and efficiency of our DCP-NAS. All deployment results are significant for the computer vision on real-world edge devices.

5 Conclusion and future work

This paper proposes Discrepant Child-Parent Neural Architecture search (DCP-NAS), which consider the architecture discrepancy between real-valued and binarized models, resulting in a unified binary neural architecture search framework with novel decoupled optimization. We incorporate tangent propagation, GGN method and decoupled backtracking optimization into the binary NAS in an efficient and effective manner.

Extensive experiments on CIFAR and ImageNet demonstrate that DCP-NAS achieves the best classification performance comparing with both previous binary NAS and directly binarizing real-valued NAS, and even surpasses real-valued ResNet-18 in a clear margin. We also achieve a promising performance on person re-identification and object detection, which validate the generality of our proposed method. In the future, we will extend our method with Transformer-based architectures to build more generalized binary NAS. We will also try other optimization methods to find the optimal architecture for more compact neural networks.

6 Acknowledgement

This work was supported in part by the National Natural Science Foundation of China under Grant 62076016 and 61827901, Beijing Natural Science Foundation-Xiaomi Innovation Joint Fund L223024, Foundation of China Energy Project GJNY-19-90. "One Thousand Plan" innovation leading talent funding projects in Jiangxi Province Jxsg2023102268

References

1. Bengio, Y., Léonard, N., Courville, A.: Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv :1308.3432* (2013)
2. Brock, A., Lim, T., Ritchie, J.M., Weston, N.: Smash: one-shot model architecture search through hypernetworks. *arXiv preprint arXiv:1708.05344* (2017)
3. Bulat, A., Martinez, B., Tzimiropoulos, G.: Bats: Binary architecture search. In: *Proc. of ECCV*, pp. 309–325 (2020)
4. Bulat, A., Martinez, B., Tzimiropoulos, G.: High-capacity expert binary networks. *arXiv* (2020)

5. Bulat, A., Tzimiropoulos, G.: Xnor-net++: Improved binary neural networks. In: Proc. of BMVC, pp. 1–12 (2019)
6. Cai, H., Chen, T., Zhang, W., Yu, Y., Wang, J.: Efficient architecture search by network transformation. In: Proc. of AAAI, pp. 2787–2794 (2018)
7. Cai, H., Yang, J., Zhang, W., Han, S., Yu, Y.: Path-level network transformation for efficient architecture search. In: Proc. of ICML, pp. 678–687 (2018)
8. Cai, H., Zhu, L., Han, S.: Proxylessnas: Direct neural architecture search on target task and hardware. In: Proc. of ICLR, pp. 1–13 (2018)
9. Chen, H., Zhang, B., Zheng, X., Liu, J., Doermann, D., Ji, R., et al.: Binarized neural architecture search. In: Proc. of AAAI, pp. 10526–10533 (2020)
10. Chen, H., Zhuo, L., Zhang, B., Zheng, X., Liu, J., Ji, R., Doermann, D., Guo, G.: Binarized neural architecture search for efficient object recognition. *International Journal of Computer Vision* **129**(2), 501–516 (2021)
11. Chen, X., Xie, L., Wu, J., Tian, Q.: Progressive differentiable architecture search: Bridging the depth gap between search and evaluation. In: Proc. of CVPR, pp. 1294–1303 (2019)
12. Chen, Y., Yang, T., Zhang, X., Meng, G., Xiao, X., Sun, J.: Detnas: Backbone search for object detection. In: Proc. of NeuIPS, pp. 6642–6652 (2019)
13. Courbariaux, M., Bengio, Y., David, J.P.: Binaryconnect: Training deep neural networks with binary weights during propagations. In: Proc. of NeuIPS, pp. 3123–3131 (2015)
14. Courbariaux, M., Hubara, I., Soudry, D., El-Yaniv, R., Bengio, Y.: Binarized neural networks: Training deep neural networks with weights and activations constrained to +1 or -1. *arXiv preprint arXiv:1602.02830* (2016)
15. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: Proc. of CVPR, pp. 248–255 (2009)
16. DeVries, T., Taylor, G.W.: Improved regularization of convolutional neural networks with cutout. *arXiv :1708.04552* (2017)
17. Everingham, M., Van Gool, L., Williams, C.K., Winn, J., Zisserman, A.: The pascal visual object classes (voc) challenge. *International journal of computer vision* **88**(2), 303–338 (2010)
18. Felzenszwalb, P.F., Girshick, R.B., McAllester, D., Ramanan, D.: Object detection with discriminatively trained part-based models. *IEEE transactions on pattern analysis and machine intelligence* **32**(9), 1627–1645 (2009)
19. Feng, J.: Bolt. <https://github.com/huawei-noah/bolt> (2021)
20. Gholami, A., Kim, S., Dong, Z., Yao, Z., Mahoney, M.W., Keutzer, K.: A survey of quantization methods for efficient neural network inference. *arXiv* (2021)
21. Gu, J., Li, C., Zhang, B., Han, J., Cao, X., Liu, J., Doermann, D.: Projection convolutional neural networks for 1-bit cnns via discrete back propagation. In: Proc. of AAAI, pp. 8344–8351 (2019)
22. Guo, Y., Yao, A., Chen, Y.: Dynamic network surgery for efficient dnns. In: Proc. of NeuIPS, pp. 1379–1387 (2016)
23. Ha, D., Dai, A., Le, Q.V.: Hypernetworks. *arXiv preprint arXiv:1609.09106* (2016)
24. Han, S., Pool, J., Tran, J., Dally, W.: Learning both weights and connections for efficient neural network. In: Proc. of NeuIPS, pp. 1135–1143 (2015)
25. He, K., Zhang, X., Ren, S., Sun, J.: Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In: Proc. of ICCV, pp. 1026–1034 (2015)
26. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proc. of CVPR, pp. 770–778 (2016)
27. Howard, A.G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M., Adam, H.: Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv :1704.04861* (2017)
28. Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q.: Densely connected convolutional networks. In: Proc. of CVPR, pp. 4700–4708 (2017)
29. Hubara, I., Courbariaux, M., Soudry, D., El-Yaniv, R., Bengio, Y.: Binarized neural networks. In: Proc. of NeuIPS, pp. 1–9 (2016)
30. Juefei-Xu, F., Naresh Boddeti, V., Savvides, M.: Local binary convolutional neural networks. In: Proc. of CVPR, pp. 19–28 (2017)
31. Kim, D., Singh, K.P., Choi, J.: Learning architectures for binary networks. In: Proc. of ECCV, pp. 575–591 (2020)

32. Kim, D., Singh, K.P., Choi, J.: Bnas v2: Learning architectures for binary networks with empirical improvements. arXiv preprint arXiv:2110.08562 (2021)
33. Krizhevsky, A., Nair, V., Hinton, G.: The cifar-10 dataset. online: <http://www.cs.toronto.edu/kriz/cifar.html> (2014)
34. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. Proc. of NeuIPS **25**, 1097–1105 (2012)
35. Li, H., Kadav, A., Durdanovic, I., Samet, H., Peter Graf, H.: Pruning filters for efficient convnets. In: Proc. of ICLR, pp. 1–13 (2017)
36. Li, W., Zhao, R., Xiao, T., Wang, X.: Deepreid: Deep filter pairing neural network for person re-identification. In: Proc. of CVPR, pp. 152–159 (2014)
37. Li, Y., Xu, S., Zhang, B., Cao, X., Gao, P., Guo, G.: Q-vit: Accurate and fully quantized low-bit vision transformer. arXiv (2022)
38. Li, Z., Ni, B., Zhang, W., Yang, X., Gao, W.: Performance guaranteed network acceleration via high-order residual quantization. In: Proc. of ICCV, pp. 2584–2592 (2017)
39. Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection. In: Proc. of CVPR, pp. 2117–2125 (2017)
40. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: Proc. of ECCV, pp. 740–755 (2014)
41. Lin, X., Zhao, C., Pan, W.: Towards accurate binary convolutional neural network. In: Proc. of NeuIPS, pp. 344–352 (2017)
42. Liu, C., Zoph, B., Neumann, M., Shlens, J., Hua, W., Li, L.J., Fei-Fei, L., Yuille, A., Huang, J., Murphy, K.: Progressive neural architecture search. In: Proc. of ECCV, pp. 19–34 (2018)
43. Liu, H., Simonyan, K., Yang, Y.: Darts: Differentiable architecture search. In: Proc. of ICLR, pp. 1–13 (2018)
44. Liu, Z., Luo, W., Wu, B., Yang, X., Liu, W., Cheng, K.T.: Bi-real net: Binarizing deep network towards real-network performance. International Journal of Computer Vision **128**(1), 202–219 (2020)
45. Liu, Z., Shen, Z., Li, S., Helweggen, K., Huang, D., Cheng, K.T.: How do adam and training strategies help bnns optimization. In: Proc. of ICML, pp. 6936–6946 (2021)
46. Liu, Z., Shen, Z., Savvides, M., Cheng, K.T.: Reactnet: Towards precise binary neural network with generalized activation functions. In: Proc. of ECCV, pp. 143–159 (2020)
47. Liu, Z., Wu, B., Luo, W., Yang, X., Liu, W., Cheng, K.T.: Bi-real net: Enhancing the performance of 1-bit cnns with improved representational capability and advanced training algorithm. In: Proc. of ECCV, pp. 722–737 (2018)
48. Novikov, A., Podoprikin, D., Osokin, A., Vetrov, D.P.: Tensorizing neural networks. In: Proc. of NeuIPS, pp. 442–450 (2015)
49. Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., Lerer, A.: Automatic differentiation in pytorch. In: Proc. of NeuIPS Workshops, pp. 1–4 (2017)
50. Petersen, K.B., Pedersen, M.S., et al.: The matrix cookbook. Technical University of Denmark **7**(15), 510 (2008)
51. Pham, H., Guan, M.Y., Zoph, B., Le, Q.V., Dean, J.: Efficient neural architecture search via parameter sharing. In: Proc. of ICML, pp. 4095–4104 (2018)
52. Phan, H., Liu, Z., Huynh, D., Savvides, M., Cheng, K.T., Shen, Z.: Binarizing mobilenet via evolution-based searching. In: Proc. of CVPR, pp. 13420–13429 (2020)
53. Rastegari, M., Ordonez, V., Redmon, J., Farhadi, A.: Xnor-net: Imagenet classification using binary convolutional neural networks. In: Proc. of ECCV, pp. 525–542 (2016)
54. Real, E., Aggarwal, A., Huang, Y., Le, Q.V.: Regularized evolution for image classifier architecture search. In: Proc. of AAAI, pp. 4780–4789 (2019)
55. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: Proc. of ICLR, pp. 1–15 (2015)
56. Tan, M., Le, Q.: Efficientnet: Rethinking model scaling for convolutional neural networks. In: Proc. of ICML, pp. 6105–6114 (2019)
57. Wang, N., Gao, Y., Chen, H., Wang, P., Tian, Z., Shen, C., Zhang, Y.: Nas-fcos: Fast neural architecture search for object detection. In: Proc. of CVPR, pp. 11943–11951 (2020)

58. Wang, X., Zhang, B., Li, C., Ji, R., Han, J., Cao, X., Liu, J.: Modulated convolutional networks. In: Proc. of CVPR (2018)
59. Wang, Z., Wu, Z., Lu, J., Zhou, J.: Bidet: An efficient binarized object detector. In: Proc. of CVPR, pp. 2049–2058 (2020)
60. Xie, S., Zheng, H., Liu, C., Lin, L.: Snas: stochastic neural architecture search. In: Proc. of ICLR, pp. 1–17 (2018)
61. Xu, S., Li, Y., Lin, M., Gao, P., Guo, G., Lu, J., Zhang, B.: Q-detr: An efficient low-bit quantized detection transformer. arXiv preprint arXiv:2304.00253 (2023)
62. Xu, S., Li, Y., Ma, T., Lin, M., Dong, H., Zhang, B., Gao, P., Lv, J.: Resilient binary neural network (2023)
63. Xu, S., Li, Y., Wang, T., Ma, T., Zhang, B., Gao, P., Qiao, Y., Lu, J., Guo, G.: Recurrent bilinear optimization for binary neural networks. arXiv (2022)
64. Xu, S., Li, Y., Zeng, B., Ma, T., Zhang, B., Cao, X., Gao, P., Lu, J.: Ida-det: An information discrepancy-aware distillation for 1-bit detectors. arXiv (2022)
65. Xu, S., Liu, C., Zhang, B., Lü, J., Guo, G., Doermann, D.: Bire-id: Binary neural network for efficient person re-id. ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM) **18**(1s), 1–22 (2022)
66. Xu, S., Zhao, J., Lu, J., Zhang, B., Han, S., Doermann, D.: Layer-wise searching for 1-bit detectors. In: Proc. of CVPR, pp. 5682–5691 (2021)
67. Xu, Y., Xie, L., Zhang, X., Chen, X., Qi, G.J., Tian, Q., Xiong, H.: Pc-darts: Partial channel connections for memory-efficient architecture search. In: Proc. of ICLR, pp. 1–13 (2019)
68. Xu, Z., Lin, M., Liu, J., Chen, J., Shao, L., Gao, Y., Tian, Y., Ji, R.: Recu: Reviving the dead weights in binary neural networks. In: Proc. of ICCV, pp. 5198–5208 (2021)
69. Xue, S., Wang, R., Zhang, B., Wang, T., Guo, G., Doermann, D.: Idarts: Interactive differentiable architecture search. In: Proc. of ICCV, pp. 1163–1172 (2021)
70. Zagoruyko, S., Komodakis, N.: Wide residual networks. In: Proc. of BMVC, pp. 1–13 (2016)
71. Zhang, X., Zou, J., He, K., Sun, J.: Accelerating very deep convolutional networks for classification and detection. IEEE Transactions on Pattern Analysis and Machine Intelligence **38**(10), 1943–1955 (2016)
72. Zhao, J., Xu, S., Zhang, B., Gu, J., Doermann, D., Guo, G.: Towards compact 1-bit cnns via bayesian learning. International Journal of Computer Vision **130**(2), 201–225 (2022)
73. Zhao, T., Ning, X., Shi, X., Yang, S., Liang, S., Lei, P., Chen, J., Yang, H., Wang, Y.: Bars: Joint search of cell topology and layout for accurate and efficient binary architectures. arXiv (2020)
74. Zheng, L., Shen, L., Tian, L., Wang, S., Wang, J., Tian, Q.: Scalable person re-identification: A benchmark. In: Proc. of ICCV, pp. 1116–1124 (2015)
75. Zheng, X., Ji, R., Tang, L., Wan, Y., Zhang, B., Wu, Y., Wu, Y., Shao, L.: Dynamic distribution pruning for efficient network architecture search. arXiv preprint arXiv:1905.13543 (2019)
76. Zheng, X., Ji, R., Tang, L., Zhang, B., Liu, J., Tian, Q.: Multinomial distribution learning for effective neural architecture search. In: Proc. of ICCV, pp. 1304–1313 (2019)
77. Zheng, Z., Zheng, L., Yang, Y.: A discriminatively learned cnn embedding for person reidentification. ACM transactions on multimedia computing, communications, and applications (TOMM) **14**(1), 1–20 (2017)
78. Zhou, S., Wu, Y., Ni, Z., Zhou, X., Wen, H., Zou, Y.: Dorefa-net: Training low bitwidth convolutional neural networks with low bitwidth gradients. arXiv preprint arXiv:1606.06160 (2016)
79. Zhu, B., Al-Ars, Z., Hofstee, H.P.: Nasb: Neural architecture search for binary convolutional neural networks. In: 2020 International Joint Conference on Neural Networks (IJCNN), pp. 1–8. IEEE (2020)
80. Zhuo, L., Zhang, B., Chen, H., Yang, L., Chen, C., Zhu, Y., Doermann, D.: Cp-nas: Child-parent neural architecture search for binary neural networks. In: Proc. of IJCAI, pp. 1033–1039 (2020)
81. Zoph, B., Le, Q.V.: Neural architecture search with reinforcement learning. In: Proc. of ICLR, pp. 1–16 (2017)

82. Zoph, B., Vasudevan, V., Shlens, J., Le, Q.V.: Learning transferable architectures for scalable image recognition. In: Proc. of CVPR, pp. 8697–8710 (2018)