

# DisWOT: Student Architecture Search for Distillation WithOut Training

Peijie Dong<sup>1†</sup> Lujun Li<sup>2†\*</sup> Zimian Wei<sup>1†</sup>

<sup>1</sup> National University of Defense Technology, <sup>2</sup> Chinese Academy of Sciences

<sup>1</sup>{dongpeijie, weizimian16}@nudt.edu.cn, <sup>2</sup>lilujunai@gmail.com

## Abstract

Knowledge distillation (KD) is an effective training strategy to improve the lightweight student models under the guidance of cumbersome teachers. However, the large architecture difference across the teacher-student pairs limits the distillation gains. In contrast to previous adaptive distillation methods to reduce the teacher-student gap, we explore a novel training-free framework to search for the best student architectures for a given teacher. Our work first empirically show that the optimal model under vanilla training cannot be the winner in distillation. Secondly, we find that the similarity of feature semantics and sample relations between random-initialized teacher-student networks have good correlations with final distillation performances. Thus, we efficiently measure similarity matrixs conditioned on the semantic activation maps to select the optimal student via an evolutionary algorithm without any training. In this way, our student architecture search for Distillation WithOut Training (DisWOT) significantly improves the performance of the model in the distillation stage with at least 180× training acceleration. Additionally, we extend similarity metrics in DisWOT as new distillers and KD-based zero-proxies. Our experiments on CIFAR, ImageNet and NAS-Bench-201 demonstrate that our technique achieves state-of-the-art results on different search spaces. Our project and code are available at <https://lilujunai.github.io/DisWOT-CVPR2023/>.

## 1. Introduction

Despite the remarkable achievements of Deep Neural Networks (DNNs) in numerous visual recognition tasks [58, 71–74, 76], they usually lead to heavy costs of memory, computation, and power at model inference due to their large numbers of parameters. To address this issue, Knowledge Distillation (KD) has been proposed as a means of transferring knowledge from a high-capacity teacher model to a low-capacity target student model, providing a more optimal accuracy-efficiency trade-off during runtime [5, 8, 84]. The

\*Corresponding author, † equal contribution, PD conducted main experiments, LL proposed ideas and led the project & writing.

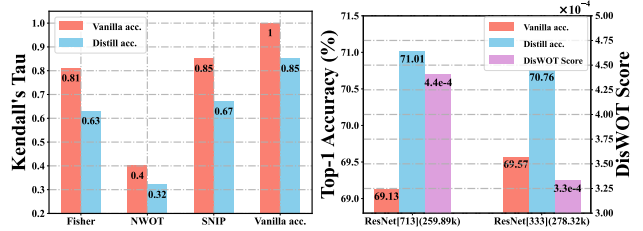


Figure 1. Left: Ranking correlation of proxies in zero-cost NAS with vanilla and distillation accuracy. Right: Vanilla accuracy, distillation accuracy, prediction scores of DisWOT for ResNet[7,1,3] and ResNet[3,3,3] on search space  $S_0$ .

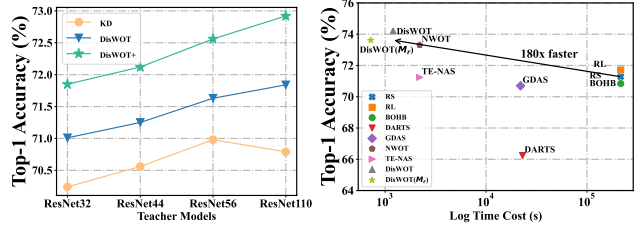


Figure 2. Left: KD [26], DisWOT, DisWOT+ results for ResNet20 under different teachers. Right: Comparison of distill accuracy & training efficiency with other NAS methods on NAS-Bench-201.

original KD method [26] utilizes the logit outputs of the teacher network as the source of knowledge. Subsequent studies [2, 25, 29, 31, 61, 68, 81] have focused on extracting informative knowledge based on intermediate feature representations. However, as the gap in capacity between students and teachers increases, existing KD methods are unable to improve results, particularly in tasks that depend on large-scale visual models such as ViT and GPT-3 [6, 18]. For example, as shown in Figure 2 (Left), the large teacher (e.g., ResNet110) lead to worse performance for the fixed student than the relatively smaller one (e.g., ResNet56).

To solve this issue, adaptive KD methods have been proposed in terms of training paradigms (e.g., early stop [11]) and architectural adaptations (e.g., assistant teacher [49] and architecture search [46]), respectively. However, they are ineffective in improving distillation performance or involve enormous training costs in the additional model training and

search process. In sharp contrast to these methods, we tackle this challenging problem from a new perspective regarding training-free architecture search. To achieve this goal, we construct a search space  $S_0$  for ResNet-like models with different depth configurations and obtain vanilla and distill performance for each candidate in  $S_0$  by individual training. Then, we evaluate the ranking correlation between predicted scores of training-free search methods and the actual performance of each student model. Surprisingly, as shown in Figure 1 (Left), there are common ranking correlation loss (10%  $\downarrow \sim$  20%  $\downarrow$ ) for these methods in predicting distillation accuracy than vanilla accuracy. To clarify this, we carefully analyze the disparities in vanilla and distillation performance for each model: (1) for overall search space, vanilla accuracy only preserves 85% correlations with actual distillation performance. (2) for a particular instance, as shown in Figure 1 (Right), ResNet20 with 3 res-blocks in each stage (i.e., ResNet[3,3,3]) has more parameters and better standalone performance but is weaker than ResNet[7,1,3] in the distillation process. Considering that ResNet[7,1,3] has more layers than ResNet20, we seek to understand the above phenomenon regarding the vanilla-distillation accuracy gap from the perspective of semantic matching [42]. ResNet[7,1,3] enjoys a larger effective receptive field and more excellent matched knowledge with teacher, resulting in significant distillation gains. Encouraged by this understanding, we strive to design a new zero-proxy regarding the semantic matching of teacher-student. As a result, we find that the similarity scores of feature semantics and sample relations can outperform conventional zero-cost NAS in predicting final distillation accuracy (see the comparison of ranking correlation on search space  $S_0$  in Table 8). As shown in Figure 1(Right), similarity scores are also consistent with distillation performance.

Drawing on the aforementioned observations, we introduce DisWOT, a simple yet effective training-free framework that finds the best student architectures for distilling the given teacher model. For better semantic matching in distillation, DisWOT leverages novel zero-cost metrics regarding the feature semantics and sample relations to select better student model. For the feature semantic similarity metric, we remark that randomly initialized models can localize objects well [7] and generate localization heatmaps via Grad-CAM [63] as reliable semantic information. Then, we measure the channel-wise similarity matrix of localization heatmaps and take the  $L_2$  distance of the similarity matrix for the teacher-student model as the metric. For input samples, different models have diverse abilities to discriminate their relationships. To improve relational knowledge matching ability, we use the  $L_2$  distance of sample-relation correlation matrix as a relation similarity metric. Finally, we search for student architectures using an evolutionary algorithm with semantic and relations similarity metrics. Then,

the distillation process is implemented between the searched student and the pre-defined teacher. In addition, we leverage these metrics directly as new distillers to enhance the student, as the DisWOT $^\dagger$ . Equipped with our train-free search and distillation design, our DisWOT and DisWOT $^\dagger$  framework significantly improve the model’s accuracy-latency tradeoff in inference with at least  $180\times$  training acceleration.

In principle, our DisWOT use higher-order statistics of teacher-student models to optimize the student architecture to fit a given teacher model. Its merits can be highlighted in three aspects: (1) In contrast to training-based student architecture search requires the individual or weight-sharing training, our DisWOT does not require the training of student models in the search phase. In addition, DisWOT is efficient to compute and easy to implement as it uses only the mini-batch data at initialization. (2) DisWOT is a teacher-aware search for distillation, which has better predictive distill accuracy than conventional NAS. (3) DisWOT exploits the distance of higher-order knowledge between the neural networks, bridging knowledge distillation and zero-proxy NAS. We further demonstrate the competitive ranking correlation of DisWOT among 10 knowledge distances in KD as zero-proxy for predicting vanilla accuracy in NAS-Bench-201. We anticipate that our work on KD-based zero-proxy can offer some assistance in furthering research endeavors related to KD and NAS.

We conduct extensive experiments on CIFAR-100, ImageNet, and the NAS-Bench-201 [16] dataset, demonstrating the superiority of our proposed approach. In contrast to experiments in traditional architectural search, we focus on final distillation accuracy instead of the vanilla accuracy for the student. The results show that our DisWOT can achieve better accuracy than traditional Zero-shot NAS in the same search space. Besides, by switching to a larger space, our DisWOT can obtain new state-of-the-art architectures. For example, in the same ResNet-like search space, we significantly improved 1.62% Top-1 accuracy over KD for ResNet50-ResNet18 pair under the same training settings. We also conducted comprehensive ablation studies to investigate how our method can use the predictability of zero-cost metrics to boost the distillation performance.

#### Main Contributions:

- By analyzing and exploring the discrepancy between teacher-student capability, we empirically show that their semantic similarities have a stronger correlation with the final distillation accuracy. This motivates us to propose a new student architecture search for the Distillation without Training (DisWOT) framework to reduce the teacher-student capability gap, which, to the best of our knowledge, is not achieved in the area of knowledge distillation.
- DisWOT proposes novel zero-cost metrics on similarity

of feature semantics and sample relations and ensemble these metrics to select the optimal student via an evolutionary algorithm at the initial time. In the distillation stage, DisWOT achieves state-of-the-art performances in multiple datasets and search spaces.

- We further expand 10 kinds of knowledge distances including DisWOT as new universal KD-based zero proxies, which enjoy competitive predictive power with actual performance of models. We hope that our contributions in this endeavor may aid to some degree in advancing future research on KD and NAS.

## 2. Related Work and Background

In this section, we summarize existing knowledge distillation and architecture search methods and clarify their differences to our method.

### 2.1. General Formulation of Knowledge Distillation

The fundamental concept underlying Knowledge Distillation (KD) involves utilizing acquired knowledge (e.g., logits [36], feature values [35, 37, 38, 40, 78], and sample relations [51, 66]) from a high-capacity teacher to guide the training of a student model. The training dataset  $(X, Y)$  comprises training samples  $X = x_{i=1}^n$  and their corresponding labels  $Y = y_{i=1}^n$ . Let  $f_T$  be the output logits of the fixed teacher  $T$  and let  $f_S$  be the output of student  $S$ , respectively. In KD, the student network  $f_S$  is trained by minimizing:

$$\mathcal{L}_S = \mathcal{L}_{CE}(f_S, Y) + \mathcal{L}_{KL}(f_S, f_T) + \mathcal{D}_f(\phi_S(x), \phi_T(x)), \quad (1)$$

where  $\mathcal{L}_{CE}$  is the regular cross-entropy loss.  $\mathcal{L}_{KL}$  represents Kullback-Leibler (KL) divergence.  $\mathcal{D}_f(\cdot, \cdot)$  is the distance function measuring the difference of intermediate feature representations (see Table 1 for particular distillers).

Table 1. Comparison of recent distillers.

Method	Knowledge	Distance $\mathcal{D}_f(\cdot, \cdot)$
FitNets [61]	Feature representation	$\mathcal{L}_2$
AT [81]	Attention maps	$\mathcal{L}_2$
CC [55]	Instance relation	$\mathcal{L}_2$
NST [29]	Neuron selectivity patterns	$\mathcal{L}_{MMD}$
PKT [52]	Similarity probability distribution	$\mathcal{L}_{KL}$

**Comparison with Other Adaptive KDs for Distillation Gap.** DisWOT is the first train-free architecture search solution to reduce the teacher-student gap. Unlike training manners [11] and KD-loss designs, DisWOT utilizes the classic KD training configurations and distillers. In addition, DisWOT is free from the assistant teacher in ATKD [49], which involves a complex training routine and budget. AKD [46] searches student via reinforcement learning based on feedback from individual training of lots of models. As a completely alternative technical route to these training-based

Table 2. Formulation of NAS methods.  $\mathcal{A}$  is the search space. A candidate architecture in the search space is denoted as  $\alpha \in \mathcal{A}$ , which corresponds to a neural architecture  $S(\alpha, w)$  with weight  $w$ .  $\mathcal{W}$  is the weight of the supernet.  $\mathcal{L}_{train}$  and  $\mathcal{L}_{val}$  are the loss functions on the training and validation sets, respectively.

Type	Evaluation	Formula
Training-based	Multi-trial Training	$\alpha^* = \arg \min_{\alpha \in \mathcal{A}} \mathcal{L}_{val}(S(\alpha, w_\alpha)),$ s.t. $w_\alpha = \arg \min_w \mathcal{L}_{train}(S(\alpha, w))$
	Weight sharing	$\alpha^* = \arg \min_{\alpha \in \mathcal{A}} \mathcal{L}_{val}(S(\alpha, \mathcal{W}_\mathcal{A}(\alpha))),$ s.t. $\mathcal{W}_\mathcal{A} = \arg \min_{\mathcal{W}} \mathcal{L}_{train}(S(\mathcal{A}, \mathcal{W}))$
Training-free	Zero-cost Proxy	$\alpha^* = \text{ZeroProxy}(S(\alpha, w))$ $\alpha \in \mathcal{A}$

NAS [20], our training-free DisWOT builds on new zero-proxy and achieves  $180\times \sim 1000\times$  training acceleration, which greatly improves its easy-to-use and flexibility.

Table 3. Comparison with different training-free NAS.

Type	Method	Teacher-aware	Objective
Prune-based	SNIP [33], Fisher [1], Synflow [65]	✗	Vanilla acc.
Activation-based	NWOT [48], Zen-NAS [41]	✗	Vanilla acc.
KD-based	DisWOT (ours)	✓	Distill acc.

### 2.2. Revisiting Architecture Search Methods

Neural Architecture Search (NAS) is emerged to reduce human efforts in architecture design and automate the discovery of high-performance networks. As formalized in Tab. 2, Multi-trial NAS methods [46, 87] train a large number of candidates individually, which leads to extensive resource consumption. To alleviate this, many NAS [9, 13, 28, 57] methods adopt a weight-sharing strategy within a single supernet to facilitate the simultaneous training of candidates. The supernet is trained for hundred of epochs by path sampling [12, 21] or compound optimization with architecture representations [44, 77]. As an orthogonal direction, zero-cost NAS methods [48, 79] focus on identifying well-performed architectures with training-free metrics. For example, NWOT [48] calculates the architecture score based on the kernel matrix of binary activations between small batches of samples.

**Comparison with NAS with Teacher.** Some training-based NAS [34, 56, 85] employ a teacher model to supervise supernet training to improve predictive ability in the search stage. However, these methods aim to improve vanilla accuracy, not for distillation, and they do not use the teacher model in the full training stage. In addition, without any training

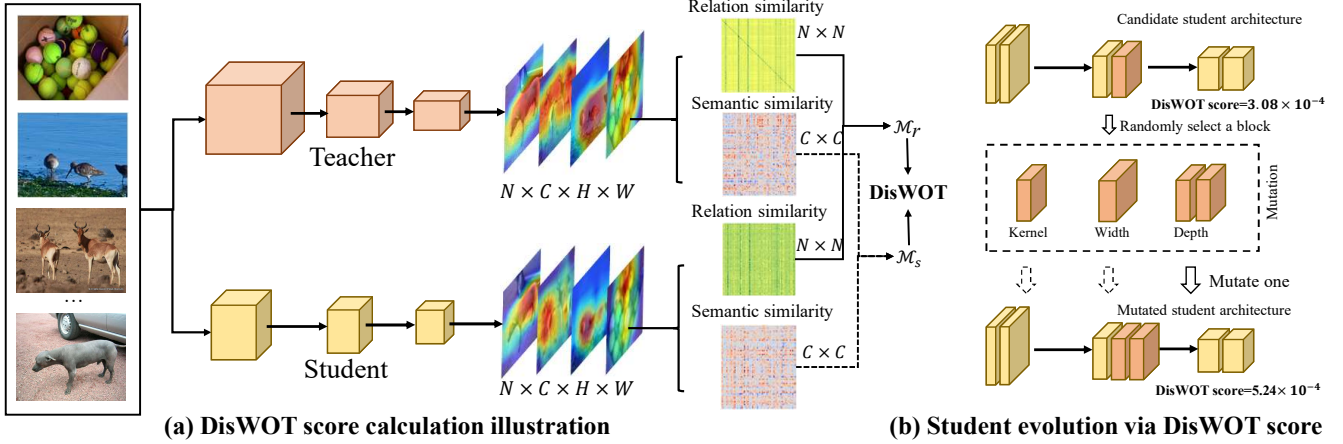


Figure 3. A schematic overview of our DisWOT, including (a) detailed calculation of the DisWOT scores and (b) evolution of the student architecture via the DisWOT scores. In search phase, DisWOT use semantic similarity metrics and relations similarity metrics to select good student for a given teacher. The semantic similarity metric is measured by  $l_2$  distance of the channel-wise correlation matrix for Grad-cam activation maps. Similarly, the relation similarity matrix statistics the sample-wise correlation matrix distance of the randomly initialized teacher-student pairs. With the feedback from these metrics, the evolutionary search in DisWOT automatically imitates good student from weak ones. In distillation phase, this searched student is distilled via teacher model and achieves superior gains.

costs, our training-free DisWOT enjoys obvious differences than these methods and advantages in efficiency.

**Compared to Other Training-free NAS.** Table 3 clearly summarizes the differences between DisWOT and other zero-cost methods [1, 33, 41, 48, 65]. Moreover, DisWOT outperforms these methods on distillation performance prediction and boosting in our sufficient experiments dealing with diverse datasets and search spaces.

### 3. Methodology

Figure 3 provides an overview of the DisWOT framework, which is comprised of two main stages: optimal student network search and distillation with high-order knowledge. In the search stage, we employ the neural architecture search technique to obtain an optimal student network for a pre-defined teacher network. Notably, we propose a training-free proxy called DisWOT to accurately rank enormous student networks and prevent expensive evaluation processes with high efficiency. In the distillation stage, the searched student network is retrained with distillation to imitate high-order knowledge in the teacher network. We give the details of these two designs in the following sections.

#### 3.1. Search for Optimal Student Network

We first present the training-free metrics we designed to score a student architecture, which indicates its final accuracy when distilled with a pre-defined teacher network. Then we depict the details of the evolutionary process to obtain an optimal student candidate.

**Semantic Similarity Metric.** The semantic information is meaningful for neural networks to perceive as humans. In

distillation, the teacher network always has more convolutional operations than the student, resulting in a teacher feature map with a larger receptive field and greater richness of semantic information. In contrast to distiller designs to alleviate semantic gaps, we aim for train-free student architecture to better match the teacher model with computational constraints. We notice that the network with random initial weights also has some semantic localization capability. Thus, we start to analyze the localization performance of the randomly initialized teacher-student model. Specifically, we utilize Grad-CAM maps [86] to localize semantic object regions, which explains the model decisions using gradient information. Given a mini-batch of input images, we define the high-level feature map before the Global Average Pooling (GAP) layer of the teacher network  $T$  as  $A_T \in \mathbf{R}^{B \times C_T \times H_T \times W_T}$ , where  $B$  represents the batch size,  $C_T$  denotes the number of output channels, and  $H_T$  and  $W_T$  are the spatial dimensions. Additionally, we introduce  $A_T^c \in \mathbf{R}^{N \times H_T \times W_T}$  as the  $c$ -th spatial map along the channel dimension. For the student network  $S_i$ , we have feature map  $A_{S_i}^c \in \mathbf{R}^{B \times C_S \times H_S \times W_S}$  and spatial map  $A_{S_i}^c \in \mathbf{R}^{B \times H_S \times W_S}$ , respectively. To compute the Grad-CAM maps of the  $n$ -th class for both the teacher and student networks, we can use the following formulations:

$$G_T = \sum_{c=1}^{C_T} w_{n,c}^T A_T^c, \quad G_{S_i} = \sum_{c=1}^{C_S} w_{n,c}^S A_{S_i}^c, \quad (2)$$

where  $w^T \in \mathbf{R}^{N \times C_T}$  and  $w^S \in \mathbf{R}^{N \times C_S}$  are weights of the last fully-connected layer in the teacher and student network.  $N$  represents the number of classes.  $w_{n,c}^T$  and  $w_{n,c}^S$  refer to the element located in the  $n$ -th row and  $c$ -th column of weight matrices  $w^T$  and  $w^S$ , respectively. To quantify



the intersection of class-discriminative localization maps, we formulate semantic similarity metric  $\mathcal{M}_s$  as the inter-correlation on the accumulated Grad-CAM maps for both teacher and student networks as follows:

$$\mathcal{G}^T = \frac{(G_T) \cdot (G_T)^\top}{\|(G_T) \cdot (G_T)^\top\|_2}, \mathcal{G}^S = \frac{(G_S) \cdot (G_S)^\top}{\|(G_S) \cdot (G_S)^\top\|_2}, \quad (3)$$

$$\mathcal{M}_s = \|\mathcal{G}^T - \mathcal{G}^S\|_2. \quad (4)$$

**Relation Similarity Metric.** The relationships between input samples are non-trivial for knowledge transfer. To reduce the teacher-student gap and improve the relation-distillation performance, we use the correlation matrix as the sample-wise metric to search for an optimal student network. For the random teacher network  $T$  and student network  $S_i$  with activation maps  $A_T \in \mathbf{R}^{N \times C_T \times H_T \times W_T}$  and  $A_S^i \in \mathbf{R}^{N \times C_i \times H_i \times W_i}$ , the correlation matrix of the mini-batch samples in the teacher network is formulated as follows:

$$\mathcal{A}^T = \frac{(\tilde{A}_T) \cdot (\tilde{A}_T)^\top}{\|(\tilde{A}_T) \cdot (\tilde{A}_T)^\top\|_2}, \mathcal{A}^{S_i} = \frac{(\tilde{A}_S) \cdot (\tilde{A}_S)^\top}{\|(\tilde{A}_S) \cdot (\tilde{A}_S)^\top\|_2}, \quad (5)$$

where  $\tilde{A}_T \in \mathbf{R}^{N \times CHW}$  is a reshaping of  $A_T$ , and  $M_T$  is a  $N \times N$  matrix. Thus, the  $(i, j)$  entry in matrix  $C_T$  represents the similarity between the  $i$ -th and  $j$ -th images within the mini-batch. Based on this, the sample similarity metric  $\mathcal{M}_r$  for a potential student model  $S_i$  is defined as follows:

$$\mathcal{M}_r = \|\mathcal{A}^T - \mathcal{A}^{S_i}\|_2. \quad (6)$$

**Training-Free Evolutionary Search.** Based on the above metric, we conduct a training-free evolutionary search algorithm to efficiently discover the optimal student  $\alpha^*$  from search space  $\mathcal{A}$ , as:

$$\alpha^* = \arg \min_{\alpha \in \mathcal{A}} (\mathcal{M}_s + \mathcal{M}_r). \quad (7)$$

**Theoretical Understanding.** According to the VC theory [70], the classification error of the vanilla teacher-student network can be decomposed as follows:

$$R(f_s) - R(f_r) \leq O\left(\frac{|\mathcal{F}_s|_C}{n^{\alpha_{sr}}}\right) + \epsilon_{sr}; R(f_t) - R(f_r) \leq O\left(\frac{|\mathcal{F}_t|_C}{n^{\alpha_{tr}}}\right) + \epsilon_{tr}, \quad (8)$$

where  $f_s \in \mathcal{F}_s$  is the student function,  $f_t \in \mathcal{F}_t$  is the teacher function, and  $f_r \in \mathcal{F}_r$  is the target function.  $R$  is the error.  $O(\cdot)$  and  $\epsilon_{sr}$  terms are the estimation and approximation error, respectively.  $O(\cdot)$  is related to the statistical procedure when given the number of data points. In contrast,  $\epsilon_{sr}$  is the approximation error of the student function class  $\mathcal{F}_s$  for  $f_r \in \mathcal{F}_r$ .  $|\cdot|_C$  is a function class capacity measure, and  $n$  is the number of data point. During distillation, the student

network is supervised purely with the teacher network as follows:

$$R(f_s) - R(f_t) \leq O\left(\frac{|\mathcal{F}_s|_C}{n^{\alpha_{st}}}\right) + \epsilon_{st}, \quad (9)$$

where  $\alpha_{st}$  and  $\epsilon_{st}$  are associated to student learning from teacher. By combining Equations 3.1 and 9, we obtain:

$$R(f_s) - R(f_r) \leq O\left(\frac{|\mathcal{F}_t|_C}{n^{\alpha_{tr}}}\right) + \epsilon_{tr} + O\left(\frac{|\mathcal{F}_s|_C}{n^{\alpha_{st}}}\right) + \epsilon_{st}. \quad (10)$$

When student obtains gains in KDs, its upper bound of error in distillation is smaller than vanilla training, which satisfies the following inequality:

$$O\left(\frac{|\mathcal{F}_t|_C}{n^{\alpha_{tr}}} + \frac{|\mathcal{F}_s|_C}{n^{\alpha_{st}}}\right) + \epsilon_{tr} + \epsilon_{st} \leq O\left(\frac{|\mathcal{F}_s|_C}{n^{\alpha_{sr}}}\right) + \epsilon_{sr}. \quad (11)$$

Based on the assumption in [26] that  $\epsilon_{tr} + \epsilon_{st} \leq \epsilon_{sr}$  holds consistently, we focus on minimizing  $O\left(\frac{|\mathcal{F}_t|_C}{n^{\alpha_{tr}}} + \frac{|\mathcal{F}_s|_C}{n^{\alpha_{st}}}\right)$  to improve the distillation performance. As noted in Lopez-Paz et al [47], a better representation allows for a faster learning rate with a fixed amount of data. Hence, when there is a larger gap between the capacities of the student and teacher networks, the value of  $\alpha_{st}$  tends to be lower. Thus we aim to search for an optimal student network that meets the requirement of  $\alpha_{s_{it}} \leq \alpha_{s_{ot}}$ , where  $s_i$  is all candidate student networks, and  $s_o$  is our searched student network. In this case, the inequality becomes more effective, and we improve the knowledge distillation by injecting a larger  $\alpha_{s_{ot}}$ . Specifically, we present the overall procedure for discovering optimal student in algorithm 1.

**Effects of Search Strategies.** We compare the evolution search algorithm and the random search algorithm in search space  $\mathcal{S}_2$  with the same number of iterations, as shown in Figure 4. We find that the evolution search algorithm can consistently find architectures with lower DisWOT, especially when the search space is relatively large, and the evolutionary search can explore better architectures.

### 3.2. Distillation with High-order Knowledge

In the distillation stage, teacher model  $T$  is employed to distill the optimal student network  $f_s$ . To verify the superiority of our search architecture, we adopt the existing distillers (e.g., KD) as the default distillation setting. In addition, we observe that the metrics we searched for actually serve as minimization optimization goals in the distillation process to transfer the teacher's privileged semantic and sample relational knowledge as the semantic distillation and sample distillation:

$$\mathcal{L}_{\mathcal{M}_s} = \frac{1}{c^2} \|\mathcal{G}^T - \mathcal{G}^S\|_2, \mathcal{L}_{\mathcal{M}_r} = \frac{1}{b^2} \|\mathcal{A}^T - \mathcal{A}^S\|_2, \quad (12)$$

Finally, we involve these advanced distillers in our framework, called DisWOT<sup>†</sup>. The total loss for DisWOT and

**Algorithm 1** Evolution Search for DisWOT

**Input:** Search space  $\mathcal{S}$ , population  $\mathcal{P}$ , architecture constraints  $\mathcal{C}$ , max iteration  $\mathcal{N}$ , sample ratio  $r$ , sampled pool  $\mathcal{Q}$ , topk  $k$ , teacher network  $\mathcal{T}$ .

**Output:** Highest DisWOT score architecture.

```

1:  $\mathcal{P}_0 := \text{Initialize population}(\mathcal{P}, \mathcal{C})$ ;
2: sample pool  $\mathcal{Q} := \emptyset$ ;
3: for  $i = 1 : \mathcal{N}$  do
4:   Clear sample pool  $\mathcal{Q} := \emptyset$ ;
5:   Randomly select  $r \times \mathcal{P}$  subnets  $\hat{P}_i \in \mathcal{P}$  to get  $\mathcal{Q}$ ;
6:   Candidates  $\{A_i\}_k := \text{GetTopk}(\mathcal{Q}, k)$ ;
7:   Parent  $A_i := \text{RandomSelect}(\{A_i\}_k)$ ;
8:   Mutate  $\hat{P}_i := \text{MUTATE}(A_i)$ ;
9:   if  $\hat{P}_i$  do not meet the constraints  $\mathcal{C}$  then
10:    Do nothing;
11:   else
12:    Get DisWOT-Score  $z := \text{DisWOT}(\hat{P}_i, \mathcal{T})$ ;
13:    Append  $\hat{P}_i$  to  $\mathcal{P}$ ;
14:   end if
15:   Remove network of smallest DisWOT-score;
16: end for

```

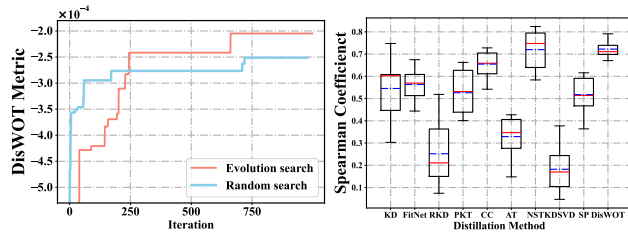


Figure 4. Left: comparison of random search and evolution search. Right: ranking correlation of different distillation methods on NAS-Bench-201.

Table 4. Spearman correlation  $\rho$  (%) on NAS-Bench-201.

Type	Method	$\rho$	Method	$\rho$
Zero-cost Proxies	Grad_Norm [1]	58.70 $\pm$ 0.11	Synflow [65]	<b>74.61<math>\pm</math>0.08</b>
	SNIP [33]	58.17 $\pm$ 0.15	Jacob [64]	73.42 $\pm$ 0.03
	Fisher [1]	35.91 $\pm$ 0.09	Zen-NAS [41]	41.36 $\pm$ 0.06
	NWOT [48]	64.41 $\pm$ 0.08	FLOPs [1]	63.38 $\pm$ 0.06
KD-based Proxies	KD [26]	54.43 $\pm$ 0.09	PKT [53]	52.65 $\pm$ 0.09
	FitNets [62]	56.18 $\pm$ 0.09	CC [54]	65.90 $\pm$ 0.08
	SP [69]	51.24 $\pm$ 0.08	NST [30]	72.35 $\pm$ 0.09
	RKD [50]	25.71 $\pm$ 0.17	DisWOT	<b>72.36<math>\pm</math>0.02</b>

DisWOT $^\dagger$  as:

$$\begin{aligned} \mathcal{L}_{\text{DisWOT}} &= \mathcal{L}_{CE}(f_S, Y) + \mathcal{L}_{KL}(f_S, f_T), \\ \mathcal{L}_{\text{DisWOT}^\dagger} &= \mathcal{L}_{\text{DisWOT}} + \mathcal{L}_{\mathcal{M}_s} + \mathcal{L}_{\mathcal{M}_r}. \end{aligned} \quad (13)$$

### 3.3. Bridging Distiller and Zero-proxy

In our DisWOT framework, we use the semantic and relational similarity metrics as a distillation performance predictor and distiller. In addition, DisWOT also enjoys good performance for vanilla performance predictions. Encour-

aged by this intriguing observation, we employ the knowledge function in of different KDs as zero-proxies and evaluate their ranking consistency with vanilla accuracy. As shown in Table 4, these KD-based zero-proxies enjoy competitive rankings with other NAS methods. Detailed results in Figure 4 illustrate that our DisWOT and NST [29] are the winners in the family of KD-based proxies. These attempts reveal the close connections between KD and NAS, and augment 10+ new universal proxies from the teacher-student learning perspective for training-free NAS research.

## 4. Experimental results

In this section, we present the experimental results of our DisWOT on different datasets. First, we describe the four datasets used in our experiments and three search spaces  $S_0, S_1, S_2$  in Sec. 4.1. Then, we conduct a comprehensive set of experiments to evaluate the effectiveness of DisWOT.

### 4.1. Experimental Setup

We perform experiments on four datasets, namely CIFAR-10, CIFAR-100, ImageNet-16-120, and ImageNet-1k. In the search process, we only use one batch of training data to get the statistic at nearly no cost. Following previous works, our experiments are conducted on the following three search spaces:

**Search Space  $S_0$ :** Following cifar-ResNet [23], the search space consists of three residual blocks and is based on CIFAR-100 datasets. The depth of each residual block is searched in set  $\{1, 3, 5, 7\}$ .

**Search Space  $S_1$ :** Following NAS-Bench-201 [16], this Darts-like search space is a cell-based search space consisting of stacked directed acyclic graphs.  $S_1$  is conducted on CIFAR-10, CIFAR-100, and ImageNet-16-120 datasets.

**Search Space  $S_2$ :** Following NDS [60], this search space consists of residual and bottleneck blocks defined in ResNet.  $S_2$  is based on CIFAR-100 and ImageNet-1k dataset.

### 4.2. Experiments on CIFAR-100

**Implementation Details.** We compare distillation gains with other zero-nas on search space  $S_1$ . In the search phase, we configure 48k evolution iters with 512 population sizes. In distillation, All searched student networks are trained via CRD’s settings [66] with ResNet56 as the teacher model.

**Distillation Results of Zero-cost Proxies.** We conduct detailed experiments on other zero-cost proxies with different knowledge distillation methods. Note that we search the student network under constraints of 1M parameters. The results in Table 5 demonstrated that our proposed DisWOT achieved superior results compared with other zero-cost proxies with different knowledge distillation methods. The DisWOT outperforms its counterparts vanilla networks by around 2%, while achieving consistent improvements among

Table 5. Distillation results (%) of different zero-cost proxies with knowledge distillation methods under 1M parameters.

Method	Random	FLOPs	Synflow	NWOT	DisWOT
Baseline	69.52	71.37	72.88	71.80	73.12
KD	70.45	72.13	73.72	72.57	74.73
FitNets	70.12	72.40	73.55	72.72	74.85
AT	70.16	72.97	73.52	72.08	74.50
SP	70.46	72.14	73.50	72.16	74.95
RKD	71.19	72.22	73.69	72.63	74.62
CRD	71.59	72.78	73.99	73.12	75.25

different distillation methods, such as KD [26], FitNets [62], AT [82], SP [69], RKD [50], and CRD [67].

Table 6. Distillation results(%) of zero-cost proxies under {0.5,1,2}M parameters.

Param.	FLOPs	NWOT	DisWOT	DisWOT <sup>†</sup>
0.5M	69.88	70.38	72.89	<b>73.75</b>
1M	72.13	72.57	74.23	<b>75.25</b>
2M	73.27	73.86	75.95	<b>76.67</b>

Table 7. Distillation results(%) of zero-cost proxies under {50,100}M FLOPs on space  $S_1$ .

FLOPs	NWOT	Synflow	DisWOT	FLOPs	NWOT	Synflow	DisWOT
50M	63.19	64.28	65.98	100M	70.38	72.12	72.89

Table 8. Ranking correlation (%) of zero-cost proxies on  $S_0$  space on CIFAR-100.

Method	Kendall’s Tau	Spearman	Pearson
FLOPs [1]	51.61	72.92	76.40
Fisher [1]	62.86	81.37	20.90
Grad_Norm [1]	63.75	82.35	39.35
SNIP [33]	67.22	85.07	51.09
NWOT [48]	31.87	45.66	48.99
DisWOT (ours)	<b>73.98</b>	<b>91.38</b>	<b>84.83</b>

**Analysis on Varying Parameter Constraints.** We analyze the performance of student models under different parameter constraints obtained by DisWOT on CIFAR-100. As shown in the Table 6, we compared our method with two zero proxies, a.k.a. FLOPs [1] and NWOT [48], under the parameter constraints of 0.5, 1, and 2M, respectively, and the results demonstrate that our method still achieves excellent results. As shown in Tab. 7, DisWOT also outperforms previous SOTA methods with 0.8%~1.7%<sup>†</sup> gains under same FLOPs constraints,

**Ranking Correlation with Distill Accuracy.** Based on search space  $S_0$ , we perform vanilla training and distillation for each candidate with CRD’s settings [66]. Then, we collect these vanilla results as GT and analyze the different zero-proxy’s correlation with them. As shown in Table 8, the results illustrate that our DisWOT achieves higher than Fisher, GradNorm, SNIP, FLOPs, and NWOT by a large margin, and achieve results that are on par with the best zero-cost proxy, a.k.a. Zen-NAS and Synflow, on Kendall’s Tau, Pearson, and Spearman coefficient.

### 4.3. Experiments on NAS-Bench-201

**Implementation Details.** For search trials, we first adopt ResNet110/56 as the teachers and then Conduct an evolution search with the DisWOT metric and get the best student network. We randomly sampled 50 candidate architectures to evaluate sequencing consistency. The distillation settings are the same as the Sec.4.2

**Comparison results** As shown in Table 9, some training-free methods can achieve good results with much faster speedups, such as NWOT and TE-NAS. Our proposed method DisWOT achieves a speedup ratio of 180 $\times$ , where if semantic similarity metric is removed, we can achieve a 300 $\times$  speedup ratio at the expense of some accuracy.

### 4.4. Experiments on ImageNet

**Implementation Details.** We searched the ResNet18 level network regarding the search space in NDS [60]. Specifically, we limit the number of parameters to less than 13M and the depth of the network to up to 20 layers and find the optimal network by evolution algorithm with the DisWOT metric. As shown in Table 10, guided by three different sizes of networks, we used DisWOT to find the optimal student network. We trained the student network obtained by the search using the distillation strategy in DisWOT. Implementation details are available in supplementary materials.

**Comparison Results.** Table 10 reports the performance of DisWOT on ImageNet with ResNet34/50 as teacher network. The results demonstrate that the student architecture of the ResNet18-level obtained by DisWOT under different teacher guidance and using different distillation strategies yielded significantly better results than its counterparts.

### 4.5. Ablation Studies of DisWOT

We perform ablation experiments to verify the validity of each component of DisWOT in search space  $S_0$ . As shown in Table 11, for semantic knowledge, similarity matrix obtains a more robust ranking improvement than simple FitNet [61]. For  $\mathcal{M}_r$ , similarity matrix performs better on relational knowledge than RKD [51]. DisWOT integrates semantic and relational knowledge to obtain an additional ranking improvement than stand-alone scores. The weight initialization scheme plays an important role in zero-proxy.

Table 9. Distillation results on CIFAR-10, CIFAR-100, and ImageNet-16 in NAS-Bench-201 [14]. Dis. Acc. (%) represents the accuracy of the searched architecture after distillation training. Time (s) denotes the time cost (GPU-seconds) during the search phase. The results of NWOT and TE-NAS come from their original papers. Our DisWOT achieves competitive results with the lowest costs.

Type	Model	CIFAR-10			CIFAR-100			ImageNet-16-120		
		Dis. Acc(%)	Time (s)	Speed-up	Dis.Acc(%)	Time (s)	Speed-up	Dis. Acc(%)	Time (s)	Speed-up
Multi-trial	RS	93.63	216K	1.0×	71.28	460K	1.0×	44.88	1M	1.0×
	RL [4]	92.83	216K	1.0×	71.71	460K	1.0×	44.35	1M	1.0×
	BOHB [19]	93.49	216K	1.0×	70.84	460K	1.0×	44.33	1M	1.0×
	RSPS [39]	91.67	10K	21.6×	57.99	46K	21.6×	36.87	104K	9.6×
Weight-sharing	GDAS [17]	93.39	22K	12.0×	70.70	39K	11.7×	42.35	130K	7.7×
	DARTS [43]	89.22	23K	9.4×	66.24	80K	5.8×	43.18	110K	9.1×
Training-free	NWOT [48]	93.73	2.2K	100×	73.31	4.6K	100×	45.43	10K	100×
	TE-NAS [10]	93.92	2.2K	100×	71.24	4.6K	100×	44.38	10K	100×
DisWOT	$\mathcal{M}_s$ & $\mathcal{M}_r$	93.55	1.2K	180×	74.21	9.2K	180×	47.30	20K	180×
	$\mathcal{M}_r$	93.49	0.72K	<b>300×</b>	73.62	18.4K	<b>300×</b>	45.63	40K	<b>300×</b>

Table 10. The accuracy (%) of ResNet18 on ImageNet-1k with various teachers. Results of other KD methods refer to the papers of CRD [67] and ESKD [11]. ATKD  $A_{R34}$  [49] denotes ResNet34 used as the assistant teacher. N/A means no available results. Our DisWOT obtains better performance than other methods and improves students’ performance positively correlated with that of the teacher.

Teacher	Student	Acc.	Teacher	Student	KD [26]	ESKD [11]	ATKD $A_{R18}$ [49]	ONE [32]	DML [83]	CRD [67]	DisWOT
ResNet34	ResNet18	Top-1	73.40	69.75	70.66	70.89	70.78	70.55	71.03	71.17	<b>72.08</b>
		Top-5	91.42	89.07	89.88	90.06	89.99	89.59	90.28	90.32	<b>90.38</b>
Teacher	Student	Acc.	Teacher	Student	KD [26]	ATKD $A_{R18}$ [49]	ATKD $A_{R34}$ [49]	Seq. ESKD [11]	ESKD [11]	SRRL [80]	DisWOT
ResNet50	ResNet18	Top-1	76.16	69.75	70.68	70.65	70.85	70.65	70.95	71.20	<b>72.30</b>
		Top-5	92.86	89.07	N/A	N/A	N/A	N/A	N/A	N/A	<b>90.51</b>

Table 11. Spearman correlation (“mean±std”) of DisWOT on search space  $S_0$ .

Knowledge	Metric	Spearman (%)
$\mathcal{M}_s$	FitNets [61]	64.06±6.11
$\mathcal{M}_s$	Similarity matrix	73.68±5.45
$\mathcal{M}_r$	RKD [66]	13.52±11.51
$\mathcal{M}_r$	Similarity matrix	72.36±3.42
$\mathcal{M}_s$ & $\mathcal{M}_r$	Similarity matrix	<b>77.51±2.76</b>

We verify the effect of the initialization strategy of the network on the ranking consistency. The results in Tab. 12 demonstrate that the Gaussian initialization strategy is detrimental to  $\mathcal{M}_s$ , but beneficial to  $\mathcal{M}_r$ .

## 5. Conclusion

In this paper, we present DisWOT, a new teacher-aware student architecture search without training framework for distillation. Based on key observations about the difference between vanilla and distillation accuracy, DisWOT measures the new zero-cost proxy conditioned on the similarity of feature semantics and sample relations between random-initialized teacher-student network. Then, DisWOT search for the best student architectures for the given teacher using an evolutionary algorithm with these metrics. Thorough

Table 12. “mean±std %” Spearman of proxies via Kaiming and Gaussian initialization on search space  $S_0$  and NAS-Bench-201 with various seeds.

Space	Initial	Fisher	GradNorm	NWOT	DisWOT
$S_0$	Kaim.	81.37±0.01	82.35±0.01	45.66±0.05	84.08±0.03
	Gauss.	80.99±0.01	75.50±0.01	45.36±0.03	91.38±0.03
NB-201	Kaim.	54.63±0.15	58.70±0.11	64.41±0.08	65.57±0.02
	Gauss.	45.91±0.09	45.70±0.11	62.24±0.07	72.36±0.02

evaluations are performed on diverse datasets and search spaces, and DisWOT achieves significant performance gains in various neural networks with at least 180× training acceleration. We experimentally and theoretically explained the relationship between similarity difference and distillation performance. In addition, we also extend DisWOT to new distillers and general zero proxy to predict the performance of models. By doing this, we bridge the higher-order knowledge between distillation and network architecture search. This approach represents an elegant and practical solution, which we hope will inspire future research on knowledge distillation and architecture search design.

**Limitations.** Following most zero-cost NAS, we evaluate DisWOT in classification tasks. In the future work, we will make efforts to expand the DisWOT for downstream tasks (e.g., object detection and semantic segmentation).



## A. More Comparisons and Discussions

In this section, we provide more analysis and discussion about DisWOT from different aspects.

### A.1. DisWOT under different teacher models.

When the teacher model becomes larger, the fixed hand-designed model would have huge teacher-student gaps, limiting the performance gain. DisWOT aims to solve this problem by searching the suitable student architecture for different teacher models. According to the results in Table 13, the accuracy of the student network ResNet20 [24] is unable to make consistent gains as the size of the teacher network increases. Our proposed DisWOT enables a consistent increase in student network performance as the teacher network capacity increases on search space  $S_0$ . In addition, DisWOT $\dagger$  achieves a performance gain of about 2% when stronger distillers are adopted.

Table 13. Top-1 accuracy (%) of ResNet20 with KD [27], student (DisWOT) with KD [27], student (DisWOT) with DisWOT $\dagger$  on search space  $S_0$  under different teachers.

Teacher	ResNet20	DisWOT	DisWOT $\dagger$
ResNet32	70.24	71.01	71.85
ResNet44	70.56	71.25	72.12
ResNet56	70.98	71.63	72.56
ResNet110	70.79	71.84	72.92

### A.2. Ranking correlation metrics

We denote the ground-truth (GT) performance and approximated scores of architectures  $\alpha_i (i = 1, \dots, N)$  as  $\beta_i (i = 1, \dots, N)$  and  $\gamma_i (i = 1, \dots, N)$ , respectively, and the ranking of the GT and estimated score  $\beta_i, \gamma_i$  as  $r_i, k_i \in \{1, \dots, N\}$ . Three correlation criteria is adopted in this paper. Pearson coefficient ( $r$ ), Kendall’s Tau ( $\tau$ ), Spearman coefficient ( $\rho$ ).

- Pearson correlation coefficient (Linear Correlation):  

$$r = \text{corr}(\beta, \gamma) / \sqrt{\text{corr}(\beta, \beta) \text{corr}(\gamma, \gamma)}.$$
- Kendall’s Tau correlation coefficient: The relative difference of concordant pairs and discordant pairs  

$$\tau = \sum_{i < j} \text{sgn}(\beta_i - \beta_j) \text{sgn}(\gamma_i - \gamma_j) / \binom{M}{2}.$$
- Spearman correlation coefficient: The Pearson correlation coefficient between the ranking variables  $\rho = \text{corr}(r, k) / \sqrt{\text{corr}(r, r) \text{corr}(k, k)}.$

Pearson measures the linear relationship between two variables, while Kendall’s Tau and Spearman measure the monotonic relationship. They return a value between -1 and 1, with -1 indicating an inverse correlation, 1 indicating a positive correlation, and 0 representing no relationship.

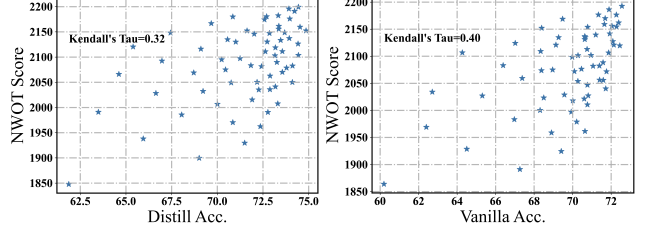


Figure 5. Left: Correlation of distill accuracy & NWOT score on search space  $S_0$ . Right: Correlation of vanilla accuracy & NWOT score on search space  $S_0$ .

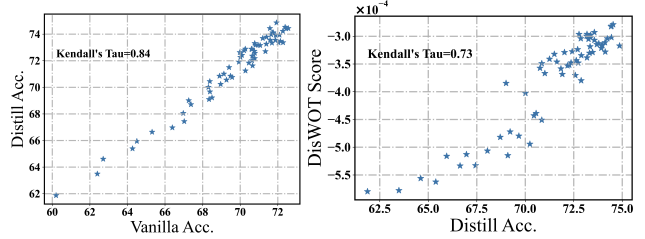


Figure 6. Left: Correlation of vanilla accuracy & distill accuracy on search space  $S_0$ . Right: Correlation of distilling accuracy & DisWOT score on search space  $S_0$ .

In search space  $S_0$ , we evaluate and verify the ranking consistency for all the architectures in the search space. In search space  $S_1$ , we randomly sampled 50 sub-networks from the search space to calculate the ranking consistency results due to the excessive time overhead of the full measurement and repeated each experiment 10 times.

### A.3. Detailed analysis about vanilla-distilling training disparity

In this paper, we notice an interesting and non-trivial observation: the discrepancy between the model’s performance under vanilla training and under distillation training. We present more analysis in detail here from three aspects.

**Ranking correlation degrade.** In Table 14, we tabulate the ranking correlations for different zero-proxies with distilling and vanilla training on the search space  $S_0$ . The existing zero-proxies’ distillation correlations are reduced by 5%  $\sim$  69% than the vanilla correlation. This common issue shows that existing zero-shot NAS methods score sub-optimal student architectures for a given teacher model. DisWOT not only has a better correlation for distillation but also is free of vanilla-distill ranking consistency degradation.

**Correlation visualization of different scores.** Figure 5 demonstrates the ranking consistency of NWOT [48] for distillation accuracy and vanilla accuracy, and there is a large discrepancy between them, with an 8% difference in Kendall’s Tau. Figure 6 (left) demonstrates the ranking correlation between vanilla accuracy and distillation accuracy.

Table 14. Details experiments of the discrepancy between vanilla accuracy and distillation accuracy on search space  $S_0$ .

Method	Ranking with distill accuracy			Ranking with vanilla accuracy			Ranking gap for distill and vanilla accuracy		
	Kendall's Tau	Spearman	Pearson	Kendall's Tau	Spearman	Pearson	Kendall's Tau	Spearman	Pearson
FLOPs [1]	51.61	72.92	76.40	58.74	79.47	79.19	7.13 (↓)	6.55 (↓)	2.79 (↓)
Fisher [1]	62.86	81.37	20.90	81.68	95.28	70.24	18.82 (↓)	13.91 (↓)	49.34 (↓)
Grad_norm [1]	63.75	82.35	39.35	84.76	96.55	76.07	21.01 (↓)	14.2 (↓)	36.72 (↓)
NWOT [48]	31.87	45.66	48.99	40.29	56.46	56.23	8.42 (↓)	10.80 (↓)	7.24 (↓)
Plain [1]	10.72	13.57	-0.91	54.98	77.12	67.55	44.26 (↓)	63.55 (↓)	68.46 (↓)
SNIP [33]	67.22	85.07	51.09	84.66	96.38	77.83	17.44 (↓)	11.31 (↓)	26.74 (↓)
<b>DisWOT</b>	73.98	91.38	84.83	73.02	91.26	82.98	0.96 (↑)	0.12 (↑)	1.85 (↑)

Surprisingly, their correlation is only 0.84, which indicates that there is a non-negligible gap between distillation results and vanilla results. On the one hand, it indicates that a new zero-cost metric needs to be designed to improve the ranking consistency for the distillation. On the other hand, vanilla accuracy can be used as a rough measure of ranking consistency when distillation accuracy is unavailable. Figure 6 (right) illustrates that our proposed DisWOT achieve better ranking consistency of distillation accuracy on search space  $S_0$ .

**Analysis of detailed examples.** Table 15 summarizes 4 groups of specific student pairs with vanilla-distill gaps. For groups A and B, despite student models A2, B2 having more parameters and better vanilla accuracy, their distillation accuracy is inferior to A1, B1. This indicates an important effect on the overall depth of students for distillation. For groups C and D, the results show that the model with more blocks in stage-2 enjoys better distillation accuracy. In addition, DisWOT predicts the correct scores for these models.

Table 15. Parameters (K), vanilla accuracy (%), distillation accuracy (%), and prediction scores ( $10^{-4}$ ) of DisWOT for the student on search space  $S_0$ .

Group	Student	Param.	Vanilla Acc.	Distill Acc.	DisWOT
A1	ResNet[7,1,3]	259.89	69.13	71.01	4.41
A2	ResNet[3,3,3]	278.32	69.57	70.76	3.34
B1	ResNet[7,5,3]	334.13	70.76	72.58	5.44
B2	ResNet[1,7,3]	343.22	70.77	72.18	5.20
C1	ResNet[5,5,7]	620.72	71.93	74.86	7.37
C2	ResNet[3,7,7]	648.50	72.45	74.42	7.33
D1	ResNet[7,3,5]	444.98	72.04	73.36	4.85
D2	ResNet[5,5,5]	472.76	72.09	73.94	8.17

#### A.4. Semantic properties of random networks

DisWOT leverages the semantic similarity metric of a randomly initialized teacher-student model to predict distillation performance, abandoning the training-based NAS paradigm. Models with various architectures have different semantic features because of their different effective receptive fields. To represent semantic and localization information, Gradient-weighted Class Activation Mapping (Grad-CAM [63]) meth-

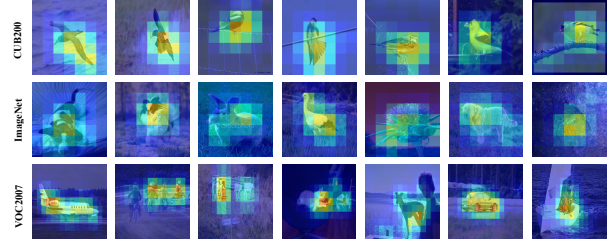


Figure 7. Localization results [7] of a randomly initialized network on ImageNet, VOC2007, and CUB-200. The network can localize the objects in an image, with a small standard deviation between different trials. Note that this figure is from Tobias [7].

ods have been widely adopted in weakly supervised object localization and model interpretability [3]. Recently, several studies [3, 75] reveal that randomly initialized models also have favorable semantic localization capabilities. As shown in the visual localization heatmaps of Figure 7, we can intuitively observe that randomly initialized networks can locate a single object without any training. The visualization results demonstrate that semantic information exists even in random networks, which can localize the objects in an image. In addition, we evaluate different strategies (e.g. CAM, Grad-CAM, and SCDA [75]) for semantic localization maps and find that Grad-CAM achieves stable prediction performance for the semantic similarity metric. Thus, we adopt the Grad-CAM of deeper layers to capture the informative relation similarity in this paper.

#### A.5. Comparisons with different NAS approaches.

We presents the results of more NAS methods in this section and find that our proposed DisWOT is better able to distinguish good architectures than both One-shot NAS and Zero-shot NAS. As shown in Table 19 and 16, we compare DisWOT with one-shot NAS (e.g. ENAS [57], SETN [15], SPOS [21]) and zero-shot NAS(a.k.a. Zen-NAS [41]). Results demonstrate that DisWOT achieves better performance than its counterparts in distillation and classification results. In addition, we conduct more correlation evaluation in Tab. 17. The results show that DisWOT

surpasses other proxies in the NAS-Bench-101/101-KD/201-KD and DisWOT ( $\mathcal{M}_s$ ) achieves superior correlations than DisWOT ( $\mathcal{M}_r$ ), which are consistent with findings in space  $S_0$ .

Table 16. Top-1 accuracy (%) of different NAS algorithms under distill training on NAS-Bench-201 [14].

Datasets	ENAS [57]	SETN [15]	SPOS [21]	Zen-NAS [41]	DisWOT
CIFAR-10	34.94	81.61	92.98	89.60	<b>93.55</b>
CIFAR-100	11.14	59.78	72.91	71.86	<b>74.21</b>
ImageNet16-120	11.53	29.91	43.50	39.44	<b>47.30</b>

## A.6. About KD-based zero-cost proxies.

In this section, we present DisWOT as a new universal zero proxy and propose a series of KD-based zero proxies based on this motivation. As shown in Table 18, we further provide the ranking correlation of various KD-based zero-cost proxies on three datasets. We adopt the optimal architecture in the search space of NAS-Bench-201 as a teacher network and conduct 10 independent experiments of three knowledge distillation methods (a.k.a., CC [54], KD [27], and NST [30]). We observe that DisWOT achieves an acceptable ranking correlation on three datasets. NST [30] show impressive ranking ability, which is the best KD-based zero-proxies in DisWOT framework. DisWOT reveals better performance than most of zero-cost proxies under vanilla training, as shown in Table 19.

## B. Details of Search Space and Settings

In this section, we introduce the implementation details in the three search spaces and the detailed training settings.

### B.1. $S_0$ search space

**Search space.** As illustrated in Figure 20, we construct the search space  $S_0$  based on ResNet20, a simple resnet designed for CIFAR-10/100, where each building block consists of two  $3 \times 3$  convolutional layers and the depth of each residual block is searched in set  $\{1, 3, 5, 7\}$ . The search space size is  $4^3 = 64$  in total.

**Implementation details.** As for results of vanilla classification results, we train each architecture in the search space  $S_0$  with the same strategy. For each architecture in the search space  $S_0$ , we adopt ResNet110 as a teacher network. Specifically, we train each architecture via momentum SGD, using cross-entropy loss for 240 epochs. We set the weight decay as  $5e-4$  and adopted a multi-stage scheduler to decay the learning rate from 0.1 to 0. We use the random flip with the probability of 0.5, the random crop  $32 \times 32$  patch with 4 pixels paddings, and the normalization over RGB channels. All of the experiments are based on CIFAR-100 datasets. As for the distillation results, the vanilla knowledge distillation

Table 17. “mean $\pm$ std %” Spearman correlation on NAS-Bench-101 and NAS-Bench-201. NAS-Bench-101/201-KD denotes to distill accuracies of the architectures on NAS-Bench-101/201.

Method	NAS-Bench-101	NAS-Bench-201	NAS101-KD	NAS201-KD
FLOPs	30.81% $\pm$ 0.00	63.38% $\pm$ 0.06	15.56% $\pm$ 0.04	64.55% $\pm$ 0.01
Fisher	-38.81% $\pm$ 0.14	35.91% $\pm$ 0.09	-33.92% $\pm$ 0.14	4.45% $\pm$ 0.08
Grad_Norm	-39.23% $\pm$ 0.08	58.70% $\pm$ 0.11	-39.16% $\pm$ 0.01	-10.01% $\pm$ 0.11
SNIP	-29.01% $\pm$ 0.09	58.17% $\pm$ 0.15	-21.78% $\pm$ 0.02	16.91% $\pm$ 0.10
Synflow	43.69% $\pm$ 0.12	<b>74.61%<math>\pm</math>0.08</b>	20.36% $\pm$ 0.08	74.63% $\pm$ 0.02
NWOT	32.84% $\pm$ 0.51	64.41% $\pm$ 0.08	22.97% $\pm$ 0.04	35.27% $\pm$ 0.03
DisWOT ( $\mathcal{M}_s$ )	<b>49.61%<math>\pm</math>0.05</b>	65.74% $\pm$ 0.07	50.16% $\pm$ 0.09	53.88% $\pm$ 0.06
DisWOT ( $\mathcal{M}_r$ )	30.74% $\pm$ 0.06	56.46% $\pm$ 0.08	42.94% $\pm$ 0.11	45.27% $\pm$ 0.07

Table 18. Ranking correlation of our KD-based zero-cost proxies on NAS-Bench-201.

Datasets	Method	Kendall’s Tau	Spearman	Pearson
CIFAR-10	CC [54]	0.48	0.68	0.56
	KD [27]	0.35	0.50	0.40
	NST [30]	0.64	0.83	0.72
	DisWOT	0.41	0.61	0.54
CIFAR-100	CC [54]	0.43	0.65	0.58
	KD [27]	0.38	0.54	0.55
	NST [30]	0.57	0.72	0.64
	DisWOT	0.56	0.72	0.65
ImageNet16	CC [54]	0.53	0.71	0.66
	KD [27]	0.44	0.61	0.65
	NST [30]	0.54	0.74	0.74
	DisWOT	0.49	0.69	0.55

Table 19. Top-1 accuracy (%) of different NAS algorithms under vanilla training on NAS-Bench-201 [14].

Datasets	ENAS [57]	SETN [15]	SPOS [21]	Zen-NAS [41]	DisWOT
CIFAR-10	53.89	87.64	93.23	90.70	<b>93.37</b>
CIFAR-100	13.96	59.05	71.03	68.26	<b>71.53</b>
ImageNet16-120	14.84	32.52	42.19	40.60	<b>45.50</b>

Table 20. Supernet architecture of the  $S_0$  search space. Each line describes a sequence of 1 or more identical layers, repeated *repeat* times. All layers in the same sequence have the same number of output channels.

input	block	channels	repeat	stride
$32^2 \times 3$	$3 \times 3$ conv	16	1	2
$32^2 \times 16$	Res Block	16	[1,3,5,7]	2
$16^2 \times 16$	Res Block	32	[1,3,5,7]	2
$8^2 \times 32$	Res Block	64	[1,3,5,7]	2
$8^2 \times 64$	Global Avgpool	-	1	-
64	FC	100	1	-

methods [27] are adopted. Specifically, we conduct experiments based on the CRD [67]. For KD [27], we follow the Equation 14 and set  $\alpha = 0.9$  and  $\rho = 4$ .

$$\mathcal{L}_{KL} = \alpha \rho^2 CE(\sigma(z^T/\rho), \sigma(z^S/\rho)) \quad (14)$$

where  $z^T$  and  $z^S$  denote the logits of teacher and student, respectively.  $\rho$  is the temperature,  $\alpha$  is a balancing weight, and  $\sigma$  is a softmax function. CE denotes the cross entropy loss.

## B.2. $S_1$ search space

**Search space.** The search space  $S_1$  is following the cell-based search space NAS-Bench-201 [14], where a cell is represented as a directed acyclic graph (DAG). Each edge in search space  $S_1$  is associated with an operation selected from a predefined operation set, which consists of (1) zero, (2) skip connection, (3)  $1 \times 1$  convolution, (4)  $3 \times 3$  convolution, and (5)  $3 \times 3$  average pooling layer. The DAG has 4 nodes, each representing the sum of all features from previous nodes. The search space size of  $S_1$  is 15,625 in total.

**Implementation Details.** We randomly sampled 50 candidate architectures to evaluate ranking consistency. All experiments are implemented on a single NVIDIA 3090Ti GPU, with the baseline from the AutoDL [14]. We recommend using a network with higher complexity or better performance in the search space as the teacher network. The process of DisWOT is divided into three steps: (1) Determine a specific teacher network (deeper or more complex). (2) Perform an evolutionary search with DisWOT metrics to obtain the best student network. (3) Distill the student network with vanilla KD [26] based on a specific teacher network. The distillation setting is the same as Section B.1.

**Searched architectures.** For other NAS methods, the searched architectures (see Table 21) of RS [14], ENAS [57], SETN [15] and SPOS [21] are borrowed from the official implementation [14], and the remaining zero-shot NAS methods utilize the same evolutionary search algorithm. Expressly, we set the initial population size as 20, and the sample size as 10. The total evolution search cycle is set as 5,000. We calculate the zero-proxy score with only one batch of data as fitness during evolution. The teacher network used in DisWOT is the best architecture in the search space.

## B.3. $S_2$ search space

**Search space.** Following NDS [59], we design the search space for CIFAR and ImageNet, respectively. The search space designed for CIFAR consists of a stem, followed by 6 stages, and a head, as shown in Table 22. The  $i$ -th stage consists of  $d_i$  blocks with  $c_i$  channels and stride of  $s_i \in \{1, 2\}$ . The number of channels  $c_i$  needs to be divisible by 8, and the minimal number of channels should be larger than 8. The candidate blocks can be residual blocks or bottleneck blocks defined in ResNet, and the kernel size can be chosen from set  $\{3, 5, 7\}$ . As shown in Table 23, the search space designed for ImageNet consists of 4 stages, following the configuration of ResNet18.

Table 22. Design space parameterization of  $S_2$  for CIFAR-10/100. "POOL" denotes the global average pooling, and "FC" denotes a fully connected network.

stage	block	channels	repeat	stride
steam	$3 \times 3$ conv	$c_0$	1	1
stage1	{block}	$c_1$	$d_1$	$s_1$
stage2	{block}	$c_2$	$d_2$	$s_2$
stage3	{block}	$c_3$	$d_3$	$s_3$
stage4	{block}	$c_4$	$d_4$	$s_4$
stage5	{block}	$c_5$	$d_5$	$s_5$
stage6	{block}	$c_6$	$d_6$	$s_6$
head	POOL + FC	10	-	-

Table 23. Design space parameterization of  $S_2$  for ImageNet. "POOL" denotes the global average pooling, and "FC" denotes fully connected network.

stage	block	channels	repeat	stride
steam	$7 \times 7$ conv	$c_0$	1	2
stage1	{block}	$c_1$	$d_1$	$s_1$
stage2	{block}	$c_2$	$d_1$	$s_2$
stage3	{block}	$c_3$	$d_1$	$s_3$
stage4	{block}	$c_4$	$d_1$	$s_4$
head	POOL + FC	1000	-	-

**Training settings.** We search the ResNet18 level network regarding the search space in NDS [59]. Specifically, we limit the number of parameters to less than 13M and the depth of the network to up to 20 layers and find the optimal network by evolution algorithm with the DisWOT metric. Please refer to Section C.1 for more details about the evolutionary search for the search space  $S_2$ . Specifically, we adopt ResNet34 as a teacher network and conduct a vanilla knowledge distillation process [26] with  $\rho = 1$  and  $\alpha = 3$  as shown in Equation 14. For ImageNet, we follow the standard PyTorch practice, and the batch size is 256.

**Searched architectures.** After the evolutionary search, we presented the optimal student network obtained for CIFAR and ImageNet, as shown in the Table 24 and Table 25, respectively. We observe that the searched architecture of DisWOT has very different characteristics from the artificially designed student architecture, i.e., DisWOT prefers student networks with larger convolutional kernels in shallow layers. Generally speaking, teacher networks tend to have deeper layers and thus have a larger receptive field. Guided by the teacher network, DisWOT favors larger convolutional kernels in the shallow layers so that the receptive field of the student network is as close to that of the teacher as possible. However, the network searched on ImageNet only changed the number of channels under the parameter restriction. We infer that expanding the kernel size leads to a



Table 21. Searched architectures of NAS algorithms. The searched results are denoted by a string from NAS-Bench-201 API [14].

	Searched architectures
RS [14]	lskip_connect~0l+lnor_conv_3x3~0lskip_connect~1l+lnor_conv_3x3~0lnor_conv_1x1~1lavg_pool_3x3~2l
ENAS [57]	lskip_connect~0l+avg_pool_3x3~0lskip_connect~1l+avg_pool_3x3~0lskip_connect~1lskip_connect~2l
SETN [15]	lnor_conv_3x3~0l+lskip_connect~0lskip_connect~1l+lskip_connect~0lskip_connect~1lavg_pool_3x3~2l
SPOS [21]	lskip_connect~0l+lnor_conv_1x1~0lnor_conv_3x3~1l+lnor_conv_1x1~0lavg_pool_3x3~1lnor_conv_3x3~2l
Zen-NAS [41]	lskip_connect~0l+lnor_conv_3x3~0lnor_conv_3x3~1l+lskip_connect~0lskip_connect~1lnor_conv_3x3~2l
NWOT [48]	lnor_conv_1x1~0l+lnor_conv_3x3~0lnor_conv_1x1~1l+lnor_conv_1x1~0lnor_conv_3x3~1lnor_conv_1x1~2l
DisWOT(ours)	lskip_connect~0l+lnor_conv_3x3~0lnor_conv_1x1~1l+lnor_conv_1x1~0lnor_conv_3x3~1lnor_conv_3x3~2l

massive amount of additional parameters, which will lead to exceeding the budget. We infer that the network will prefer a larger kernel if a sufficient budget is available.

Table 24. Search results of the ResNet-like search space for CIFAR-10/100. "Basic" denotes the basic block proposed in ResNet [24].

input	block	channels	repeat	stride
$32^2 \times 3$	$3 \times 3$ conv	88	1	1
$32^2 \times 88$	Basic $7 \times 7$	96	3	1
$32^2 \times 120$	Basic $5 \times 5$	192	2	2
$16^2 \times 192$	Basic $5 \times 5$	176	2	1
$16^2 \times 96$	Basic $5 \times 5$	168	3	2
$8^2 \times 168$	Basic $3 \times 3$	112	3	2
$4^2 \times 112$	Basic $3 \times 3$	512	1	1
512	POOL + FC	1000	1	-

Table 25. Search results of the ResNet-like search space for ImageNet under 13M parameter limit. "Basic" denotes the basic block proposed in ResNet [24].

input	block	channels	repeat	stride
$224^2 \times 3$	$7 \times 7$ conv	96	1	2
$112^2 \times 96$	Basic $3 \times 3$	64	3	2
$56^2 \times 64$	Basic $3 \times 3$	128	2	2
$28^2 \times 128$	Basic $3 \times 3$	256	2	2
$14^2 \times 256$	Basic $3 \times 3$	512	2	2
512	POOL + FC	1000	1	-

## C. Details of Algorithm for DisWOT

In this section, we describe the implementation details of the evolutionary algorithm and the implementation code of DisWOT.

### C.1. Implementation of Evolutionary Algorithm

Here we adopt Evolutionary Algorithm (EA) as an architecture generator to find optimal student network. In this section, we further describe the mutation process of the evolutionary algorithm in details. In the evolutionary algorithm, we randomly generate  $P$  architectures with constraints  $C$  and then select the  $top - k$  architectures by DisWOT metric

Table 26. DisWOT-11.7M based on ImageNet (Searched for segmentation task.)

block	kernel	in	out	stride	bottleneck	# layers
Conv	7	3	64	2	-	1
Res	3	64	64	2	64	2
Res	3	64	128	2	128	2
Res	3	128	256	2	256	2
Res	3	256	512	2	512	2
Conv	1	512	2384	1	-	1

from the population. Then we randomly select the parent architecture from the  $top - k$  architectures and mutate it. The mutation algorithm is presented in Algorithm 2. Specifically, for  $S_2$  search space, we provide basic blocks with a kernel size of  $\{3, 5, 7\}$  and choose  $\{0.67, 0.8, 1.25, 1.5\}$  as the mutation range for channels. The number of channels should be divisible by 8, and the max number of channels is 2048. The depth of chosen block can be mutated in range  $\{+1, -1\}$ . After mutating the architecture, we check whether it is valid, e.g., the parameter is meet the predefined constraint.

#### Algorithm 2 Mutation Algorithm for DisWOT

**Input:** Parent architecture  $A_i$ , Search space  $S$ .

**Output:** Mutated architecture  $\hat{P}_i$

- 1: Randomly select a block  $a_i$  from Parent architecture  $A_i$ ;
- 2: Randomly mutate the kernel size of  $a_i$  from  $S(\{block\})$ ;
- 3: Randomly mutate the width of  $a_i$  from  $S(c_i)$ ;
- 4: Randomly mutate the depth of  $a_i$  from  $S(d_i)$ ;
- 5: Check whether the mutated architecture is valid;
- 6: Return the mutated architecture  $\hat{P}_i$ ;

### C.2. Implementation of metric in DisWOT

The section presents the implementation of semantic similarity metric and relation similarity metric in DisWOT. The semantic similarity metric measures the inter-correlation on the accumulated Grad-CAM for teacher and student net-

works, whose calculation needs one forward and one backward to get the localization information. The relation similarity metric measures the relationship between input samples whose activations of teacher and student networks are needed.

**Implementation of semantic similarity metric.** Here, we present the implementation code of the relation similarity metric, as shown in List C.2. We need the Grad-CAM maps of all classes for calculation, which needs at least one forward and one backward to get the Grad-CAM similarity. Different from ICKD [45], there are mainly two differences: (1) The Gaussian initialized teacher network and student network is backpropagated only once. (2) We only use the grad of the fully-connected layer for calculation.

Listing 1. The PyTorch implementation of semantic similarity metric.

---

```
import torch
import torch.nn as nn
import torch.nn.functional as F

def semantic_similarity_metric(teacher,
    student, batch_data):
    criterion = nn.CrossEntropyLoss()
    image, label = batch_data
    # Forward once.
    t_logits = teacher.forward(image)
    s_logits = student.forward(image)
    # Backward once.
    criterion(t_logits, label).backward()
    criterion(s_logits, label).backward()
    # Grad-cam of fc layer.
    t_grad_cam = teacher.fc.weight.grad
    s_grad_cam = student.fc.weight.grad
    # Compute channel-wise similarity
    return -1 *
        channel_similarity(t_grad_cam,
            s_grad_cam)

def channel_similarity(f_t, f_s):
    bsz, ch = f_s.shape[0], f_s.shape[1]
    # Reshape
    f_s = f_s.view(bsz, ch, -1)
    f_t = f_t.view(bsz, ch, -1)
    # Get channel-wise similarity matrix
    emd_s = torch.bmm(f_s, f_s.permute(0,
        2, 1))
    emd_s = F.normalize(emd_s, dim=2)
    emd_t = torch.bmm(f_t, f_t.permute(0,
        2, 1))
    emd_t = F.normalize(emd_t, dim=2)
    # Produce L2 distance
    G_diff = emd_s - emd_t
    return (G_diff * G_diff).view(bsz,
        -1).sum() / (ch * bsz)
```

---

**Implementation of relation similarity metric.** Here we present the implementation code of relation similarity, as shown in List C.2. With only the logits of the teacher and student network, our relation similarity metric is easy to implement. There are mainly two differences compared with SP [69]: (1) The teacher and student networks are initialized with kaiming initialization [22], which means the teacher network did not undergo any backpropagation. (2) Here, we only used the activation before global average pooling as input, and the activations of shallow layers are not utilized. In fact, we find that activations closer to the output are more informative. When using activations from shallow layers, experiments demonstrate that the relation similarity ranks poorly.

Listing 2. The PyTorch implementation of relation similarity metric.

---

```
import torch
import torch.nn as nn
import torch.nn.functional as F

def relation_similarity_metric(teacher,
    student, batch_data):
    image, label = batch_data
    # Forward pass
    t_feats =
        teacher.forward_features(image)
    s_feats =
        student.forward_features(image)
    # Get activation before average pooling
    t_feat = t_feats[-2]
    s_feat = s_feats[-2]
    # Compute batch similarity
    return -1 * batch_similarity(t_feat,
        s_feat)

def batch_similarity(f_t, f_s):
    # Reshape
    f_s = f_s.view(f_s.shape[0], -1)
    f_t = f_t.view(f_t.shape[0], -1)
    # Get batch-wise similarity matrix
    G_s = torch.mm(f_s, torch.t(f_s))
    G_s = F.normalize(G_s)
    G_t = torch.mm(f_t, torch.t(f_t))
    G_t = F.normalize(G_t)
    # Produce L2 distance
    G_diff = G_t - G_s
    return (G_diff * G_diff).view(-1,
        1).sum() / (bsz * bsz)
```

---

## References

- [1] Mohamed S Abdelfattah, Abhinav Mehrotra, Łukasz Dudziak, and Nicholas Donald Lane. Zero-cost proxies for lightweight nas. In *ICLR*, 2020. 3, 4, 6, 7, 10

- [2] Sungsoo Ahn, Shell Xu Hu, Andreas Damianou, Neil D Lawrence, and Zhenwen Dai. Variational information distillation for knowledge transfer. In *CVPR*, 2019. 1
- [3] Wonho Bae, Junhyug Noh, and Gunhee Kim. Rethinking class activation mapping for weakly supervised object localization. In *ECCV*, 2020. 10
- [4] Bowen Baker, Otkrist Gupta, Nikhil Naik, and Ramesh Raskar. Designing neural network architectures using reinforcement learning. In *ICLR*, 2017. 8
- [5] Lucas Beyer, Xiaohua Zhai, Amélie Royer, Larisa Markeeva, Rohan Anil, and Alexander Kolesnikov. Knowledge distillation: A good teacher is patient and consistent. In *CVPR*, 2022. 1
- [6] Tom B. Brown, Benjamin Mann, Nick Ryder and Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. *arXiv preprint*, arXiv:2005.14165, 2020. 1
- [7] Yun-Hao Cao and Jianxin Wu. A random cnn sees objects: One inductive bias of cnn and its applications. In *AAAI*, 2022. 2, 10
- [8] Defang Chen, Jian-Ping Mei, Hailin Zhang, Can Wang, Yan Feng, and Chun Chen. Knowledge distillation with the reused teacher classifier. In *CVPR*, 2022. 1
- [9] Kunlong Chen, Liu Yang, Yitian Chen, Kunjin Chen, Yidan Xu, and Lujun Li. Gp-nas-ensemble: a model for the nas performance prediction. In *CVPRW*, 2022. 3
- [10] Wuyang Chen, Xinyu Gong, and Zhangyang Wang. Neural architecture search on imagenet in four gpu hours: A theoretically inspired perspective. In *ICLR*, 2020. 8
- [11] Jang Hyun Cho and Bharath Hariharan. On the efficacy of knowledge distillation. In *ICCV*, 2019. 1, 3, 8
- [12] Peijie Dong, Xin Niu, Lujun Li, Zhiliang Tian, Xiaodong Wang, Zimian Wei, Hengyue Pan, and Dongsheng Li. Rd-nas: Enhancing one-shot supernet ranking ability via ranking distillation from zero-cost proxies. *arXiv preprint* arXiv:2301.09850, 2023. 3
- [13] Peijie Dong, Xin Niu, Lujun Li, Linzhen Xie, Wenbin Zou, Tian Ye, Zimian Wei, and Hengyue Pan. Prior-guided one-shot neural architecture search. *arXiv preprint* arXiv:2206.13329, 2022. 3
- [14] Xuanyi Dong and Yi Yang. Nas-bench-201: Extending the scope of reproducible neural architecture search. In *ICLR*, 2019. 8, 11, 12, 13
- [15] Xuanyi Dong and Yezhou Yang. One-shot neural architecture search via self-evaluated template network. *2019 ICCV*, 2019. 10, 11, 12, 13
- [16] Xuanyi Dong and Yi Yang. Searching for a robust neural architecture in four gpu hours. In *CVPR*, 2019. 2, 6
- [17] Xuanyi Dong and Yezhou Yang. Searching for a robust neural architecture in four gpu hours. *CVPR*, 2019. 8
- [18] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint* arXiv:2010.11929, 2020. 1
- [19] Stefan Falkner, Aaron Klein, and Frank Hutter. Bohb: Robust and efficient hyperparameter optimization at scale. In *ICML*, 2018. 8
- [20] Jindong Gu and Volker Tresp. Search for better students to learn distilled knowledge. In *ECAI*, 2020. 3
- [21] Zichao Guo, Xiangyu Zhang, Haoyuan Mu, Wen Heng, Zechun Liu, Yichen Wei, and Jian Sun. Single path one-shot neural architecture search with uniform sampling. *arXiv preprint* arXiv:1904.00420, 2019. 3, 10, 11, 12, 13
- [22] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. *2015 IEEE International Conference on Computer Vision (ICCV)*, 2015. 14
- [23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 6
- [24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 9, 13
- [25] Byeongho Heo, Minsik Lee, Sangdoo Yun, and Jin Young Choi. Knowledge transfer via distillation of activation boundaries formed by hidden neurons. In *AAAI*, 2019. 1
- [26] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint* arXiv:1503.02531, 2015. 1, 5, 6, 7, 8, 12
- [27] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. In *arXiv:1503.02531*, 2015. 9, 11
- [28] Yiming Hu, Xingang Wang, Lujun Li, and Qingyi Gu. Improving one-shot nas with shrinking-and-expanding supernet. *Pattern Recognition*, 2021. 3
- [29] Zehao Huang and Naiyan Wang. Like what you like: Knowledge distill via neuron selectivity transfer. *arXiv preprint* arXiv:1707.01219, 2017. 1, 3, 6
- [30] Zehao Huang and Naiyan Wang. Like what you like: Knowledge distill via neuron selectivity transfer. *arXiv:1707.01219*, 2017. 6, 11
- [31] Jangho Kim, SeoungUk Park, and Nojun Kwak. Paraphrasing complex network: Network compression via factor transfer. In *NeurIPS*, 2018. 1
- [32] Xu Lan, Xiatian Zhu, and Shaogang Gong. Knowledge distillation by on-the-fly native ensemble. In *NeurIPS*, 2018. 8
- [33] Namhoon Lee, Thalaiyasingam Ajanthan, and Philip Torr. Snip: Single-shot network pruning based on connection sensitivity. In *ICLR*, 2018. 3, 4, 6, 7, 10
- [34] Changlin Li, Jiefeng Peng, Liuchun Yuan, Guangrun Wang, Xiaodan Liang, Liang Lin, and Xiaojun Chang. Block-wisely supervised neural architecture search with knowledge distillation. *CVPR*, 2020. 3

- [35] Lujun Li. Self-regulated feature learning via teacher-free feature distillation. In *ECCV*, 2022. 3
- [36] Lujun Li and Zhe Jin. Shadow knowledge distillation: Bridging offline and online knowledge transfer. In *NeurIPS*, 2022. 3
- [37] Lujun Li, Liang Shiuan-Ni, Ya Yang, and Zhe Jin. Boosting online feature transfer via separable feature fusion. In *IJCNN*, 2022. 3
- [38] Lujun Li, Liang Shiuan-Ni, Ya Yang, and Zhe Jin. Teacher-free distillation via regularizing intermediate representation. In *IJCNN*, 2022. 3
- [39] Liam Li and Ameet S. Talwalkar. Random search and reproducibility for neural architecture search. *ArXiv*, 2019. 8
- [40] Lujun Li, Yikai Wang, Anbang Yao, Yi Qian, Xiao Zhou, and Ke He. Explicit connection distillation. In *ICLR*, 2020. 3
- [41] Ming Lin, Pichao Wang, Zhenhong Sun, Hesen Chen, Xiuyu Sun, Qi Qian, Hao Li, and Rong Jin. Zen-nas: A zero-shot nas for high-performance image recognition. 2021. 3, 4, 6, 10, 11, 13
- [42] Sihao Lin, Hongwei Xie, Bing Wang, Kaicheng Yu, Xiaojun Chang, Xiaodan Liang, and Gang Wang. Knowledge distillation via the target-aware transformer. In *CVPR*, 2022. 2
- [43] Hanxiao Liu, Karen Simonyan, and Yiming Yang. DARTS: differentiable architecture search. In *ICLR*. 8
- [44] Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. *arXiv preprint arXiv:1806.09055*, 2018. 3
- [45] Li Liu, Qinwen Huang, Sihao Lin, Hongwei Xie, Bing Wang, Xiaojun Chang, and Xiao-Xue Liang. Exploring inter-channel correlation for diversity-preserved knowledge distillation. *2021 ICCV*, 2021. 14
- [46] Yu Liu, Xuhui Jia, Mingxing Tan, Raviteja Vemulapalli, Yukun Zhu, Bradley Green, and Xiaogang Wang. Search to distill: Pearls are everywhere but not the eyes. In *CVPR*, 2020. 1, 3
- [47] David Lopez-Paz, Léon Bottou, Bernhard Schölkopf, and Vladimir Vapnik. Unifying distillation and privileged information. *arXiv preprint arXiv:1511.03643*, 2015. 5
- [48] Joe Mellor, Jack Turner, Amos Storkey, and Elliot J Crowley. Neural architecture search without training. In *ICML*, 2021. 3, 4, 6, 7, 8, 9, 10, 13
- [49] Seyed Iman Mirzadeh, Mehrdad Farajtabar, Ang Li, Nir Levine, Akihiro Matsukawa, and Hassan Ghasemzadeh. Improved knowledge distillation via teacher assistant. In *AAAI*, 2020. 1, 3, 8
- [50] Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. Relational knowledge distillation. In *CVPR*, 2019. 6, 7
- [51] Wonpyo Park, Yan Lu, Minsu Cho, and Dongju Kim. Relational knowledge distillation. In *CVPR*, 2019. 3, 7
- [52] Nikolaos Passalis and Anastasios Tefas. Learning deep representations with probabilistic knowledge transfer. In *ECCV*, 2018. 3
- [53] Nikolaos Passalis and Anastasios Tefas. Learning deep representations with probabilistic knowledge transfer. In *ECCV*, 2018. 6
- [54] Baoyun Peng, Xiao Jin, Jiaheng Liu, Dongsheng Li, Yichao Wu, Yu Liu, Shunfeng Zhou, and Zhaoning Zhang. Correlation congruence for knowledge distillation. In *ICCV*, 2019. 6, 11
- [55] Baoyun Peng, Xiao Jin, Jiaheng Liu, Shunfeng Zhou, Yichao Wu, Yu Liu, Dong-sheng Li, and Zhaoning Zhang. Correlation congruence for knowledge distillation. In *ICCV*, 2019. 3
- [56] Houwen Peng, Hao Du, Hongyuan Yu, Qi Li, Jing Liao, and Jianlong Fu. Cream of the crop: Distilling prioritized paths for one-shot neural architecture search. *NeurIPS*, 2020. 3
- [57] Hieu Pham, Melody Y. Guan, Barret Zoph, Quoc V. Le, and Jeff Dean. Efficient neural architecture search via parameter sharing. In *ICML*, 2018. 3, 10, 11, 12, 13
- [58] Jie Qin, Jie Wu, Xuefeng Xiao, Lujun Li, and Xingang Wang. Activation modulation and recalibration scheme for weakly supervised semantic segmentation. In *AAAI*, 2022. 1
- [59] Ilija Radosavovic, Justin Johnson, Saining Xie, Wan-Yen Lo, and Piotr Dollár. On network design spaces for visual recognition. *2019 ICCV*, 2019. 12
- [60] Ilija Radosavovic, Raj Prateek Kosaraju, Ross B. Girshick, Kaiming He, and Piotr Dollár. Designing network design spaces. *CVPR*, 2020. 6, 7
- [61] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. In *ICLR*, 2015. 1, 3, 7, 8
- [62] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. *ICLR*, 2015. 6, 7
- [63] R. R. Selvaraju, Abhishek Das, Ramakrishna Vedantam, Michael Cogswell, D. Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. 2019. 2, 10
- [64] Suraj Srinivas and François Fleuret. Knowledge transfer with jacobian matching. In *ICML*, 2018. 6
- [65] Hidenori Tanaka, Daniel Kunin, Daniel L Yamins, and Surya Ganguli. Pruning neural networks without any data by iteratively conserving synaptic flow. *NeurIPS*, 2020. 3, 4, 6
- [66] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive representation distillation. In *ICLR*, 2020. 3, 6, 7, 8
- [67] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive representation distillation. In *ICLR*, 2020. 7, 8, 11
- [68] Frederick Tung and Greg Mori. Similarity-preserving knowledge distillation. In *ICCV*, 2019. 1
- [69] Frederick Tung and Greg Mori. Similarity-preserving knowledge distillation. In *ICCV*, 2019. 6, 7, 14
- [70] Vladimir Vapnik. *Statistical learning theory*. 1998. 5
- [71] Likang Wang and Lei Chen. Dionysus: Recovering scene structures by dividing into semantic pieces. 1
- [72] Likang Wang and Lei Chen. Ftso: Effective nas via first topology second operator. 2023. 1
- [73] Likang Wang, Yue Gong, Xinjun Ma, Qirui Wang, Kaixuan Zhou, and Lei Chen. Is-mvsnet: Importance sampling-based mvsnet. In *ECCV*, 2022. 1
- [74] Likang Wang, Yue Gong, Qirui Wang, Kaixuan Zhou, and Lei Chen. Flora: dual-frequency loss-compensated real-time monocular 3d video reconstruction. In *AAAI*, 2023. 1



- [75] Xiu-Shen Wei, Jian-Hao Luo, Jianxin Wu, and Zhi-Hua Zhou. Selective convolutional descriptor aggregation for fine-grained image retrieval. TIP, 2017. 10
- [76] Zimian Wei, Hengyue Pan, Lujun Li Li, Menglong Lu, Xin Niu, Peijie Dong, and Dongsheng Li. Convformer: Closing the gap between cnn and vision transformers. arXiv preprint arXiv:2209.07738, 2022. 1
- [77] Bichen Wu, Xiaoliang Dai, Peizhao Zhang, Yanghan Wang, Fei Sun, Yiming Wu, Yuandong Tian, Peter Vajda, Yangqing Jia, and Kurt Keutzer. Fbnet: Hardware-aware efficient convnet design via differentiable neural architecture search. In CVPR, 2019. 3
- [78] Liu Xiaolong, Li Lujun, Li Chao, and Anbang Yao. Norm: Knowledge distillation via n-to-one representation matching. In ICLR, 2023. 3
- [79] Jingjing Xu, Liang Zhao, Junyang Lin, Rundong Gao, Xu Sun, and Hongxia Yang. Knas: Green neural architecture search. In ICML, 2021. 3
- [80] Jing Yang, Brais Martinez, Adrian Bulat, Georgios Tzimiropoulos, et al. Knowledge distillation via softmax regression representation learning. ICLR, 2021. 8
- [81] Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In ICLR, 2017. 1, 3
- [82] Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. ICLR, 2017. 7
- [83] Ying Zhang, Tao Xiang, Timothy M Hospedales, and Huchuan Lu. Deep mutual learning. In CVPR, 2018. 8
- [84] Borui Zhao, Quan Cui, Renjie Song, Yiyu Qiu, and Jiajun Liang. Decoupled knowledge distillation. In CVPR, 2022. 1
- [85] Xiawu Zheng, Xiang Fei, Lei Zhang, Chenglin Wu, Fei Chao, Jianzhuang Liu, Wei Zeng, Yonghong Tian, and Rongrong Ji. Neural architecture search with representation mutual information. CVPR, 2022. 3
- [86] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In CVPR, 2016. 4
- [87] Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. Learning transferable architectures for scalable image recognition. In CVPR, 2018. 3