

---

# Sparse Structure Search for Delta Tuning

---

Shengding Hu<sup>1\*</sup>, Zhen Zhang<sup>1\*</sup>, Ning Ding<sup>1</sup>, Yadao Wang<sup>3</sup>,  
Yasheng Wang<sup>3</sup>, Zhiyuan Liu<sup>1,2,4†</sup>, Maosong Sun<sup>1,2,4</sup>

<sup>1</sup>Dept. of Comp. Sci. & Tech., Institute for AI, Tsinghua University, Beijing, China  
Beijing National Research Center for Information Science and Technology

<sup>2</sup>Institute Guo Qiang, Tsinghua University, Beijing, China <sup>3</sup>Noah’s Ark Lab, Huawei

<sup>4</sup>International Innovation Center of Tsinghua University, Shanghai, China

{hsd20, zhen-zha19}@mails.tsinghua.edu.cn

## Abstract

Adapting large pre-trained models (PTMs) through fine-tuning imposes prohibitive computational and storage burdens. Recent studies of delta tuning (DT), i.e., parameter-efficient tuning, find that only optimizing a small portion of parameters conditioned on PTMs could yield on-par performance compared to conventional fine-tuning. Generally, DT methods exquisitely design delta modules (DT modules) which could be applied to arbitrary fine-grained positions inside PTMs. However, the effectiveness of these fine-grained positions largely relies on sophisticated manual designation, thereby usually producing sub-optimal results. In contrast to the manual designation, we explore constructing DT modules in an automatic manner. We automatically Search for the Sparse Structure of **Delta** Tuning (S<sup>3</sup>Delta). Based on a unified framework of various DT methods, S<sup>3</sup>Delta conducts the differentiable DT structure search through bi-level optimization and proposes shifted global sigmoid method to explicitly control the number of trainable parameters. Extensive experiments show that S<sup>3</sup>Delta surpasses manual and random structures with less trainable parameters. The searched structures preserve more than 99% fine-tuning performance with 0.01% trainable parameters. Moreover, the advantage of S<sup>3</sup>Delta is amplified with extremely low trainable parameters budgets (0.0009%~0.01%). The searched structures are transferable and explainable, providing suggestions and guidance for the future design of DT methods. Our codes are publicly available at <https://github.com/thunlp/S3Delta>.

## 1 Introduction

Increasingly large pre-trained models (PTMs) [6, 27, 30, 31, 12] building upon Transformers [36] have been emerging and achieving state-of-the-art results on a variety of downstream tasks. Despite the blessing of effectiveness, these big models also bring the curse of prohibitive costs on computation and storage during the adaptation because of the gradient computation of the whole model and the giant size of the fine-tuned checkpoint.

To alleviate such costs, studies of delta tuning (DT) [7], also known as parameter-efficient tuning [15, 28, 42, 16, 25, 9, 20, 18], have been developed, which only train a small portion of PTMs and keep the vast majority of parameters frozen. Studies have verified that delta tuning could achieve competitive performance compared to conventional fine-tuning with very few trainable parameters, resulting in considerable savings in model adaptation costs. Generally, these approaches manually design delta modules (DT modules) to complete model adaptation. For example, adapter-based

---

\* Equal contribution, ordered alphabetically.

† Corresponding authors: Z.Liu (liuzy@tsinghua.edu.cn)

methods [15, 28, 25] inject two newly-introduced feed-forward layers to Transformers and only re-tune 0.5%-8% parameters to yield promising results; BitFit [42] only re-tunes the bias terms (0.04% - 0.1% parameters) within Transformers; LoRA [16] inserts trainable rank decomposition matrices to each layer of Transformers and is successfully adopted on GPT-4 with 175 billion parameters.

While early research focused on how to design practically effective DT modules, more recent research has advanced the understanding of delta tuning more deeply. He et al. [13] bridge connections among different approaches to form a unified framework. And Ding et al. [7] indicate that the combination of different trainable modules could bring different levels of gain on downstream tasks. The above empirical evidence implies that there may exist optimal mixture of DT modules that is more effective than manually designed structures. In fact, considering the fine-grained structure inside PTMs, the positions where the DT modules could be applied are numerous, but not all DT modules at all positions contribute equally to the task performance. How to find the optimal structure of DT modules and remove the redundancy in the trainable parameters is essential for a more efficient adaptation method. Predictably, such optimal structure is difficult to construct artificially and may vary with specific tasks and models. Therefore, we propose to automatically search the optimal structure that contains a mixture of DT modules at diverse positions inside PTMs. Also importantly, the structure should be sparse to ensure the parameter efficiency.

We present Sparse Structure Search for Delta Tuning ( $\mathcal{S}^2$ Delta) to automatically search such optimal trainable structure, which could flexibly control the number of trainable parameters according to practical requirements. The searching process and the optimization are guided by performance to ensure the effectiveness on specific tasks. Moreover, the structures change automatically to suit the preset limitation of the number of trainable parameters. In contrast, heuristically designed structures are usually coarse-grained and independent of performance and budget, making them neither optimal nor flexible to adjust the number of trainable parameters.

In terms of the specific methodology, we firstly construct a unified search space by applying probabilistic gating controlled by structural parameters to all potential DT modules. Then, we develop a framework of differentiable DT structure search by treating the problem as a constrained neural architecture search problem. In our framework, the structural parameters are updated via bi-level optimization [22]. Unlike the traditional neural architecture search that learns from scratch, we implement the first neural structure search based on a pre-defined backbone and under the delta tuning scenario. To search under a pre-defined budget of trainable parameters, we develop a global sigmoid to explicitly control the number of activated DT modules in the searching phase.

We conducted extensive experiments to study the effectiveness of  $\mathcal{S}^2$ Delta. Firstly, the experiments show that with 0.01% parameters, we are able to recover 99% and 98% re-tuning performance on GLUE [38] and SuperGLUE [7], respectively. Secondly, the searched structure surpasses the human-designed structures considerably while consuming less trainable parameters. Moreover, the advantage enlarges when the number of trainable parameters is minimal (0.0009%). Furthermore, the searched structures are transferable across tasks, which significantly strengthens the usefulness of the searched structures. Apart from the performance boost, we visualize and explain the searched structures, which is beneficial to the future design of new DT methods.

## 2 Related Work

Delta Tuning (DT). Our work is related to the studies of delta tuning (DT) for pre-trained models [7]. Generally, DT only optimizes a small portion of parameters and leaves the vast majority of parameters untouched for the adaptation to downstream tasks. Pioneer work selects parts of the PTMs to be trainable [35, 11, 10]. Adapter [15] is one of the earliest methods that apply the concept of parameter-wise efficiency to pre-trained language models, which inserts linear neural modules to every Transformer layer and achieves on-par results to full re-tuning. As the PTMs scaling in recent years, DT is valued for its efficiency in computing and storage. This has spawned not only empirical studies [28, 14] and variants on the adapter [28, 25, 34], but also a range of other approaches. Prefix tuning [20] prepends embeddings to the hidden states of the Transformer model, and prompt tuning [18] further simplifies the strategy and only prepends such embeddings to the input layer. There are also approaches which specify some of the parameters inside PTMs that can be trainable to achieve good results, such as Mask4it [43], BitFit [42], DiffPruning [9], etc. LoRA [16] assumes that the change in model weights is intrinsically low-rank after re-tuning, and uses trainable rank-decomposition matrices for model adaptation. In addition to specific methods, some

studies have comprehensively investigated DT methods. He [13] models multiple methods in a unified manner, Ding et al [7] provides a theoretical discussion and comprehensive empirical study of these methods. Our work proposes to automatically search for trainable structures in the context of DT, which is a different perspective from all the aforementioned work. In terms of the structure of DT, AdapterDrop [3] explores dropping a fraction of Adapter modules based on manual trials. A concurrent work [26] learns switches on Adapter modules to select the beneficial adapter modules. However, it is not optimized under a preset number of trainable parameters. They are also both limited to adapter-based methods. On the contrary, our proposed method can search within a mixture of almost all DT modules under a constrained trainable parameter budget.

Neural Architecture Search (NAS). Our work conducts structure search in the scope of DT, which is related to the Neural Architecture Search algorithms. A line of NAS algorithms uses Reinforcement Learning or Evolutionary Algorithms to explore the best structure with reward from training the structure from scratch [4, 45, 32, 29], which usually consumes prohibitive computation resources. Another line of NAS algorithms [21, 5] approaches the problem with gradient-based optimization. DARTS [2] relaxes the discrete structure using continuous structural parameters, which are optimized with gradient-based optimization. DARTS achieves competitive performances with much fewer computational resources. We take inspirations from DARTS in optimizing the structural parameters of  $\Delta$ . We are also the first to conduct NAS conditioned on a pre-trained backbone model. We also take inspiration from the NAS algorithms with binary gates [3, 40].

### 3 Method

In this section, we firstly introduce the preliminaries of pre-trained model adaptation, transformer architecture, and the delta tuning. Then we introduce our method  $\Delta$  in detail.

#### 3.1 Preliminaries

Pretrained Model Adaptation. The recent prevalent pre-train then fine-tune paradigm in deep learning takes advantage of a pre-trained model with parameters  $\theta$  and continues to optimize on a downstream task  $\mathcal{D} = \{D_{train}, D_{val}, D_{test}\}$  under an objective function. In fine-tuning, all the parameters of the pre-trained model are optimized using the train split to minimize,

$$\min_{\theta} L(M(\theta); D_{train}) \quad (1)$$

Transformer Architecture. The pre-trained models typically adopt the Transformer model [36] as their backbone. The Transformer model is composed of multiple stacked Transformer layers that processes the hidden state sequentially through different computation modules, such as Self-Attention module (SelfAttn), Cross-Attention module (CrossAttn), Feed-Forward module (FFN), and Layer Normalization module (LN), etc., and details of each module are in Appendix A. The computation process in the Transformer can be abstracted by a sequence of transformations of the hidden representation. In each computation step, the input hidden representation  $H^{in} \in \mathbb{R}^{s \times d_1}$  is transformed into an output hidden representation  $H^{out} \in \mathbb{R}^{s \times d_2}$ , where  $s$  is the sequence length of the input and  $d_1, d_2$  are the hidden dimensions,

$$H^{out} = m(H^{in}) \quad (2)$$

Delta Tuning (DT). DT methods only train a small portion of parameters conditioned on the backbone PTMs to improve the adaptation efficiency [5, 28, 42, 16, 25, 9, 20, 18]. Although the specific forms of the various DT modules are substantially different, He [13] unify them as modifications of the hidden state,

$$H^{out} = m(H^{in}) + \Delta \quad (3)$$

The formulas of some DT methods under the unified view are listed in Table 1. The DT modules can be applied to extensive positions on the backbone PTMs, which are listed in the rightmost column in Table 1. In training, we freeze all the parameters in the backbone module,  $m$ , and set

<sup>3</sup>We use a little bit more flexible notation than [3], which takes into account the frozen backbone module  $m$  and thus can distinguish the DT modules that take either  $H^{out}$  or  $H^{in}$  as their input.

Table 1: Different DT methods are the specializations of the uni ed view (Equation (3)) and can be applied to extensive positions on the PTMs.

Method	Transformation		Potential Positions
LoRA [16]	$H^{out} = H^{in}(W + AB)$	$H_0AB$	Weight matrices
Adapter [15]	$H^{out} = m(H^{in}) + f(m(H^{in})W_{down})W_{up}$	$f(m(H^{in})W_{down})W_{up}$	After any modules
Parallel Adapter [13]	$H^{out} = m(H^{in}) + f(H^{in}W_{down})W_{up}$	$f(H^{in}W_{down})W_{up}$	Between any two modules
BitFit [42]	$H^{out} = m(H^{in}) + b$	$b$	Linear layers
LNFit <sup>4</sup>	$H^{out} = \frac{H^{in}}{\sqrt{\text{Var}(H^{in})}}(s + s') + b$	$\frac{H^{in}}{\sqrt{\text{Var}(H^{in})}}s$	Layer Normalization modules

Figure 1: The framework of  $\Delta^3$ . We propose a uni ed search space with probabilistic gating to enable search among a mixture of DT methods. We find the optimal sparse structure using the differentiable DT structure search and explicit sparsity control.

parameters introduced in computing denoted by  $\theta$ , as the only trainable parameters. Therefore, the adaptation objective in DT is

$$\min_{\theta} L(M(\theta; \mathcal{D}_{train})); \quad (4)$$

For the convenience of notation, we simplify Equation (4) into

$$\min_{\theta} L_{train}(\theta); \quad (5)$$

### 3.2 Sparse Structure Search for Delta Tuning

Our target is to search for the optimal structure of DT constrained by a pre-defined and limited trainable parameter budget. To achieve this, we design Sparse Structure Search for Delta Tuning ( $\Delta^3$ ), which is driven by three essential components: uni ed search space with probabilistic gating, an efficient differentiable DT structure search algorithm, and an explicit sparsity control algorithm using shifted global sigmoid.

**Uni ed Search Space with Probabilistic Gating.** The potential positions in the backbone models to which DT modules could be applied are extensive, especially when we consider a mixture of different types of DT modules through the uni ed view (Equation (3)). However, not all positions contribute to the task performance equally, and only a fraction of positions should be activated to avoid redundancy of the trainable parameters. To this end, we design a probabilistic gating mechanism over all possible DT modules' positions. Specifically, for each DT module that computes  $H_i$  for the hidden representation, we activate the modification  $H_i$  with probability  $p_i \in [0, 1]$ ,

$$H_i^{out} = m(H_i^{in}) + z_i \cdot i; \quad (6)$$

where  $z_i \sim \text{Bernoulli}(p_i)$  is a random sample from Bernoulli distribution.

**Differentiable DT Structure Search.** Finding the optimal structure from a search space with abundant potential positions is challenging due to the compositionality of activated positions. Furthermore,

<sup>4</sup>LNFit only trains the variance vector in the Layer Normalization module of the PTMs, which is inspired by Frankle et al. [8] who only train the Batch Normalization module in Convolutional Neural Networks.

directly comparing the fully trained model with each DT structure is intractable. In our work, we propose to optimize the gating probability  $p_i$  with gradient-based optimization. To make the sampling process differentiable, we use the Binary Concrete Distribution [17] as a soft and differentiable approximation to the Bernoulli distribution,

$$z_i = \frac{1}{\tau} \log \frac{u p_i}{(1-u)(1-p_i)}; \quad (7)$$

where  $u \sim U(0, 1)$  is a random sample from the uniform distribution  $[0, 1]$  and  $\tau$  is the temperature to control the sharpness of the distribution<sup>5</sup>. Similar distributions are used in learning sparse networks [23] or pruning dense networks [9]. By replacing the hard sample with the soft approximation  $z_i$ , we can back-propagate through training. However, directly optimizing  $p_i$  in the probability space  $[0, 1]$  may lead to numerical instability, therefore, we parameterize it with structural parameters  $s_i \in \mathbb{R}$ , i.e.,  $p_i = g(s_i)$ . And we denote all the structure parameters as

We optimize  $s$  through bi-level optimization [1, 22], i.e., optimizing  $s$  conditioned on the optimized parameters  $\theta$  of the DT modules. The inner and outer level of optimization are conducted on separate splits of training data, denoted  $D_1; D_2$ , which is analogous to validating structures trained on  $D_1$  using a different split  $D_2$  to avoid over-fitting to  $D_1$ . Thus, the optimization objective is

$$\min_s L(M(\theta; s)); \quad (8)$$

$$s: t. \theta = \operatorname{argmin}_{\theta} L(M(\theta; s)); \quad (9)$$

Following DARTS [22], we make approximations to the gradient of the structural parameters by applying chain rule and taking finite difference approximations

$$r_{s_i} L(\theta; s) \quad (10)$$

$$r_{s_i} L(\theta; s) = r_{s_i} L(\theta; s); \quad (11)$$

$$r_{s_i} L(\theta; s) \approx r_{s_i}^2 L(\theta; s) + r_{s_i} L(\theta; s) \quad (12)$$

$$r_{s_i} L(\theta; s) \approx \frac{r_{s_i} L(\theta; s^+) + r_{s_i} L(\theta; s^-)}{2}; \quad (13)$$

where the optimal  $\theta$  is approximated by the parameters of one-step update  $r_{s_i} L(\theta; s)$ .  $\alpha$  is the learning rate of parameters and  $\epsilon$  is a small scalar used in the finite difference approximation.

Explicit Sparsity Control with Shifted Global Sigmoid. Most of the DT modules in the search space are redundant and contribute little to the performance. However, the search algorithm may not be aware of the sparsity target and degenerate to greedily adding more DT modules. As opposed to previous sparse network learning methods [23, 9] which punish the dense structures with regularization, we explicitly control the sparsity of structure at the target level during the search through a shifted global sigmoid parameterization. (See Section 4.9 for comparing the two methods),

$$p_i = p_i \frac{\operatorname{Detach}(p_i)}{p_i}; \quad (14)$$

$$\text{where } p_i = \operatorname{Sigmoid}\left(\frac{s_i}{\tau}\right); \quad (15)$$

The  $\operatorname{Detach}(\cdot)$  operator turns a parameter that requires gradient into a scalar that is free from gradient computation. Equation (14) doesn't change the value of  $p_i$ , but it enforces the competition among different positions and DT modules, which is similar to Softmax operation (See Appendix C.2 for details).

<sup>5</sup>The distribution of  $z_i$  has the property that  $\mathbb{P}(z_i > 0.5) = p_i$ , and when  $\tau$  approaches 0, the distribution of  $z_i$  converges to  $\mathcal{B}(1; p_i)$  (See Appendix C.1), which makes it a suitable surrogate for Bernoulli distribution.

<sup>6</sup>We use the same sample to compute  $r_{s_i} L(\theta; s^+)$ ,  $r_{s_i} L(\theta; s^-)$ ,  $r_{s_i} L(\theta; s)$ , and  $r_{s_i} L(\theta; s)$ .

In Equation(15),  $p_i$  is a scalar. Increasing  $p_i$  value will monotonically reduce  $p_i$  to 0 while keeping  $p_i$  in  $[0; 1]$ . So the expected number of trainable parameters  $E[N]$  is a monotonic function w.r.t.,

$$E[N] = E \sum_i |z_i = 1|_{ij} = E \sum_i |z_i > 0.5|_{ij} = \sum_i p_{ij} |j|; \quad (16)$$

where the  $|j|$  is the number of parameters introduced in computing  $z_i$ . Thus, we can dynamically adjust  $p_i$  to make  $E[N]$  approach  $B$  via monotonic optimization,

$$p_i = \operatorname{argmin} \left( \sum_i p_{ij} |j| B \right); \text{ where } \sum_i p_{ij} |j| = B; \quad (17)$$

**Evaluation of the Searched Structure.** To determine the final structure of DT, instead of sampling from  $p_i$ , we choose the set of positions where the sum is the highest while still being within the budget  $B$ . This deterministic algorithm reduces the variance of the final structures. After obtaining the final structure, we re-initialize and re-train the parameters in the DT modules to converge on

---

#### Algorithm 1 Algorithm of $\mathcal{S}^3\Delta$

---

Initialize all DT modules in the search space, and initialize  
while not converged  
do  
1. Calculate  $p_i$ , and sample  $z_i$ .  
2. Compute the each loss terms by forward and backward propagation.  
3. Update  $p_i$  according to Equation (13).  
4. Update using  $L(\theta; z_i)$ .  
end while  
Determine and evaluate the final structure.

---

## 4 Experiments

### 4.1 Datasets and PTMs

We apply  $\mathcal{S}^3\Delta$  to multitask benchmarks GLUE [38] and SuperGLUE [37] following previous works. All datasets are downloaded from the HuggingFace Datasets. Since the test splits of these datasets are held out and invisible to the researchers, we conduct random splits from either train set or validation set to make the new train, validation, and test splits, which is critical to ensure fair evaluations according to Chen et al. [41]. We repeat 4 times using different random seeds for experiments in Table 2, and 8 times for experiments in Figure 2. The details are in Appendix B. We use the  $T_{\text{large}}$  model (703M parameters) as the backbone PTMs.

### 4.2 Baselines

We compare  $\mathcal{S}^3\Delta$  with several widely used baselines (See Appendix B for details).

**Fine-tune.** Traditional fine-tuning trains all parameters in the PTMs.

**LoRA.** We apply LoRA linear layer to the Self-Attention’s query modules and value modules as Hu et al. [16] suggest. We include two rank levels (8 and  $r = 1$ ) in our experiments.

**Adapter.** We adopt the first adapter method proposed by Houlsby et al. [15]. Their method requires more parameters than the other methods but achieves good empirical results.

**Low Rank Adapter (Adapter-LR).** We adopt the Low Rank Adapter as an efficient variant of the adapter-based method. It is proposed in [25] as a simple but effective baseline. The rank is set to 1.

**BitFit.** BitFit proposes to only adapt the bias layer in the model. We adopt the same setting as Zaken et al. [42] that tunes the bias inside all linear modules, and the Layer Normalization layer

**LNFit.** We train the variance vector of all Layer Normalization layers’, including the Layer Normalization after the whole transformer encoder.

<sup>7</sup>Although T5 has no bias in linear modules, we can treat it as bias vectors with zero initialization.

Table 2: Results on GLUE [38] benchmark (above) and SuperGLUE [7] benchmark (below).

Green and blue represent the best and second best scores, respectively, among the methods in our search space. The first three rows represent the results of ne-tuning and other DT methods, which are not used in our search space due to high trainable parameter ratios. On SuperGLUE tasks, since the results on COPA vary dramatically (±6.00), the average results of SuperGLUE become easily dominated by the results on COPA. Therefore we also report the average results that exclude COPA (AVG<sub>COPA</sub>). The widths of the yellow rectangles are proportional to the trainable parameter ratios.

GLUE																	
Parameter Ratios	Method	CoLA		SST2		MRPC		QQP		STS B		MNLI	QNLI	AVG			
10000%%	Fine-tune	62.25	3.96	95.87	0.42	91.86	1.19	89.50	0.22	91.86	0.46	89.61	0.30	94.22	0.35	87.88	
65.33%%	Adapter	59.03	3.06	95.90	0.29	93.02	0.28	88.39	0.06	91.77	0.25	89.53	0.07	94.17	0.19	87.40	
21.32%%	LoRA(r=8)	58.43	4.16	95.79	0.27	92.21	0.88	88.35	0.25	91.78	0.31	89.38	0.32	94.14	0.12	87.15	
Methods in the Search Space																	
8.13%%	BitFit	56.98	3.89	96.24	0.33	92.16	0.68	88.12	0.07	91.59	0.08	89.10	0.09	94.07	0.21	86.90	
4.12%%	Adapter-LR	56.78	4.80	95.90	0.14	92.76	0.67	88.08	0.13	91.26	0.31	89.30	0.14	93.94	0.07	86.86	
2.67%%	LoRA(r=1)	56.77	2.29	95.81	0.27	92.45	1.00	88.08	0.11	91.54	0.33	89.16	0.17	94.10	0.05	86.84	
1.70%%	LNFit	56.15	4.06	95.81	0.20	91.71	0.39	88.17	0.10	91.37	0.24	89.11	0.09	93.99	0.20	86.62	
1.39%%	S <sup>3</sup> Delta-M	59.34	4.75	95.84	0.14	92.13	2.09	88.04	0.23	91.58	0.25	89.14	0.13	94.12	0.12	87.17	
1.39%%	S <sup>3</sup> Delta-L	56.71	3.03	95.93	0.15	93.27	1.39	88.14	0.08	91.58	0.49	88.81	0.44	93.95	0.11	86.91	
0.35%%	S <sup>3</sup> Delta-M	54.56	3.66	95.93	0.24	92.14	1.10	88.02	0.20	91.38	0.34	89.04	0.25	93.93	0.14	86.43	
SuperGLUE																	
Parameter Ratios	Method	BoolQ	CB	COPA		MultiRC	ReCORD	RTE	WIC	AVG	AVG <sub>COPA</sub>						
10000%%	Fine-tune	86.67	0.21	96.43	2.92	73.50	5.26	76.65	1.01	85.03	0.67	88.49	2.12	73.12	1.71	82.84	84.40
65.33%%	Adapter	85.98	0.68	94.64	6.19	63.00	7.75	77.60	0.84	85.96	0.37	89.21	2.94	71.63	0.90	81.15	84.17
21.32%%	LoRA(r=8)	85.06	0.70	91.96	3.42	51.00	4.16	76.94	1.16	85.84	0.21	87.05	0.59	72.10	1.31	78.56	83.16
Methods in the Search Space																	
8.13%%	BitFit	85.02	0.48	89.29	2.92	75.00	8.08	75.79	1.15	85.85	0.32	86.15	1.48	72.34	1.61	81.35	82.41
4.12%%	Adapter-LR	84.53	0.37	84.82	8.44	49.50	6.81	76.67	1.37	86.04	0.09	85.61	2.42	71.39	0.70	76.94	81.51
2.67%%	LoRA(r=1)	85.60	0.45	84.82	1.79	67.50	5.00	76.71	1.05	85.95	0.36	86.87	1.08	71.32	1.29	79.82	81.88
1.70%%	LNFit	84.07	0.50	82.14	2.92	49.00	1.15	75.52	1.16	86.14	0.11	86.69	1.81	69.28	1.49	76.12	80.64
1.39%%	S <sup>3</sup> Delta-M	84.92	0.68	92.86	2.92	70.50	3.79	76.38	0.92	86.10	0.11	86.69	1.90	71.63	1.07	81.30	83.10
1.39%%	S <sup>3</sup> Delta-L	85.00	0.67	90.18	6.10	60.00	9.52	76.17	1.41	86.02	0.15	85.79	1.36	71.63	1.39	79.26	82.46
0.35%%	S <sup>3</sup> Delta-M	83.56	0.53	87.50	4.61	54.00	4.32	76.09	0.97	86.10	0.26	85.79	1.89	68.42	1.89	77.35	81.24

We do not include Prompt Tuning [8] as our baseline because it takes much longer steps to converge and doesn't achieve competitive performance on T<sub>Large</sub> [18].

### 4.3 S<sup>3</sup>Delta Search Spaces and Budgets

The search space has a noteworthy influence on the performance of S<sup>3</sup>Delta. In our experiment, we define two kinds of search space.

**Mix.** The first search space considers a mixture of LoRA, Adapter-LR, BitFit, and LNFit modules. LoRA can be applied to any linear modules in the transformer block, including the query(Q), key(K), value(V), output(O) sub-modules of the attention module(ATTN), and the two sub-layer W1 and W2 in Feed Forward modules(FFN). The Adapter-LR can theoretically be applied to any position in the computational graph. However, to avoid overcomplicating the search space, we limit it to the outputs of the ATTN module and the FFN modules. For BitFit, the potential applied positions are all the linear modules (Q, K, V, O, W1, W2) and the Layer Normalization modules (LN). For LNFit, the potential positions are all the LN modules. In this search space, the potential DT modules to be selected and the total number of candidate structures is 2816. We do not consider the budget constraint. We denote the structures searched on this search space as S<sup>3</sup>Delta-M.

**LoRA.** We narrow down the search space into a single type of DT module. We choose LoRA as an example. The potential positions are the same as the LoRA modules in Mix search space. There are 288 potential positions in total. We denote the structures searched on this search space as S<sup>3</sup>Delta-L.

We also explore different numbers of trainable parameters. Experiments in Table 2 are conducted on 1.39%% and 0.35%% trainable parameters ratios. More sparsity levels are tested in section 4.5.

### 4.4 Results on GLUE and SuperGLUE

Table 2 shows the performance of different methods on GLUE and SuperGLUE tasks. Comparing S<sup>3</sup>Delta with the manual structures within the search space, we find that S<sup>3</sup>Delta-M (1.39%%) achieves the highest average score on GLUE and SuperGLUE (without COPA) despite using the least number of trainable parameters (1=5 compared to BitFit). S<sup>3</sup>Delta-M (0.35%%) also

(a) RTE (b) MultiRC (c) MRPC

Figure 2: Performances under different trainable parameters ratios. The x-axis represents the ratio of the number of trainable parameters to the backbone PTM's parameters. The scaled y-axis represents the scores. The accuracy of fine-tuning is in gray horizontal line. The result of LoRA (and Low Rank Adapter) are plotted in grey dot.

surpasses Adapter-LR, LoRA (1), LNFit using approximately 1/12, 1/8, 1/5 trainable parameters, respectively. Narrowing the search space from Mix to LoRA leads to a moderate decrease in performance, which justifies the need to search among a mixture of DT modules. It may also hint that the combination of different DT modules could lead to stronger performance. However, even though the performance of Delta-L is not optimal, it compares favorably to LoRA (1), which also shows that the human designed structures, though benefit from applying DT modules uniformly on the PTMs, is sub-optimal. Compared with fine-tuning, Delta-M (1.39%) preserves 99.2% and 98.1% performance on GLUE and SuperGLUE, respectively. In fact, we must emphasize that Delta is orthogonal to specific DT modules. The performance of Delta can benefit from the future invention of better DT modules, thus potentially achieving comparable or even superior performance to fine-tuning with extremely limited trainable parameters.

#### 4.5 Performance under Different Sparsity Levels

To explore the limit of trainable parameter reduction, we train different methods with decreasing sparsity levels from 5.6% to 0.086%. To apply baseline methods on target numbers of trainable parameters, we randomly sample a set of potential positions in their corresponding search space to reach the target sparsity level. In Figure 2, we demonstrate the results on three datasets, RTE, MultiRC and MRPC. We can see that Delta-M and Delta-L have considerable advantages in extremely low trainable parameter budgets. For example, Delta-M trains only 0.086% parameters whereas recovering 96.8%, 98.7%, 97.5% of the FT performances on RTE, MultiRC, MRPC, respectively. With 5.6% trainable parameters on MultiRC and MRPC, all the methods saturate to FT performance, proving the feasibility of removing redundant parameters.

#### 4.6 Transferability of the Searched Structures

Another essential characteristic of Delta is the transferability of the searched structure. In Table 4, we split the GLUE benchmark into source datasets and target datasets. We search on the source dataset (Mix search space) and train the searched structure on the target datasets. We can see that the searched structures are highly transferable, even surpassing the structures directly searched on the target datasets sometimes. The transferability guarantees the reusability of the searched structures.

#### 4.7 Efficiency of the Search Process

Although Delta focuses on the parameter-efficiency of the searched structures, we also analyze the searching efficiency in Table 3. Generally speaking, the search for an optimal structure consumes 5-8 times training time and 2 times GPU memory (Due to bi-level optimization). However, it is affordable compared to manually designing different structures and running numerous evaluations.

#### 4.8 Visualization and Explanations of the Search Structures

To understand the searched structures, we draw the heat maps of the different datasets in Appendix D.2. We find obvious patterns and similarities in most datasets. Therefore, we average the heat maps across datasets to see the overall pattern of the searched structure. Figure 3 shows the heatmap of



Table 3: The computational resources in the searching phase and re-training phase, Computation time, memory consumption are listed.

Dataset	Search	Re-train	Ratio
Time/min			
RTE	148.6	30.0	5.0
STSB	139.3	26.3	5.3
CoLA	145.6	17.0	8.6
Memory/GB			
RTE	27.7	16.6	1.7
STSB	28.9	10.6	2.7
CoLA	16.1	8.9	1.8

Figure 3: Visualization of  $\phi_i$  of  $S^3$ Delta-M. The numbers on the squares are the average  $\phi_i$  across all datasets and all seeds. The deeper the color is, the more activated is the DT module. The x-axis represents different layers of PTMs (E denotes Encoder, and D denotes Decoder), and the y-axis represents different positions (modules) of PTMs.

Figure 5: Visualization of  $\phi_i$  of DT modules in  $S^3$ Delta-L.

Table 4: Structure transfer from source datasets and target datasets. The target datasets are in the row name, and the source datasets are in the column names. "No transfer" means the structure is searched on the target dataset.

Source	Target Datasets							
	STSB		QQP		QNLI		CoLA	
No transfer	91.58	0.25	88.03	0.23	94.11	0.12	59.34	4.75
MRPC	91.63	0.31	88.16	0.08	93.96	0.06	56.41	3.81
MNLI	91.39	0.67	88.06	0.08	94.13	0.10	56.38	3.98
SST2	91.37	0.23	88.02	0.10	94.14	0.16	55.58	4.13

Figure 4: Comparing shifted global sigmoid to  $L_0$  regularization. Performance on different datasets is in different colors. Shifted global sigmoid and  $L_0$  regularization are in solid lines and dotted lines, respectively.

$S^3$ Delta-M. We can see that (1) The BitFit modules in the Self-Attention modules and Cross-Attention modules in the higher decoder layers are highly preferred, proving that the BitFit modules are simple and effective. This observation is also beyond the intuition of human experts, as most previous work ignores the contribution of training or applying DT methods to Cross-Attention modules; (2) The last layers of the encoder and decoder are emphasized, which is close to the traditional use of PTMs as feature extractors by training only the last layer; (3) We also observe that BitFit modules tend to be distributed approximately evenly across the higher layers (See Appendix D.2 for details). Figure 5 shows the  $\phi_i$  of  $S^3$ Delta-L. The trend of choosing higher layers still exists. Interestingly, the query sub-modules are prioritized in the encoder, while the value sub-modules are stressed in the decoder.

#### 4.9 Ablation Study

To explicitly control sparsity, we propose shifted global sigmoid, which differs from the  $L_0$  regularization used in previous work [23]. We compare the results of regularization using shifted global sigmoid and  $L_0$  on three datasets. From Figure 4, it is clear that shifted global sigmoid has an advantage over  $L_0$  regularization at almost all sparsity, and the advantage increases with increasing sparsity.

## 5 Conclusion

In this paper, we propose Sparse Structure Search for Delta Tuning ( $S^3$ Delta), which conducts differentiable DT structure search with explicit sparsity control in a unified search space of a mixture of various DT modules. Experiments demonstrate the effectiveness of  $S^3$ Delta to find the optimal structure of DT modules and push the limit of trainable parameter reduction. For future works, there are open questions that are worth investigating. (1) Better search spaces or better DT modules could be designed to further explore the potential of structure search. (2) The current NAS algorithms are not tailored for the scenario where a pre-trained backbone model exists. Therefore, more specialized search algorithms could be developed for DT structure search.

## 6 Acknowledgements

This work is supported by National Key R&D Program of China (No. 2020AAA0106502), Institute Guo Qiang at Tsinghua University, Beijing Academy of Artificial Intelligence (BAAI), International Innovation Center of Tsinghua University, Shanghai, China.

Shengding Hu proposed the idea and framework. Shengding Hu and Zhen Zhang designed the methods and experiments. Zhen Zhang conducted the experiments. Shengding Hu and Ning Ding wrote the paper. Zhiyuan Liu and Maosong Sun advised the project and participated in the discussion. Yadao Wang and Yasheng Wang participated in the discussion and provided computational resources.

## References

- [1] G Anandalingam and Terry L Friesz. Hierarchical optimization: An introduction. *Annals of Operations Research*, 1992.
- [2] Tom B. Brown, Benjamin Mann, Nick Ryder, 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. In *Proceedings of NeurIPS*, 2020.
- [3] Han Cai, Ligeng Zhu, and Song Han. Proxylessnas: Direct neural architecture search on target task and hardware. In *Proceedings of ICLR*, 2019.
- [4] Guanzheng Chen, Fangyu Liu, Zaiqiao Meng, and Shangsong Liang. Revisiting parameter-efficient tuning: Are we really there yet? *ArXiv preprint*, 2022.
- [5] Xin Chen, Lingxi Xie, Jun Wu, and Qi Tian. Progressive differentiable architecture search: Bridging the depth gap between search and evaluation. *Proceedings of ICCV*, 2019.
- [6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. *Proceedings of ACL*, 2019.
- [7] Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin Chen, Chi-Min Chan, Weize Chen, et al. Delta tuning: A comprehensive study of parameter efficient methods for pre-trained language models. *ArXiv preprint*, 2022.
- [8] Jonathan Frankle, David J. Schwab, and Ari S. Morcos. Training batchnorm and only batchnorm: On the expressive power of random features in cnns. *Proceedings of ICLR*, 2021.
- [9] Demi Guo, Alexander Rush, and Yoon Kim. Parameter-efficient transfer learning with diff pruning. In *Proceedings of ACL*, 2021.
- [10] Yunhui Guo, Honghui Shi, Abhishek Kumar, Kristen Grauman, Tajana Rosing, and Rogério Schmidt Feris. Spottune: Transfer learning through adaptive re-tuning. *Proceedings of CVPR*, 2019.
- [11] Yunhui Guo, Yandong Li, Liqiang Wang, and Tajana Rosing. AdaIter: Adaptive Iter re-tuning for deep transfer learning. *Proceedings of AAAI*, 2020.
- [12] Xu Han, Zhengyan Zhang, Ning Ding, Yuxian Gu, Xiao Liu, Yuqi Huo, Jiezhong Qiu, Liang Zhang, Wentao Han, Minlie Huang, Qin Jin, Yanyan Lan, Yang Liu, Zhiyuan Liu, Zhiwu Lu, Xipeng Qiu, Ruihua Song, Jie Tang, Ji-Rong Wen, Jinhui Yuan, Wayne Xin Zhao, and Jun Zhu. Pre-trained models: Past, present and future. *Open*, 2021.
- [13] Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. Towards a unified view of parameter-efficient transfer learning. *Proceedings of ICLR*, 2022.
- [14] Ruidan He, Linlin Liu, Hai Ye, Qingyu Tan, Bosheng Ding, Liying Cheng, Jiawei Low, Lidong Bing, and Luo Si. On the effectiveness of adapter-based tuning for pretrained language model adaptation. In *Proceedings of ACL*, 2021.

- [15] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for NLP. In Proceedings of ICMJ2019.
- [16] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. 2021.
- [17] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. In Proceedings of ICLR2017.
- [18] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. In Proceedings of EMNLP2021.
- [19] Quentin Lhoest, Albert Villanova del Moral, Yacine Jernite, Abhishek Thakur, Patrick von Platen, Suraj Patil, Julien Chaumond, Mariama Drame, Julien Plu, Lewis Tunstall, Joe Davison, Mario Šaško, Gunjan Chhablani, Bhavitvya Malik, Simon Brandeis, Teven Le Scao, Victor Sanh, Canwen Xu, Nicolas Patry, Angelina McMillan-Major, Philipp Schmid, Sylvain Gugger, Clément Delangue, Théo Matussière, Lysandre Debut, Stas Bekman, Pierric Cistac, Thibault Goehringer, Victor Mustar, François Lagunas, Alexander Rush, and Thomas Wolf. Datasets: A community library for natural language processing. Proceedings of EMNLP2021.
- [20] Xiang Lisa Li and Percy Liang. Pre x-tuning: Optimizing continuous prompts for generation. In Proceedings of ACL2021.
- [21] Hanwen Liang, Shifeng Zhang, Jiacheng Sun, Xingqiu He, Weiran Huang, Kechen Zhuang, and Zhenguo Li. Darts+: Improved differentiable architecture search with early stopping. ArXiv preprint, 2019.
- [22] Hanxiao Liu, Karen Simonyan, and Yiming Yang. DARTS: differentiable architecture search. In Proceedings of ICLR2019.
- [23] Christos Louizos, Max Welling, and Diederik P. Kingma. Learning sparse neural networks through l<sub>0</sub> regularization. In Proceedings of ICLR2018.
- [24] Chris J. Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. Proceedings of ICLR2017.
- [25] Rabeeh Karimi Mahabadi, James Henderson, and Sebastian Ruder. Compacter: Efficient low-rank hypercomplex adapter layers. 2021.
- [26] Na se Sadat Moosavi, Quentin Delfosse, Kristian Kersting, and Iryna Gurevych. Adaptable adapters. ArXiv preprint, 2022.
- [27] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. Proceedings of NAACL 2018.
- [28] Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. AdapterFusion: Non-destructive task composition for transfer learning. Proceedings of ACL 2021.
- [29] Hieu Pham, Melody Guan, Barret Zoph, Quoc Le, and Jeff Dean. Efficient neural architecture search via parameters sharing. Proceedings of ICMJ2018.
- [30] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. 2018.
- [31] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. Journal of Machine Learning Research, 2020.
- [32] Esteban Real, Alok Aggarwal, Yanping Huang, and Quoc V. Le. Regularized evolution for image classifier architecture search. Proceedings of AAAI2019.

- [33] Andreas Rücklé, Gregor Geigle, Max Glockner, Tilman Beck, Jonas Pfeiffer, Nils Reimers, and Iryna Gurevych. AdapterDrop: On the efficiency of adapters in transformer-based models. In Proceedings of EMNLP, 2021.
- [34] Yi-Lin Sung, Jaemin Cho, and Mohit Bansal. Vi-adapter: Parameter-efficient transfer learning for vision-and-language tasks. 2021.
- [35] Nima Tajbakhsh, Jae Y Shin, Suryakanth R Gurudu, R Todd Hurst, Christopher B Kendall, Michael B Gotway, and Jianming Liang. Convolutional neural networks for medical image analysis: Full training or fine tuning? IEEE transactions on medical imaging, 2016.
- [36] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. Proceedings of NeurIPS, 2017.
- [37] Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. Superglue: A stickier benchmark for general-purpose language understanding systems. Proceedings of NeurIPS, 2019.
- [38] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In Proceedings of ICLR, 2019.
- [39] Ziheng Wang, Jeremy Wohlwend, and Tao Lei. Structured pruning of large language models. In Proceedings of EMNLP, 2020.
- [40] 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. Proceedings of CVPR, 2019.
- [41] Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang, Yanyan Lan, Liwei Wang, and Tie-Yan Liu. On layer normalization in the transformer architecture. In Proceedings of ICLR, 2020.
- [42] Elad Ben Zaken, Shauli Ravfogel, and Yoav Goldberg. Bit t: Simple parameter-efficient fine-tuning for transformer-based masked language-models. 2022.
- [43] Mengjie Zhao, Tao Lin, Fei Mi, Martin Jaggi, and Hinrich Schütze. Masking as an efficient alternative to fine-tuning for pretrained language models. Proceedings of EMNLP, 2020.
- [44] Barret Zoph and Quoc V. Le. Neural architecture search with reinforcement learning. In Proceedings of ICLR, 2017.
- [45] Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V. Le. Learning transferable architectures for scalable image recognition. Proceedings of CVPR, 2018.

## Checklist

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [\[Yes\]](#)
  - (b) Did you describe the limitations of your work? [\[Yes\]](#) In the conclusion section 5, we describe two future works which can also be seen as the limitations of our current work.
  - (c) Did you discuss any potential negative societal impacts of your work? [\[Yes\]](#) See Appendix E.
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)
2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#) Our theoretical results are based on Transformer Architecture and Delta Tuning, which are stated in Section 3.1 and Appendix A.

- (b) Did you include complete proofs of all theoretical results? [Yes] Yes, the set of proofs and illustrations are in Section 3.2 and Appendix C.
3. If you ran experiments...
- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] In the supplemental material.
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] In Section 4.1 and Appendix B
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] We report the standard error of each method in the Table 2.
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] We include these information in Appendix B.3.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- (a) If your work uses existing assets, did you cite the creators? [Yes] We use existing datasets, which are listed in Section 4.1
  - (b) Did you mention the license of the assets? [Yes] In Section 4.1, we use the dataset from Huggingface Library [19].
  - (c) Did you include any new assets either in the supplemental material or as a URL? [Yes] We add our own code and model to supplemental material. We do not generate new data
  - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A] The data is publicly available.
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
  - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

## A Transformer Architecture

### A.1 Transformer Architecture

The pre-trained language model typically adopts a Transformer model as their backbone, which is formed by stacking multiple transformer layers. The parameters of a transformer layer are located in three sub-modules. The self-attention sub-layer captures the interactions between the tokens' hidden representations  $H_{in}^{(l)} \in \mathbb{R}^{s \times d}$

$$H_{\text{attn}}^{(l)} = \text{Softmax}\left(\frac{QK^T}{d}\right)VW_o \quad (18)$$

where

$$Q = H_{in}^{(l)}W_q^{(l)}; K = H_{in}^{(l)}W_k^{(l)}; V = H_{in}^{(l)}W_v^{(l)}; W_q^{(l)}; W_k^{(l)}; W_v^{(l)}; W_o^{(l)} \in \mathbb{R}^{d \times d} \quad (19)$$

Then the output is followed by the fully connected feed-forward layer, which is typically composed of two linear projections

$$H_{\text{ff}}^{(l)} = (H_{\text{attn}}^{(l)}W_1^{(l)} + b_1^{(l)})W_2^{(l)} + b_2^{(l)} \quad (20)$$

where  $W_1 \in \mathbb{R}^{d \times d_m}$ ;  $b_1 \in \mathbb{R}^{d_m}$ ;  $W_2 \in \mathbb{R}^{d_m \times d}$ ;  $b_2 \in \mathbb{R}^d$ , and generally  $d_m \leq d$ .  $\sigma(\cdot)$  is an activation function. Suppose we use a post-norm structure [1], the residual sublayer is a layernorm sublayer to add and normalize the hidden states

$$H_{\text{out}}^{(l)} = \text{LayerNorm}(H_{\text{ff}}^{(l)}) = \frac{H_{\text{ff}}^{(l)}}{\text{var}(H_{\text{ff}})}s^{(l)} + b^{(l)} \quad (21)$$

where  $s^{(l)}; b^{(l)} \in \mathbb{R}^d$ . Note that, for simplicity, we ignore the details of multi-head attention and skip connections between sub-layers, as they do not introduce additional trainable parameters. From a unified view, all sub-layers can be abstracted by a transformation over the sub-layers input,

$$H_{\text{out}} = m(H_{\text{in}}) \quad (22)$$

## B Details of Experiment Configurations

### B.1 Details of experiments

We apply  $\mathcal{S}\Delta$  to the datasets to multitask benchmarks GLUE [1] and SuperGLUE [17] following previous works. All datasets are downloaded from the HuggingFace Data Library. Since the test split of these datasets are held offcially and invisible to the researchers, we randomly split off 2k samples from the training set as validation set for large datasets (QQP, QNLI, ReCoRD, SST2, MNLI), and use the remaining as the training set  $D_{\text{train}}$ , and use the original validation set as the test set  $D_{\text{test}}$ . For other datasets, we randomly split the original validation set in half as the validation and the test set  $D_{\text{test}}$  and use the training set  $D_{\text{train}}$ . The same dataset is splitted differently with different random seeds. For each experiment setting, we repeat the experiment with 8 seeds. In all experiments, the maximum sequence length is 128 for the tasks in GLUE and 256 for the tasks in SuperGLUE. The batch size is 16 for SuperGLUE and 32 for GLUE. Especially, we set the maximum sequence length to 512 and batch size to 8 for ReCoRD. We use model (703M parameters) as the backbone model and we freeze the pre-trained parameters in all experiments except finetuning. We use AdamW as the optimizer with a linear learning rate decay schedule.

For  $\mathcal{S}\Delta$ , following DARTS [2], we equally split the original training set  $D_{\text{train}}$  into two parts:  $D$  for optimizing the parameters in DT modules,  $D_s$  for optimizing the structural parameter. The original validation set, is used to evaluate and save the search structure every steps. The searched structure is retrained in the original training set  $D_{\text{train}}$ , and evaluated in  $D_{\text{val}}$ . We report the average performances and standard deviations on the  $D_{\text{test}}$  across 8 seeds.

### B.2 Hyperparameters

We do pre-experiments on BoolQ and SST2 using learning rates in  $\{3e-5, 3e-4, 3e-3\}$ , learning rate in  $\{1e-3, 1e-2, 1e-1, 1\}$ , and in  $\{0.1, 0.3, 1\}$  and select the learning rate of  $3e-4$ , alpha learning rate

Table 3: Specific parameters on different tasks.

	Search			Re-train		
	Batchsize	Epoch	Validation Steps	Batchsize	Epoch	Validation Steps
GLUE						
CoLA	32	15	100	32	15	100
MNLI	32	1	200	32	3	500
MRPC	32	20	50	32	20	50
QNLI	32	4	200	32	4	200
QQP	32	1	200	32	3	500
SST2	32	5	150	32	5	150
STSB	32	40	100	32	40	100
SuperGLUE						
BoolQ	16	15	200	16	15	200
CB	16	60	20	16	60	20
COPA	16	40	20	16	40	20
MultiRC	16	5	200	16	10	200
ReCoRD	8	1	200	8	1	200
RTE	16	20	50	16	40	50
WiC	16	20	100	16	20	100

Table 4: The metrics we used to evaluate the GLUE and SuperGLUE Benchmark.

	Tasks	Metric
GLUE	CoLA	Matthew's Corr
	SST-2	Accuracy
	MRPC	F1
	QQP	F1
	STS-B	Pearson Corr
	MNLI	Accuracy
	QNLI	Accuracy
SuperGLUE	BoolQ	Accuracy
	CB	Accuracy
	COPA	Accuracy
	MultiRC	F1
	ReCoRD	F1
	RTE	Accuracy
	WiC	Accuracy

of 0.1, and of 1 which perform the best. For, we set it to 1. The specific parameters on different tasks are listed in Table 3. Specifically, for fine-tuning, we try learning rates in  $\{3e-5, 1e-4, 3e-4\}$  and find that  $3e-5$  performs the best. We apply these hyperparameters to all baselines and the re-training phase of our searched structure in Table 2 and Figure 2 and conduct no further hyperparameter-tuning. Therefore, the comparison is fair despite that better performance might be achieved with dataset-specific grid search. The metrics we used to evaluate the GLUE and SuperGLUE Benchmark are in Table 4.

### B.3 Computing Resources

We run all the experiments on NVIDIA V100 32GB GPUs.

## C Theoretical Issues about Sigmoid

### C.1 Binary Concrete Distribution

We use the Binary Concrete Distribution as a soft approximation to Bernoulli Distribution  $B(1; p)$ , where  $p$  is the probability of the sample being 1.

$$z = \frac{1}{\epsilon} \log \frac{u}{(1-u)(1-p)} \quad (23)$$

The probability of  $z > 0.5$  is  $p$ , that is when we sample  $z$ , we can use  $z$  to determine  $u$ 's value,

$$p(z > 0.5) = p \frac{u}{(1-u)(1-p)} > 1 = p(u > 1-p) = p \quad (24)$$

As  $\epsilon$  approaching 0, the Binary Concrete Distribution converges in distribution to the Bernoulli distribution. For any constant  $c < 1$ ,

$$\lim_{\epsilon \rightarrow 0} P(z < c) = 1 - p \quad (25)$$

$$\lim_{\epsilon \rightarrow 0} P(z > 1-c) = p; \quad (26)$$

Therefore  $z \stackrel{d}{\sim} z$ .

### C.2 Gradient of Global Shifted Sigmoid

In global shifted sigmoid function, we multiply a constant with value to the shifted sigmoid to enable a global comparison among the structural parameters of each DT module,

$$p_i = p_j P_i \quad (27)$$

$$i = \frac{j P_j \text{Detach}(p_j)}{j P_j} = 1; \quad (28)$$

Though the value of  $p_i$  is equal to the value of  $p_j$ , The gradient using these two parameterization is different, which result in different behaviour in the optimization of structural parameters.

Using  $p_j$  as the parameterization function, the gradient is

$$\frac{\partial}{\partial i} = \frac{\partial p_j}{\partial i} \frac{\partial}{\partial p_j}; \quad (29)$$

while using  $p_i$  as the parameterization function, the gradient is

$$\frac{\partial}{\partial i} = \sum_j \frac{\partial}{\partial p_j} \frac{\partial p_j}{\partial i} = \sum_{j \in i} \frac{\partial}{\partial p_j} \frac{\partial p_j}{\partial i} + \frac{\partial}{\partial p_i} \frac{\partial p_i}{\partial i} \quad (30)$$

$$= \sum_j \frac{\partial}{\partial p_j} \frac{\partial p_j}{\partial i} P_j + \frac{\partial}{\partial p_i} \frac{\partial p_i}{\partial i} \quad (31)$$

$$= \sum_j \frac{P_j \text{Detach}(p_k)}{(p_k P_k)^2} \frac{\partial p_j}{\partial i} \frac{\partial p_j}{\partial p} + \frac{\partial}{\partial p_i} \frac{\partial p_i}{\partial i} \quad (32)$$

$$= \frac{\partial p_j}{\partial i} \frac{\partial}{\partial i} \sum_j \frac{P_j P_j}{p_k P_k} \frac{\partial}{\partial p} + \frac{\partial}{\partial p} \quad (33)$$

The additional term

$$\sum_j \frac{P_j P_j}{p_k P_k} \frac{\partial}{\partial p} \quad (34)$$

serves as an adjustment from the global structural gradient to the local gradient. In a special case, if the gradients to all  $p_i$ , i.e.,  $\frac{\partial}{\partial p}$ , are equal, no gradient will be pass to the structural parameter which is reasonable.

We conduct ablation studies on three datasets using the global shifted sigmoid (denoted by Global) and the shifted sigmoid without global comparison (denoted by Local). Figure 8 proves the correctness of the global shifted sigmoid parameterization.



Figure 8: The performance difference between local and global shifted sigmoid.

Table 5: The computational resources in the searching phase and re-training phase, Computation time, memory consumption, and the corresponding ratio between searching phase and re-training phase are listed.

Dataset	Time/min			Memory/GB		
	Search	Re-train	Ratio	Search	Re-train	Ratio
RTE	148.6	30.0	5.0	27.7	16.6	1.7
STSB	139.3	26.3	5.3	28.9	10.6	2.7
CoLA	145.6	17.0	8.6	16.1	8.9	1.8

## D Additional Results of Experiments

### D.1 Efficiency of S<sup>3</sup>Delta

We report the detailed consumption of computational resources of S<sup>3</sup>Delta. From Table 5, we can see that the searching time is approximately 5 times of a training time, which is acceptable due to the large search space. Once a sparse structure is searched, the training and inference are as fast as or faster than the other DT methods.

### D.2 Heat maps of $f_{\rho}$ on Each Datasets

The heat maps of  $f_{\rho}$  on different datasets are presented in Figure 9. Most of the datasets follow the similar trend in the activated positions. STSB and ReCoRD's optimal solution is a bit different from the others, which can be further investigated.

## E Potential societal Impact

S<sup>3</sup>Delta focuses on developing a method to find the optimal sparse structure of DT modules. The searching process takes longer time than a single training process, consuming more energy. However, the searched structures are found to be highly transferable. Therefore, we can search for an approximately optimal solution on a group of similar tasks and reuse the structure on other unseen tasks. Consequently, S<sup>3</sup>Delta is environmental friendly. Another potential negative societal impact is the malicious injection of DT modules, which will potentially harm the adaption performance of PTMs.



Figure 9: Heat maps of  $p_i$  on each datasets