

---

# OTOv3: Towards Automatic Sub-Network Search Within General Super Deep Neural Networks

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Existing neural architecture search (NAS) methods typically rely on pre-specified  
2 super deep neural networks (super-networks) with handcrafted search spaces be-  
3 forehand. Such requirements make it challenging to extend them onto general  
4 scenarios without significant human expertise and manual intervention. To over-  
5 come the limitations, we propose the third generation of Only-Train-Once (OTOv3).  
6 OTOv3 is perhaps the first automated system that trains general super-networks and  
7 produces high-performing sub-networks in the one shot manner without pretraining  
8 and fine-tuning. Technologically, OTOv3 delivers three noticeable contributions  
9 to minimize human efforts: (i) automatic search space construction for general  
10 super-networks; (ii) a Hierarchical Half-Space Projected Gradient (H2SPG) that  
11 leverages the dependency graph to ensure the network validity during optimization  
12 and reliably produces a solution with both high performance and hierarchical group  
13 sparsity; and (iii) automatic sub-network construction based on the super-network  
14 and the H2SPG solution. Numerically, we demonstrate the effectiveness of OTOv3  
15 on a variety of super-networks, including StackedUnets, SuperResNet, and DARTS,  
16 over benchmark datasets such as CIFAR10, Fashion-MNIST, ImageNet, STL-10,  
17 and SVNH. The sub-networks computed by OTOv3 achieve competitive even  
18 superior performance compared to the super-networks and other state-of-the-arts.

## 19 1 Introduction

20 Deep neural networks (DNNs) have achieved remarkable success in various fields, which success is  
21 highly dependent on their sophisticated underlying architectures (LeCun et al., 2015; Goodfellow  
22 et al., 2016). To design effective DNN architectures, human expertise have handcrafted numerous  
23 popular DNNs such as ResNet (He et al., 2016) and transformer (Vaswani et al., 2017). However,  
24 such human efforts may not be scalable enough to meet the increasing demands for customizing  
25 DNNs for diverse tasks. To address this issue, Neural Architecture Search (NAS) has emerged to  
26 automate the network creations and reduce the need for human expertise (Elsken et al., 2018).

27 Among current NAS studies, gradient-based methods (Liu et al., 2018; Yang et al., 2020; Xu et al.,  
28 2019; Chen et al., 2021b) are perhaps the most popular because of their efficiency. Such methods  
29 build an over-parameterized super-network covering all candidate connections and operations, param-  
30 eterize operations via introducing auxiliary architecture variables with weight sharing, then search a  
31 (sub)optimal sub-network via formulating and solving a multi-level optimization problem.

32 Despite the advancements in gradient-based methods, their usage is still limited due to certain  
33 inconvenience. In particular, their automation relies on manually determining the search space for a  
34 pre-specified super-network beforehand, and requires the manual introduction of auxiliary architecture  
35 variables onto the prescribed search space. To extend these methods onto other super-networks, the  
36 users still need to manually construct the search pool, then incorporate the auxiliary architecture

37 variables along with building the whole complicated multi-level optimization training pipeline. The  
 38 whole process necessitates significant domain-knowledge and engineering efforts, thereby being  
 39 inconvenient and time-consuming for users. Therefore, it is natural to ask whether we could reach an

40 **Objective.** *Given a general super-network, automatically generate its search space, train it once, and*  
 41 *construct a sub-network that achieves a dramatically compact architecture and high performance.*

42 Achieving the objective is severely challenging in terms  
 43 of both engineering developments and algorithmic de-  
 44 signs, consequently not achieved yet by the existing  
 45 NAS works to the best of our knowledge. However, the  
 46 objective has been recently achieved in an analogous  
 47 task so-called structured pruning (Lin et al., 2019) by  
 48 the second generation of Only-Train-Once framework (OTOv2) (Chen et al., 2021a, 2023). From  
 49 the perspective of computational graph, the standard NAS could be considered as removing entire  
 50 redundant connections (cutting edges) and operations (vertices) from super-networks. Structured  
 51 pruning can be largely interpreted as a complementary NAS that removes the redundancy inside each  
 52 vertex (slims operations) but preserves all the connections. OTOv2 first achieves the objective in the  
 53 view of structured pruning that given a general DNN, automatically trains it only once to achieve  
 54 both high performance and a slimmer model architecture without pre-training and fine-tuning.

	OTOv3	OTOv2	Other NAS
General DNNs	✓	✓	✗
Autonomy	✓	✓	✗
Remove Connections	✓	✗	✓
Remove Operations	✓	✗	✓
Slim Operations	✓ <sup>†</sup>	✓	✗

<sup>†</sup> Support while is not the focus and discussed in this work.

55 We now build the third-generation of Only-Train-Once  
 56 (OTOv3) that reaches the objective from the perspective  
 57 of the standard NAS. OTOv3 automatically generates a  
 58 search space given a general super-network, trains and  
 59 identifies redundant connections and vertices, then builds  
 60 a sub-network that achieves both high performance and  
 61 compactness. As the library usage presented aside, the  
 62 whole procedure can be automatically proceeded, dramatically reduce the human efforts, and fit for  
 63 general super-networks and applications. Our main contributions can be summarized as follows.

```

OTOv3 Library Usage
1 from only_train_once import OTO
2 # General Super-Network
3 oto = OTO(super_net, cut_edges=True)
4 optimizer = oto.h2spg()
5 # Train as normal
6 optimizer.step()
7 oto.construct_subnet(cut_edges=True)

```

- 64 • **Infrastructure for Automated General Super-Network Training and Sub-Network Searching.**  
 65 We propose OTOv3 that perhaps the first automatically trains and searches within a general super-  
 66 network to deliver a compact sub-network by erasing redundant connections and operations in the  
 67 one-shot manner. As the previous OTO versions, OTOv3 trains the super-network only once without  
 68 the need of pre-training and fine-tuning and is pluggable into various deep learning applications.
- 69 • **Automated Search Space Generation.** We propose a novel graph algorithm to automatically  
 70 explore and establish a dependency graph given a general super-network, then analyze the de-  
 71 pendency to form a search space consisting of minimal removal structures. The corresponding  
 72 trainable variables are then partitioned into so-called generalized zero-invariant groups (GeZIGs).
- 73 • **Hierarchical Half-Space Projected Gradient (H2SPG).** We propose a novel H2SPG optimizer  
 74 that perhaps the first solves a hierarchical structured sparsity problem for general DNNs. H2SPG  
 75 computes a solution  $\mathbf{x}_{\text{H2SPG}}^*$  of both high performance and desired hierarchical group sparsity  
 76 in the manner of GeZIGs. Compared to other optimizers, H2SPG considers the hierarchy of  
 77 dependency graph to produce sparsity for ensuring the validness of the subsequent sub-network.
- 78 • **Automated Sub-Network Construction.** We propose a novel graph algorithm to automatically  
 79 construct a sub-network upon the super-network parameterized as  $\mathbf{x}_{\text{H2SPG}}^*$ . The resulting sub-  
 80 network returns the exact same outputs as the super-network thereby no need of further fine-tuning.
- 81 • **Experimental Results.** We demonstrate the effectiveness of OTOv3 on extensive super-networks  
 82 including StackedUnets, SuperResNet and DARTS, over benchmark datasets including CIFAR10,  
 83 Fashion-MNIST, ImageNet, STL-10, and SVNH. OTOv3 is the first framework that could auto-  
 84 matically deliver compact sub-networks upon general super-networks to the best of our knowledge.  
 85 Meanwhile the sub-networks exhibit competitive even superior performance to the super-networks.

## 86 2 Related Work

87 **Neural Architecture Search (NAS).** Early NAS works utilized reinforcement learning and evolu-  
 88 tion techniques to search for high-quality architectures (Zoph & Le, 2016; Pham et al., 2018; Zoph  
 89 et al., 2018), while they were computationally expensive. Later on, differentiable (gradient-based)

90 methods were introduced to accelerate the search process. These methods start with a super-network  
 91 covering all possible connection and operation candidates, and parameterize them with auxiliary  
 92 architecture variables. They establish a multi-level optimization problem that alternately updates  
 93 the architecture and network variables until convergence (Liu et al., 2018; Chen et al., 2019; Xu et al.,  
 94 2019; Yang et al., 2020; Hosseini & Xie, 2022). However, these methods require a significant amount  
 95 of **handcraftness** from users in advance to **manually** establish the search space, introduce additional  
 96 architecture variables, and build the multi-level training pipeline. The sub-network construction is  
 97 also network-specific and not flexible. All requirements necessitate remarkable domain-knowledge  
 98 and expertise, making it difficult to extend to general super-networks and broader scenarios.

99 **Automated Structured Pruning for General DNNs.** Structure pruning is an orthogonal but related  
 100 paradigm to standard NAS. Rather than removing entire operations and connections, it focuses on  
 101 slimming individual vertices (Han et al., 2015). Similarly, prior structure pruning methods also  
 102 required numerous handcraftness and domain knowledge, which limited their broader applicability.  
 103 However, recent methods such as OTOv2 (Chen et al., 2023) and DepGraph (Fang et al., 2023) have  
 104 made progress in automating the structure pruning process for general DNNs. OTOv2 is a one-shot  
 105 method that does not require pre-training or fine-tuning, while DepGraph involves a multi-stage  
 106 training pipeline that requires some manual intervention. In this work, we propose the third-generation  
 107 version of OTO that enables automatic sub-network searching and training for general super-networks.

108 **Hierarchical Structured Sparsity Optimization.** We formulate the underlying optimization  
 109 problem of OTOv3 as a hierarchical structured sparsity problem. Its solution possesses high group  
 110 sparsity indicating redundant structures and obeys specified hierarchy. There exist deterministic  
 111 optimizers solving such problems via introducing latent variables (Zhao et al., 2009), while are  
 112 impractical for stochastic DNN tasks. Meanwhile, stochastic optimizers rarely study such problem.  
 113 In fact, popular stochastic sparse optimizers such as HSPG (Chen et al., 2021a), DHSPG (Chen et al.,  
 114 2023), proximal methods (Xiao & Zhang, 2014) and ADMM (Lin et al., 2019) overlook the hierarchy  
 115 constraint. Incorporating them into OTOv3 typically delivers invalid sub-networks. Therefore, we  
 116 propose H2SPG that considers graph dependency to solve it for general DNNs.

### 118 3 OTOv3

119 OTOv3 is an automated one-shot system that trains a general super-network and constructs a sub-  
 120 network. The produced sub-network is not only high-performing but also has a dramatically compact  
 121 architecture that is suitable for various shipping environments. The entire process minimizes the need  
 122 for human efforts and is suitable for general DNNs. As outlined in Algorithm 1, given a general super-  
 123 network  $\mathcal{M}$ , OTOv3 first explores and establishes a dependency graph. Upon the dependency graph,  
 124 a search space is automatically constructed and corresponding trainable variables are partitioned  
 125 into generalized zero-invariant groups (GeZIGs) (Section 3.1). A hierarchical structured sparsity  
 126 optimization problem is then formulated and solved by a novel Hierarchical Half-Space Projected  
 127 Gradient (H2SPG) (Section 3.2). H2SPG considers the hierarchy inside the dependency graph and  
 128 computes a solution  $\mathbf{x}_{\text{H2SPG}}^*$  of both high-performance and desired hierarchical group sparsity over  
 129 GeZIGs. A compact sub-network  $\mathcal{M}^*$  is finally constructed via removing the structures corresponding  
 130 to the identified redundant GeZIGs and their dependent structures (Section 3.3).  $\mathcal{M}^*$  returns the exact  
 131 same output as the super-network parameterized as  $\mathbf{x}_{\text{H2SPG}}^*$ , eliminating the need of fine-tuning.

---

#### Algorithm 1 Outline of OTOv3.

---

- 1: **Input:** A general DNN  $\mathcal{M}$  as super-network to be trained and searched (no need to be pretrained).
  - 2: **Automated Search Space Construction.** Establish dependency graph and partition the trainable parameters of  $\mathcal{M}$  into generalized zero-invariant groups  $\mathcal{G}_{\text{GeZIG}}$  and the complementary  $\mathcal{G}_{\text{GeZIG}}^C$ .
  - 3: **Train by H2SPG.** Seek a high-performing solution  $\mathbf{x}_{\text{H2SPG}}^*$  with hierarchical group sparsity.
  - 4: **Automated Sub-Network  $\mathcal{M}^*$  Construction.** Construct a sub-network upon  $\mathbf{x}_{\text{H2SPG}}^*$ .
  - 5: **Output:** Constructed sub-network  $\mathcal{M}^*$  (no need to be fine-tuned).
- 

### 132 3.1 Automated Search Space Construction

133 The foremost step is to automatically construct the search space for a general super-network. However,  
 134 this process presents significant challenges in terms of both engineering developments and algorithmic  
 135 designs due to the complexity of DNN architecture and the lack of sufficient public APIs. To overcome

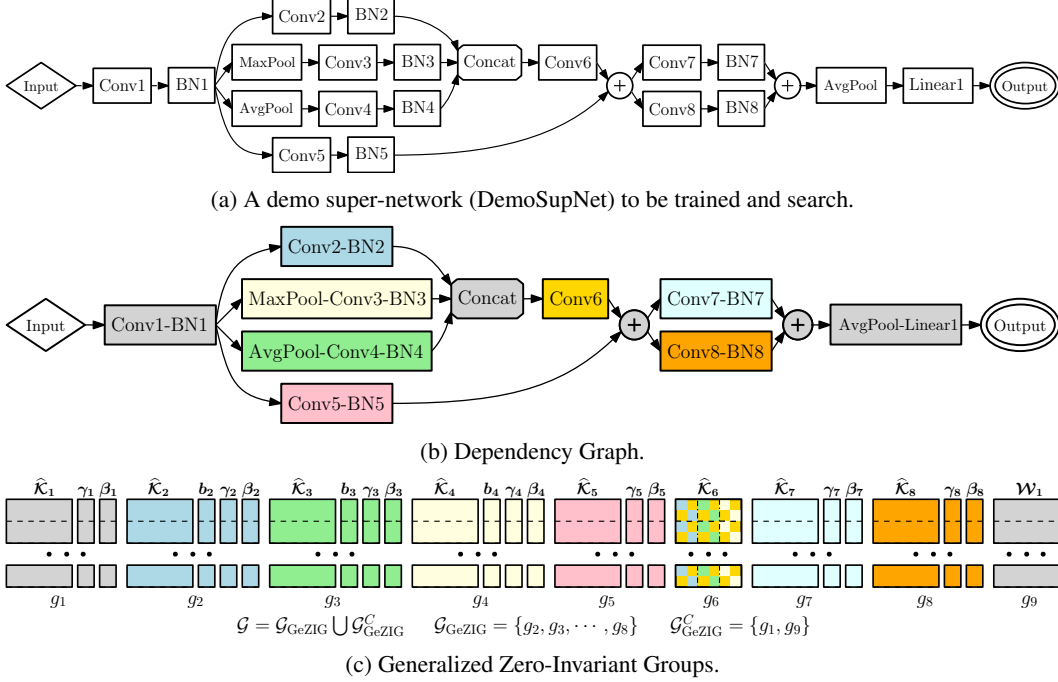


Figure 1: Automated Search Space Construction.  $\hat{\mathcal{K}}_i$  and  $b_i$  are the flatten filter matrix and bias vector for Conv- $i$ , respectively.  $\gamma_i$  and  $\beta_i$  are the weight and bias vectors for BN- $i$ .  $\mathcal{W}_i$  is the weight matrix for Linear- $i$ . The columns of  $\hat{\mathcal{K}}_6$  are marked in accordance to its incoming segments.

136 these challenges, we propose a concept called generalized zero-invariant group (GeZIG) and formulate  
 137 the search space construction as the GeZIG partition. We have also developed a dedicated graph  
 138 algorithm to automatically conduct the GeZIG partition for general super-networks.

139 **Generalized Zero-Invariant Group (GeZIG).** The key of search space construction is to figure  
 140 out the structures that can be removed from the super-network. Because of diverse roles of operations  
 141 and their complicated connections inside a DNN, removing an arbitrary structure may cause the  
 142 remaining DNN invalid. We say a structure *removal* if and only if the DNN after removing it is still  
 143 valid. A removal structure is further said *minimal* if and only if it does not contain multiple removal  
 144 structures. Zero-Invariant Group (ZIG) is proposed in (Chen et al., 2021a, 2023) that describes a  
 145 class of minimal removal structures satisfying a *zero-invariant property*, *i.e.*, if all variables in ZIG  
 146 equal to zero, then no matter what the input is, the output is always as zero. ZIG depicts the minimal  
 147 removal structure *inside each operation* and is the key for realizing automatic one-shot structured  
 148 pruning. We generalize ZIG as GeZIG that describes a class of minimal removal structures satisfying  
 149 the zero-invariant property but *consists of entire operations*. More illustrations regarding ZIG versus  
 150 GeZIG are present in Appendix. For simplicity, throughout the paper, the minimal removal structure  
 151 is referred to the counterpart consisting of operations in entirety. Consequently, automated search  
 152 space construction becomes how to automatically explore the GeZIG partition for general DNNs.

153 **Automated GeZIG Partition.** As specified in Algorithm 2, automated GeZIG partition involves  
 154 two main stages. The first stage explores the super-network  $\mathcal{M}$  and establishes a dependency graph  
 155  $(\mathcal{V}_d, \mathcal{E}_d)$ . The second stage leverages the affiliations inside the dependency graph to find out minimal  
 156 removal structures, then partitions their trainable variables to form GeZIGs. For intuitive illustrations,  
 157 we elaborate the algorithm through a small but complex demo super-network depicted in Figure 1a.

158 **Dependency Graph Construction.** Given a super-network  $\mathcal{M}$ , we first construct its trace graph  $(\mathcal{V}, \mathcal{E})$   
 159 displayed as Figure 1a (line 3 in Algorithm 2), where  $\mathcal{V}$  represents the set of vertices (operations) and  
 160  $\mathcal{E}$  represents the connections among them. As OTov2 (Chen et al., 2023), we categorize the vertices  
 161 into stem vertices, joint vertices, accessory vertices, and unknown vertices. Stem vertices refer to the  
 162 operations that contain trainable variables and can transform the input tensors into different shapes,  
 163 *e.g.*, Conv and Linear. The accessory vertices are the operations that may not have trainable  
 164 variables and have a single input, *e.g.*, BN and ReLU. Joint vertices aggregate multiple inputs into a  
 165 single output, *e.g.*, Add and Concat. The remaining vertices are considered as unknown.

---

**Algorithm 2** Automated Search Space Construction.
 

---

- 1: **Input:** A super-network  $\mathcal{M}$  to be trained and searched.
  - 2: **Dependency graph construction.**
  - 3: Construct the trace graph  $(\mathcal{E}, \mathcal{V})$  of  $\mathcal{M}$ .
  - 4: Initialize an empty graph  $(\mathcal{V}_d, \mathcal{E}_d)$ .
  - 5: Initialize queue  $\mathcal{Q} \leftarrow \{\mathcal{S}(v) : v \in \mathcal{V} \text{ is adjacent to the input of trace graph}\}$ .
  - 6: **while**  $\mathcal{Q} \neq \emptyset$  **do**
  - 7: Dequeue the head segment  $\mathcal{S}$  from  $\mathcal{Q}$ .
  - 8: Grow  $\mathcal{S}$  in the depth-first manner till meet either joint vertex or multi-outgoing vertex  $\hat{v}$ .
  - 9: Add segments into  $\mathcal{V}_d$  and connections into  $\mathcal{E}_d$ .
  - 10: Enqueue new segments into the tail of  $\mathcal{Q}$  if  $\hat{v}$  has outgoing vertices.
  - 11: **Find minimal removal structures.**
  - 12: Get the incoming vertices  $\hat{\mathcal{V}}$  for joint vertices in the  $(\mathcal{V}_d, \mathcal{E}_d)$ .
  - 13: Group the trainable variables in the vertex  $v \in \hat{\mathcal{V}}$  as  $g_v$ .
  - 14: Form  $\mathcal{G}_{\text{GeZIG}}$  as the union of the above groups, i.e.,  $\mathcal{G}_{\text{GeZIG}} \leftarrow \{g_v : v \in \hat{\mathcal{V}}\}$ .
  - 15: Form  $\mathcal{G}_{\text{GeZIG}}^C$  as the union of the trainable variables in the remaining vertices.
  - 16: **Return** trainable variable partition  $\mathcal{G} = \mathcal{G}_{\text{GeZIG}} \cup \mathcal{G}_{\text{GeZIG}}^C$  and dynamic dependency graph  $(\mathcal{V}_d, \mathcal{E}_d)$ .
- 

166 We begin by analyzing the trace graph  $(\mathcal{V}, \mathcal{E})$  to create a dependency graph  $(\mathcal{V}_d, \mathcal{E}_d)$ , wherein each  
 167 vertex in  $\mathcal{V}_d$  serves as a potential minimal removal structure candidate. To proceed, we use a queue  
 168 container  $\mathcal{Q}$  to track the candidates (line 5 of Algorithm 2). The initial elements of this queue are  
 169 the vertices that are directly adjacent to the input of  $\mathcal{M}$ , such as `Conv1`. We then traverse the graph  
 170 in the breadth-first manner, iteratively growing each element (segment)  $\mathcal{S}$  in the queue until a valid  
 171 minimal removal structure candidate is formed. The growth of each candidate follows the depth-first  
 172 search to recursively expand  $\mathcal{S}$  until the current vertices are considered as endpoints. The endpoint  
 173 vertex is determined by whether it is a joint vertex or has multiple outgoing vertices, as indicated  
 174 in line 8 of Algorithm 2. Intuitively, a joint vertex has multiple inputs, which means that the DNN  
 175 may be still valid after removing the current segment. This suggests that the current segment may  
 176 be removable. On the other hand, a vertex with multiple outgoing neighbors implies that removing  
 177 the current segment may cause some of its children to miss the input tensor. For instance, removing  
 178 `Conv1-BN1` would cause `Conv2`, `MaxPool` and `AvgPool` to become invalid due to the absence  
 179 of input in Figure 1a. Therefore, it is risky to remove such candidates. Once the segment  $\mathcal{S}$  has been  
 180 grown, new candidates are initialized as the outgoing vertices of the endpoint and added into the  
 181 container  $\mathcal{Q}$  (line 10 in Algorithm 2). Such procedure is repeated until the end of graph traversal.  
 182 Ultimately, a dependency graph  $(\mathcal{V}_d, \mathcal{E}_d)$  is created, as illustrated in Figure 1b.

183 **Form GeZIGs.** We proceed to identify the minimal removal structures in  $(\mathcal{V}_d, \mathcal{E}_d)$  to create the GeZIG  
 184 partition. The qualified instances are the vertices in  $\mathcal{V}_d$  that have trainable variables and all of their  
 185 outgoing vertices are joint vertices. This is because a joint vertex has multiple inputs and remains  
 186 valid even after removing some of its incoming structures, as indicated in line 12 in Algorithm 2.  
 187 Consequently, their trainable variables are grouped together into GeZIGs (line 13-14 in Algorithm 2  
 188 and Figure 1c). The remaining vertices are considered as either unremovable or belonging to a  
 189 large removal structure, which trainable variables are grouped into the  $\mathcal{G}_{\text{GeZIG}}^C$  (the complementary  
 190 to  $\mathcal{G}_{\text{GeZIG}}$ ). As a result, for the super-network  $\mathcal{M}$ , all its trainable variables are encompassed by the  
 191 union  $\mathcal{G} = \mathcal{G}_{\text{GeZIG}} \cup \mathcal{G}_{\text{GeZIG}}^C$ , and the corresponding structures in  $\mathcal{G}_{\text{GeZIG}}$  constitute its search space.

### 192 3.2 Hierarchical Half-Space Projected Gradient (H2SPG)

193 Given a super-network  $\mathcal{M}$  and its group partition  $\mathcal{G} = \mathcal{G}_{\text{GeZIG}} \cup \mathcal{G}_{\text{GeZIG}}^C$ , the next is to jointly search  
 194 for a valid sub-network  $\mathcal{M}^*$  that exhibits the most significant performance and train it to high  
 195 performance. Searching a sub-network is equivalent to identifying the redundant structures in  $\mathcal{G}_{\text{GeZIG}}$   
 196 to be further removed and ensures the remaining network still valid. Training the sub-network  
 197 becomes optimizing over the remaining groups in  $\mathcal{G}$  to achieve high performance. We formulate a  
 198 hierarchical structured sparsity problem to accomplish both tasks simultaneously as follows.

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} f(\mathbf{x}), \text{ s.t. } \text{Cardinality}(\mathcal{G}^0) = K, \text{ and } (\mathcal{V}_d/\mathcal{V}_{\mathcal{G}^0}, \mathcal{E}_d/\mathcal{E}_{\mathcal{G}^0}) \text{ is valid}, \quad (1)$$

199 where  $f$  is the prescribed loss function,  $\mathcal{G}^0 := \{g \in \mathcal{G}_{\text{GeZIG}} | [\mathbf{x}]_g = 0\}$  is the set of zero groups in  
 200  $\mathcal{G}_{\text{GeZIG}}$ , which cardinality measures its size.  $K$  is the target group sparsity, indicating the number of



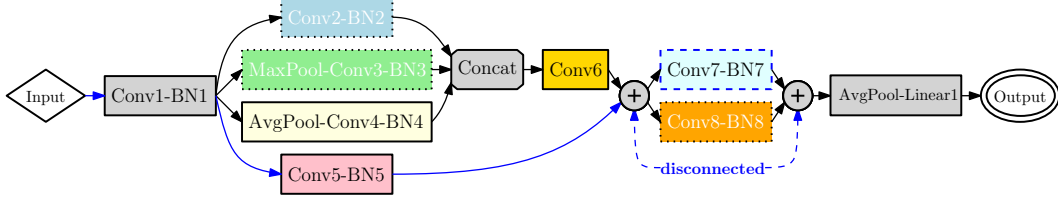


Figure 2: Check validness of redundant candidates. Target group sparsity  $K = 3$ . Conv7-BN7 has larger redundancy score than Conv2-BN2. Dotted vertices are marked as redundant candidates.

201 GeZIGs that should be identified as redundant. The redundant GeZIGs are projected onto zero, while  
 202 the important groups are preserved as non-zero and optimized for high performance. A larger  $K$   
 203 dictates a higher sparsity level that produces a more compact sub-network with fewer FLOPs and  
 204 parameters.  $(\mathcal{V}_d/\mathcal{V}_{\mathcal{G}^0}, \mathcal{E}_d/\mathcal{E}_{\mathcal{G}^0})$  refers to the graph removing vertices and edges corresponding to  
 205 zero groups  $\mathcal{G}^0$ . This graph being valid is specified for NAS that requires the zero groups distributed  
 206 obeying the hierarchy of super-network to ensure the resulting sub-network functions correctly.

207 Problem (1) is difficult to solve due to the non-differential and non-convex sparsity constraint and  
 208 the graph validity constraint. Existing optimizers such as DHSPG (Chen et al., 2023) overlook the  
 209 architecture evolution and hierarchy during the sparsity exploration, which is crucial to (1). In fact,  
 210 they are mainly applied for pruning tasks, where the connections and operations are preserved (but  
 211 become slimmer). Consequently, employing them onto (1) usually produces invalid sub-networks.

212 **Outline of H2SPG.** To effectively  
 213 solve problem (1), we propose a novel  
 214 H2SPG to consider the hierarchy and  
 215 ensure the validness of graph architec-  
 216 ture after removing redundant vertices  
 217 and connections during the optimiza-  
 218 tion process. To the best of our knowl-  
 219 edge, H2SPG is the first the optimizer  
 220 that successfully solves such hierar-  
 221 chical structured sparsity problem (1),  
 222 which outline is stated in Algorithm 3.

223 H2SPG is built upon the DHSPG in  
 224 OTOv2 but with dedicated designs  
 225 regarding the hierarchical constraint.  
 226 In general, H2SPG is a hybrid multi-  
 227 phase optimizer that first partitions  
 228 the groups of variables into impor-  
 229 tant and potentially redundant seg-  
 230 ments, then employs specified updat-  
 231 ing mechanisms onto different seg-  
 232 ments to achieve a solution with both  
 233 desired hierarchical group sparsity  
 234 and high performance. The variable partition considers the hierarchy of dependency graph  $(\mathcal{V}_d, \mathcal{E}_d)$  to  
 235 ensure the validness of the resulting sub-network graph. Vanilla stochastic gradient descent (SGD) or  
 236 its variant such as Adam (Kingma & Ba, 2014) optimizes the important variables to achieve the high  
 237 performance. Half-space gradient descent (Chen et al., 2021a) identifies redundant groups among the  
 238 candidates and projects them onto zero without sacrificing the objective function to the largest extent.

239 **Warm-Up Phase.** To proceed, H2SPG first warms up all variables by conducting SGD or its variants  
 240  $T_w$  steps (line 4-5 in Algorithm 3). During each warm-up step  $t$ , a redundancy score of each group  
 241  $g \in \mathcal{G}_{\text{GeZIG}}$  is computed upon the current iterate  $\mathbf{x}_t$  and exponentially averaged by a momentum  
 242 coefficient  $\omega$  (line 6-7 in Algorithm 3). Larger redundancy score indicates the group exhibits less  
 243 prediction power, thus may be redundant. The redundancy score calculation is modular, where we  
 244 follow DHSPG to consider the cosine similarity between negative gradient  $-\nabla f(\mathbf{x}_t)|_g$  and the  
 245 projection direction  $-\mathbf{x}_g$  as well as the average variable magnitude. After warm-up, the redundancy  
 246 scores of all groups in  $\mathcal{G}_{\text{GeZIG}}$  are sorted. We then perform a sanity check and select the groups with  
 247 top- $K$  redundancy scores as the redundant group candidates  $\mathcal{G}_r \subseteq \mathcal{G}_{\text{GeZIG}}$ . The complementary groups

---

**Algorithm 3** Hierarchical Half-Space Projected Gradient

---

- 1: **Input:** initial variable  $\mathbf{x}_0 \in \mathbb{R}^n$ , initial learning rate  $\alpha_0$ , warm-up steps  $T_w$ , target group sparsity  $K$ , momentum  $\omega$ , dependency graph  $(\mathcal{V}_d, \mathcal{E}_d)$  and group partitions  $\mathcal{G}$ .
  - 2: **Warm-up Phase.**
  - 3: **for**  $t = 0, 1, \dots, T_w - 1$  **do**
  - 4:   Calculate gradient estimate  $\nabla f(\mathbf{x}_t)$  or its variant.
  - 5:   Update next iterate  $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \alpha_t \nabla f(\mathbf{x}_t)$ .
  - 6:   Calculate redundancy score  $s_{t,g}$  for  $g \in \mathcal{G}_{\text{GeZIG}}$ .
  - 7:   Update  $s_g \leftarrow \omega s_g + (1 - \omega) s_{t,g}$  for  $g \in \mathcal{G}_{\text{GeZIG}}$ .
  - 8:   Construct  $\mathcal{G}_r$  and  $\mathcal{G}_r^C$  given scores,  $\mathcal{G}$ ,  $(\mathcal{V}_d, \mathcal{E}_d)$ , and  $K$ .
  - 9: **Hybrid Training Phase.**
  - 10: **for**  $t = T_w, T_w + 1, \dots$ , **do**
  - 11:   Compute gradient estimate  $\nabla f(\mathbf{x}_t)$  or its variant.
  - 12:   Update  $[\mathbf{x}_{t+1}]_{\mathcal{G}_r^C}$  as  $[\mathbf{x}_t - \alpha_t \nabla f(\mathbf{x}_t)]_{\mathcal{G}_r^C}$ .
  - 13:   Select proper  $\lambda_g$  for each  $g \in \mathcal{G}_r$ .
  - 14:   Compute  $[\tilde{\mathbf{x}}_{t+1}]_{\mathcal{G}_r}$  via subgradient descent of  $\psi$ .
  - 15:   Perform Half-Space projection over  $[\tilde{\mathbf{x}}_{t+1}]_{\mathcal{G}_r}$ .
  - 16:   Update  $[\mathbf{x}_{t+1}]_{\mathcal{G}_r} \leftarrow [\tilde{\mathbf{x}}_{t+1}]_{\mathcal{G}_r}$ .
  - 17: **Return** the final iterate  $\mathbf{x}_{\text{DHSPG}^+}^*$ .
-

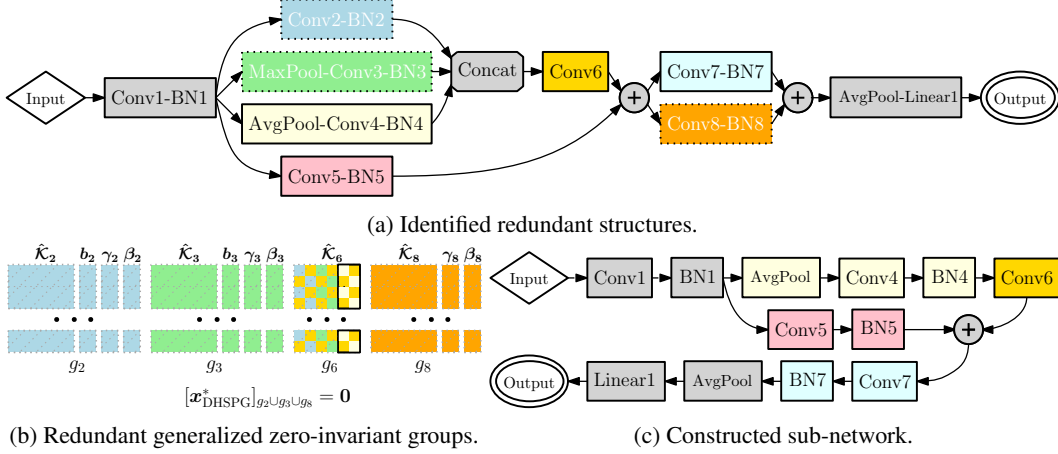


Figure 3: Redundant removal structures identifications and sub-network construction.

248 with lower redundancy scores are marked as important ones and form  $\mathcal{G}_r^C := \mathcal{G}/\mathcal{G}_r$ . The sanity check  
 249 verifies whether the remaining graph is still connected after removing a vertex. If so, the current  
 250 vertex is added into  $\mathcal{G}_r$ ; otherwise, the subsequent vertex is turned into considerations. As illustrated  
 251 in Figure 2, though Conv7-BN7 has a larger redundancy score than Conv2-BN2, Conv2-BN2 is  
 252 marked as potentially redundant but not Conv7-BN7 since there is no path connecting the input and  
 253 the output of the graph after removing Conv7-BN7. This mechanism largely guarantees that even if  
 254 all redundant candidates are erased, the resulting sub-network is still functioning as normal.

255 **Hybrid Training Phase.** H2SPG then engages into the hybrid training phase to produce desired group  
 256 sparsity over  $\mathcal{G}_r$  and optimize over  $\mathcal{G}_r^C$  for pursuing excellent performance till the convergence. This  
 257 phase mainly follows DHSPG (Chen et al., 2023), and we briefly describe the steps for completeness.  
 258 In general, for the important groups of variables in  $\mathcal{G}_r^C$ , the vanilla SGD or its variant is employed to  
 259 minimize the objective function to the largest extent (line 11-12 in Algorithm 3). For redundant group  
 260 candidates in  $\mathcal{G}_r$ , we formulate a relaxed non-constrained subproblem as (2) to gradually reduce the  
 261 magnitudes without deteriorating the objective and project groups onto zeros only if the projection  
 262 serves as a descent direction for the objective during the training process (line 13-16 in Algorithm 3).

$$\underset{[\mathbf{x}]_{\mathcal{G}_r}}{\text{minimize}} \psi([\mathbf{x}]_{\mathcal{G}_r}) := f([\mathbf{x}]_{\mathcal{G}_r}) + \sum_{g \in \mathcal{G}_r} \lambda_g \|\mathbf{x}\|_2, \quad (2)$$

263 where  $\lambda_g$  is a group-specific regularization coefficient and dedicatedly selected as DHSPG. H2SPG  
 264 then performs a subgradient descent of  $\psi$  over  $[\mathbf{x}]_{\mathcal{G}_r}$ , followed by a Half-Space projection (Chen  
 265 et al., 2021a) to effectively produce group sparsity with the minimal sacrifice of the objective function.  
 266 At the end, a high-performing solution  $\mathbf{x}_{\text{H2SPG}}^*$  with desired hierarchical group sparsity is returned.

### 267 3.3 Automated Sub-Network Construction.

268 We finally construct a sub-network  $\mathcal{M}^*$  upon the super-network  $\mathcal{M}$  and the solution  $\mathbf{x}_{\text{H2SPG}}^*$  by  
 269 H2SPG. The solution  $\mathbf{x}_{\text{H2SPG}}^*$  should attain desired target hierarchical group sparsity level and achieve  
 270 high performance. As illustrated in Figure 3, we first traverse the graph to remove the entire vertices  
 271 and the related edges from  $\mathcal{M}$  corresponding to the redundant GeZIGs being zero, e.g., Conv2-BN2,  
 272 MaxPool-Conv3-BN3 and Conv8-BN8 are removed due to  $[\mathbf{x}_{\text{H2SPG}}^*]_{g_2 \cup g_3 \cup g_8} = \mathbf{0}$ . Then, we  
 273 traverse the graph in the second pass to remove the affiliated structures that are dependent on the  
 274 removed vertices to keep the remaining operations valid, e.g., the first and second columns in  $\hat{\mathcal{K}}_6$   
 275 are erased since its incoming vertices Conv2-BN2 and MaxPool-Conv3-BN3 has been removed  
 276 (see Figure 3b). Next, we recursively erase unnecessary vertices and isolated vertices. Isolated  
 277 vertices refer to the vertices that have neither incoming nor outgoing vertices. Unnecessary vertices  
 278 refer to the skippable operations, e.g., Concat and Add (between Conv7 and AvgPool) become  
 279 unnecessary. Ultimately, a compact sub-network  $\mathcal{M}^*$  is constructed as shown in Figure 3c. By the  
 280 definition of GeZIGs, the redundant GeZIGs (have been projected onto zeros) contribute none to the  
 281 model outputs. Consequently, the  $\mathcal{M}^*$  returns the exact same output as the super-network  $\mathcal{M}$  with  
 282  $\mathbf{x}_{\text{H2SPG}}^*$ , which avoids the necessity of further fine-tuning the sub-network.<sup>1</sup>

<sup>1</sup>Remark here that the sub-network is still compatible to be fine-tuned afterwards if needed.

## 283 4 Numerical Experiments

284 In this section, we employ OTOv3 to one-shot automatically train and search within general super-  
 285 networks to construct compact sub-networks with high performance. The numerical demonstrations  
 286 cover extensive super-networks including DemoSupNet shown in Section 3, StackedUnets (Ron-  
 287 neberger et al., 2015; Chen et al., 2023), SuperResNet (He et al., 2016; Lin et al., 2021), and  
 288 DARTS (Liu et al., 2018), and benchmark datasets, including CIFAR10 (Krizhevsky & Hinton,  
 289 2009), Fashion-MNIST (Xiao et al., 2017), ImageNet (Deng et al., 2009), STL-10 (Coates et al.,  
 290 2011) and SVNH (Netzer et al., 2011). More implementation details of experiments and OTOv3  
 291 library and limitations are provided in Appendix A. The dependency graphs and the constructed sub-  
 292 networks are depicted in Appendix C. Ablation studies regarding H2SPG is present in Appendix D.

Table 1: OTOv3 on extensive super-networks and datasets.

Backend	Dataset	Method	FLOPs (M)	# of Params (M)	Top-1 Acc. (%)
DemoSupNet	Fashion-MNIST	Baseline	209	0.82	84.9
DemoSupNet	Fashion-MNIST	<b>OTOv3</b>	<b>107</b>	<b>0.45</b>	<b>84.7</b>
StackedUnets	SVNH	Baseline	184	0.80	95.3
StackedUnets	SVNH	<b>OTOv3</b>	<b>115</b>	<b>0.37</b>	<b>96.1</b>
DARTS (8 cells)	STL-10	Baseline	614	4.05	74.6
DARTS (8 cells)	STL-10	<b>OTOv3</b>	<b>127</b>	<b>0.64</b>	<b>75.1</b>

293 **DemoSupNet on Fashion-MNIST.** We first experiment with the DemoSupNet presented as Fig-  
 294 ure 1a on Fashion-MNIST. OTOv3 automatically establishes a search space of DemoSupNet and  
 295 partitions its trainable variables into GeZIGs. H2SPG then trains DemoSupNet from scratch and  
 296 computes a solution of high performance and hierarchical group-sparsity over GeZIGs, which is  
 297 further utilized to construct a compact sub-network as presented in Figure 3c. As shown in Table 1,  
 298 compared to the super-network, the sub-network utilizes 54% of parameters and 51% of FLOPs to  
 299 achieve a Top-1 validation accuracy 84.7% which is negligibly lower than the super-network by 0.2%.

300 **StackedUnets on SVNH.** We then consider a StackedUnets over SVNH. The StackedUnets is  
 301 constructed by stacking two standard Unets (Ronneberger et al., 2015) with different down-samplers  
 302 together, as depicted in Figure 5a in Appendix C. We employ OTOv3 to automatically build  
 303 the dependency graph, establish the search space, and train by H2SPG. H2SPG identifies and  
 304 projects the redundant structures onto zero and optimize the remaining important ones to attain  
 305 excellent performance. As displayed in Figure 5c, the right-hand-side Unet is disabled due to  
 306 node-72-node-73-node-74-node-75 being zero.<sup>2</sup> The path regarding the deepest depth for  
 307 the left-hand-side Unet, *i.e.*, node-13-node-14-node-15-node-19, is marked as redundant  
 308 as well. The results by OTOv3 indicate that the performance gain brought by either composing multi-  
 309 ple Unets in parallel or encompassing deeper scaling paths is not significant. OTOv3 also validates  
 310 the human design since a single Unet with properly selected depths have achieved remarkable success  
 311 in numerous applications (Ding et al., 2022; Weng et al., 2019). Furthermore, as presented in Table 1,  
 312 the sub-network built by OTOv3 uses 0.37M parameters and 115M FLOPs which is noticeably lighter  
 313 than the full StackedUnets meanwhile significantly outperforms it by 0.8% in validation accuracy.

314 **DARTS (8-Cells) on STL-10.** We next employ OTOv3 on DARTS over STL-10. DARTS is a  
 315 complicated super-network consisting of iteratively stacking multiple cells (Liu et al., 2018). Each  
 316 cell is constructed by spanning a graph wherein every two nodes are connected via multiple operation  
 317 candidates. STL-10 is an image dataset for the semi-supervising learning, where we conduct the  
 318 experiments by using its labeled samples. DARTS has been well explored in the recent years.  
 319 However, the existing NAS methods studied it based on a *handcrafted* search space beforehand to  
 320 *locally* pick up one or two important operations to connect every two nodes. We now employ OTOv3  
 321 on an eight-cells DARTS to *automatically* establish its search space, then utilize H2SPG to one shot  
 322 train it and search important structures *globally* as depicted in Figure 6c of Appendix C. Afterwards,  
 323 a sub-network is automatically constructed as drawn in Figure 6d of Appendix C. Quantitatively, the  
 324 sub-network outperforms the full DARTS in terms of validation accuracy by 0.5% by using only  
 325 about 15%-20% of the parameters and the FLOPs of the original super-network (see Table 1).

<sup>2</sup>Recall the definition of GeZIG, if one GeZIG equals to zero, its output would be always zero given whatever inputs. Therefore, node-72-node-73-node-74-node-75 only produces zero output even if its ancestor vertices may have non-zero parameters. As a result, the right-hand-side Unet is completely disabled.



326 **SuperResNet on CIFAR10.**

327 Later on, we switch to a  
 328 ResNet search space as Zen-  
 329 NAS (Lin et al., 2021), re-  
 330 ferred to as SuperResNet.  
 331 SuperResNet is constructed  
 332 by stacking several super-  
 333 residual blocks with vary-  
 334 ing depths. Each super-  
 335 residual blocks contain multi-  
 336 ple Conv candidates with  
 337 kernel sizes as 3x3, 5x5  
 338 and 7x7 separately in paral-  
 339 lel (see Figure 7a). We then  
 340 employ OTOv3 to one-shot automatically produce two sub-networks with 1M and 2M parameters. As  
 341 displayed in Table 2, the 1M sub-network by OTOv3 outperforms the counterparts reported in (Lin  
 342 et al., 2021) in terms of search cost (on an NVIDIA A100 GPU) due to the efficient single-level  
 343 optimization. The 2M sub-network could reach the benchmark over 97% validation accuracy. Remark  
 344 here that OTOv3 and ZenNAS use networks of fewer parameters to achieve competitive performance  
 345 to the DARTS benchmarks. This is because of the extra data-augmentations such as MixUp (Zhang  
 346 et al., 2017) on this experiment by ZenNAS, so as OTOv3 to follow the same training settings.

Table 2: OTOv3 over SuperResNet on CIFAR10.

Architecture	Top-1 Acc (%)	# of Params (M)	Search Cost (GPU days)
Zen-Score-1M(Lin et al., 2021)	96.2	1.0	0.4
Synflow† (Tanaka et al., 2020)	95.1	1.0	0.4
NASWOT† (Mellor et al., 2021)	96.0	1.0	0.5
Zen-Score-2M(Lin et al., 2021)	97.5	2.0	0.5
SANAS-DARTS (Hosseini & Xie, 2022)	97.5	3.2	1.2*
ISTA-NAS(He et al., 2020)	97.5	3.3	0.1
CDEP (Rieger et al., 2020)	97.2	3.2	1.3*
DARTS (2nd order) (Liu et al., 2018)	97.2	3.1	1.0
PrDARTS (Zhou et al., 2020)	97.6	3.4	0.2
P-DARTS (Chen et al., 2019)	97.5	3.6	0.3
PC-DARTS (Xu et al., 2019)	97.4	3.9	0.1
<b>OTOv3-SuperResNet-1M</b>	96.3	1.0	0.1
<b>OTOv3-SuperResNet-2M</b>	97.5	2.0	0.1

† Reported in (Lin et al., 2021).

\* Numbers are approximately scaled based on (Hosseini & Xie, 2022).

Table 3: OTOv3 over DARTS on ImageNet and comparison with state-of-the-art methods.

Architecture	Test Acc. (%)		# of Params (M)	FLOPs (M)	Search Method
	Top-1	Top-5			
Inception-v1 (Szegedy et al., 2015)	69.8	89.9	6.6	1448	Manual
ShuffleNet 2x (v2) (Ma et al., 2018)	74.9	-	5.0	591	Manual
NASNet-A (Zoph et al., 2018)	74.0	91.6	5.3	564	RL
MnasNet-92 (Tan et al., 2019)	74.8	92.0	4.4	388	RL
AmoebaNet-C (Real et al., 2019)	75.7	92.4	6.4	570	Evolution
DARTS (2nd order) (CIFAR10) (Liu et al., 2018)	73.3	91.3	4.7	574	Gradient
P-DARTS (CIFAR10) (Chen et al., 2019)	75.6	92.6	4.9	557	Gradient
PC-DARTS (CIFAR10) (Xu et al., 2019)	74.9	92.2	5.3	586	Gradient
SANAS (CIFAR10) (Hosseini & Xie, 2022)	75.2	91.7	-	-	Gradient
ProxylessNAS (ImageNet) (Cai et al., 2018)	75.1	92.5	7.1	465	Gradient
PC-DARTs (ImageNet) (Xu et al., 2019)	75.8	92.7	5.3	597	Gradient
ISTA-NAS (ImageNet) (Yang et al., 2020)	76.0	92.9	5.7	638	Gradient
<b>OTOv3 on DARTS (ImageNet)</b>	<b>75.3</b>	<b>92.5</b>	<b>4.8</b>	<b>547</b>	<b>Gradient</b>

(CIFAR10) / (ImageNet) refer to using either CIFAR10 or ImageNet for searching architecture.

347 **DARTS (14-Cells) on ImageNet.** We finally present the benchmark DARTS super-network stacked  
 348 by 14 cells on ImageNet. We employ OTOv3 over it to automatically figure out the search space which  
 349 the code base required specified handcraftness in the past, train by H2SPG to figure out redundant  
 350 structures, and construct a sub-network as depicted in Figure 8d. Quantitatively, we observe that  
 351 the sub-network produced by OTOv3 achieves competitive top-1/5 accuracy compared to other  
 352 state-of-the-arts as presented in Table 3. Remark here that it is *engineeringly* difficult yet to inject  
 353 architecture variables and build a multi-level optimization upon a search space being automatically  
 354 constructed and globally searched. The single-level H2SPG does not leverage a validation set as  
 355 others to favor the architecture search and search over the operations without trainable variables, *e.g.*,  
 356 skip connection, consequently the achieved accuracy does not outperform PC-DARTS and ISTA-NAS.  
 357 We leave further accuracy improvement based on the *automatic* search space as future work.

358 **5 Conclusion**

359 We propose the third generation of Only-Train-Once framework (OTOv3). To the best of knowledge,  
 360 OTOv3 is the first automated system that automatically establishes the search spaces for general  
 361 super-networks, then trains the super-networks via a novel H2SPG optimizer in the one-shot manner,  
 362 finally automatically produces compact sub-networks of high-performance. Meanwhile, H2SPG is  
 363 also perhaps the first stochastic optimizer that effectively solve a hierarchical structured sparsity  
 364 problem for deep learning tasks. OTOv3 further significantly reduces the human efforts upon the  
 365 existing NAS works, opens a new direction and establishes benchmarks regarding the automated  
 366 NAS for the general super-networks which currently require numerous handcraftness beforehand.

## 367 References

- 368 Han Cai, Ligeng Zhu, and Song Han. Proxylesnas: Direct neural architecture search on target task  
369 and hardware. *arXiv preprint arXiv:1812.00332*, 2018.
- 370 Tianyi Chen, Bo Ji, Tianyu Ding, Biyi Fang, Guanyi Wang, Zhihui Zhu, Luming Liang, Yixin  
371 Shi, Sheng Yi, and Xiao Tu. Only train once: A one-shot neural network training and pruning  
372 framework. In *Advances in Neural Information Processing Systems*, 2021a.
- 373 Tianyi Chen, Luming Liang, DING Tianyu, Zhihui Zhu, and Ilya Zharkov. Otov2: Automatic,  
374 generic, user-friendly. In *The Eleventh International Conference on Learning Representations*,  
375 2023.
- 376 Xin Chen, Lingxi Xie, Jun Wu, and Qi Tian. Progressive differentiable architecture search: Bridging  
377 the depth gap between search and evaluation. In *Proceedings of the IEEE/CVF international  
378 conference on computer vision*, pp. 1294–1303, 2019.
- 379 Xin Chen, Lingxi Xie, Jun Wu, and Qi Tian. Progressive darts: Bridging the optimization gap for nas  
380 in the wild. *International Journal of Computer Vision*, 129:638–655, 2021b.
- 381 Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised  
382 feature learning. In *Proceedings of the fourteenth international conference on artificial intelligence  
383 and statistics*, pp. 215–223. JMLR Workshop and Conference Proceedings, 2011.
- 384 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale  
385 hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,  
386 pp. 248–255. Ieee, 2009.
- 387 Tianyu Ding, Luming Liang, Zhihui Zhu, Tianyi Chen, and Ilya Zharkov. Sparsity-guided network  
388 design for frame interpolation. *arXiv preprint arXiv:2209.04551*, 2022.
- 389 Thomas Elsken, Jan Hendrik Metzen, and Frank Hutter. Efficient multi-objective neural architecture  
390 search via lamarckian evolution. *arXiv preprint arXiv:1804.09081*, 2018.
- 391 Gongfan Fang, Xinyin Ma, Mingli Song, Michael Bi Mi, and Xinchao Wang. Depgraph: Towards  
392 any structural pruning. *arXiv preprint arXiv:2301.12900*, 2023.
- 393 Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1.  
394 MIT press Cambridge, 2016.
- 395 Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks  
396 with pruning, trained quantization and huffman coding. *arXiv preprint arXiv:1510.00149*, 2015.
- 397 Chaoyang He, Haishan Ye, Li Shen, and Tong Zhang. Milenas: Efficient neural architecture search  
398 via mixed-level reformulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision  
399 and Pattern Recognition*, pp. 11993–12002, 2020.
- 400 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
401 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
402 2016.
- 403 Ramtin Hosseini and Pengtao Xie. Saliency-aware neural architecture search. *Advances in Neural  
404 Information Processing Systems*, 35:14743–14757, 2022.
- 405 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint  
406 arXiv:1412.6980*, 2014.
- 407 A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. *Master’s thesis,  
408 Department of Computer Science, University of Toronto*, 2009.
- 409 Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444,  
410 2015.
- 411 Ming Lin, Pichao Wang, Zhenhong Sun, Heseng Chen, Xiuyu Sun, Qi Qian, Hao Li, and Rong Jin.  
412 Zen-nas: A zero-shot nas for high-performance deep image recognition. In *2021 IEEE/CVF  
413 International Conference on Computer Vision, ICCV 2021*, 2021.

- 414 Shaohui Lin, Rongrong Ji, Yuchao Li, Cheng Deng, and Xuelong Li. Toward compact convnets via  
415 structure-sparsity regularized filter pruning. *IEEE transactions on neural networks and learning*  
416 *systems*, 31(2):574–588, 2019.
- 417 Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. *arXiv*  
418 *preprint arXiv:1806.09055*, 2018.
- 419 Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. Shufflenet v2: Practical guidelines for  
420 efficient cnn architecture design. In *Proceedings of the European conference on computer vision*  
421 *(ECCV)*, pp. 116–131, 2018.
- 422 Joe Mellor, Jack Turner, Amos Storkey, and Elliot J Crowley. Neural architecture search without  
423 training. In *International Conference on Machine Learning*, pp. 7588–7598. PMLR, 2021.
- 424 Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading  
425 digits in natural images with unsupervised feature learning. 2011.
- 426 Hieu Pham, Melody Guan, Barret Zoph, Quoc Le, and Jeff Dean. Efficient neural architecture search  
427 via parameters sharing. In *International conference on machine learning*, pp. 4095–4104. PMLR,  
428 2018.
- 429 Esteban Real, Alok Aggarwal, Yanping Huang, and Quoc V Le. Regularized evolution for image  
430 classifier architecture search. In *Proceedings of the aaai conference on artificial intelligence*,  
431 volume 33, pp. 4780–4789, 2019.
- 432 Laura Rieger, Chandan Singh, William Murdoch, and Bin Yu. Interpretations are useful: penalizing  
433 explanations to align neural networks with prior knowledge. In *International conference on*  
434 *machine learning*, pp. 8116–8126. PMLR, 2020.
- 435 Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical  
436 image segmentation. In *International Conference on Medical image computing and computer-*  
437 *assisted intervention*, pp. 234–241. Springer, 2015.
- 438 Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Du-  
439 mitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In  
440 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1–9, 2015.
- 441 Mingxing Tan, Bo Chen, Ruoming Pang, Vijay Vasudevan, Mark Sandler, Andrew Howard, and  
442 Quoc V Le. Mnasnet: Platform-aware neural architecture search for mobile. In *Proceedings of the*  
443 *IEEE/CVF conference on computer vision and pattern recognition*, pp. 2820–2828, 2019.
- 444 Hidenori Tanaka, Daniel Kunin, Daniel L Yamins, and Surya Ganguli. Pruning neural networks  
445 without any data by iteratively conserving synaptic flow. *Advances in neural information processing*  
446 *systems*, 33:6377–6389, 2020.
- 447 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz  
448 Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio,  
449 H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information*  
450 *Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- 451 Yu Weng, Tianbao Zhou, Yujie Li, and Xiaoyu Qiu. Nas-unet: Neural architecture search for medical  
452 image segmentation. *IEEE access*, 7:44247–44257, 2019.
- 453 Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking  
454 machine learning algorithms, 2017.
- 455 Lin Xiao and Tong Zhang. A proximal stochastic gradient method with progressive variance reduction.  
456 *SIAM Journal on Optimization*, 24(4):2057–2075, 2014.
- 457 Yuhui Xu, Lingxi Xie, Xiaopeng Zhang, Xin Chen, Guo-Jun Qi, Qi Tian, and Hongkai Xiong.  
458 Pc-darts: Partial channel connections for memory-efficient architecture search. *arXiv preprint*  
459 *arXiv:1907.05737*, 2019.

- 460 Yibo Yang, Hongyang Li, Shan You, Fei Wang, Chen Qian, and Zhouchen Lin. Ista-nas: Efficient and  
461 consistent neural architecture search by sparse coding. *Advances in Neural Information Processing*  
462 *Systems*, 33:10503–10513, 2020.
- 463 Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical  
464 risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.
- 465 Peng Zhao, Guilherme Rocha, and Bin Yu. The composite absolute penalties family for grouped and  
466 hierarchical variable selection. 2009.
- 467 Pan Zhou, Caiming Xiong, Richard Socher, and Steven Chu Hong Hoi. Theory-inspired path-  
468 regularized differential network architecture search. *Advances in Neural Information Processing*  
469 *Systems*, 33:8296–8307, 2020.
- 470 Barret Zoph and Quoc V Le. Neural architecture search with reinforcement learning. *arXiv preprint*  
471 *arXiv:1611.01578*, 2016.
- 472 Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. Learning transferable architectures  
473 for scalable image recognition. In *Proceedings of the IEEE conference on computer vision and*  
474 *pattern recognition*, pp. 8697–8710, 2018.