

---

# Adaptive Attention Link-based Regularization for Vision Transformers

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1        Although transformer networks are recently employed in the various vision tasks  
2        with the outperforming performance, large training data and a lengthy training  
3        time are required to train a model to disregard an inductive bias. Using trainable  
4        links between the channel-wise spatial attention of a pre-trained Convolutional  
5        Neural Network (CNN) and the attention head of Vision Transformers (ViT), we  
6        present a regularization technique to improve the training efficiency of ViT. The  
7        trainable links are referred to as the attention augmentation module, which is  
8        trained simultaneously with ViT, boosting the training of ViT and allowing it to  
9        avoid the overfitting issue caused by a lack of data. From the trained attention  
10        augmentation module, we can extract the relevant relationship between each CNN  
11        activation map and each ViT attention head, and based on this, we also propose an  
12        advanced attention augmentation module. Consequently, even with a small amount  
13        of data, the suggested method considerably improves the performance of ViT while  
14        achieving faster convergence during training.

## 15    1 Introduction

16    Convolutional Neural Networks (CNN) have become standard for solving image-related tasks using  
17    deep neural networks since the advent of large publicly available datasets [1, 2]. Recently, models  
18    with attention mechanisms mainly adopted in the area of natural language processing are becoming  
19    to take part in solving image-based tasks, which is called the Vision Transformer (ViT) [3, 4]. ViT is  
20    a transformer-based neural network fed by the patches of images with class-token for classification,  
21    replacing its input of the embedded words in natural language processing. Although ViT shows  
22    impressive accuracy compared to modern CNNs by ignoring the inductive bias of locality, a significant  
23    amount of data is required to train the model to achieve satisfactory performance without overfitting  
24    issues. Furthermore, most researchers with the limited computing hardware are not affordable to train  
25    the ViT due to its lengthy training time.

26    The overfitting and lengthy training issues must be solved to broaden the usability of ViT, so many  
27    recent studies have tried to solve the problem [5, 6, 7, 8, 9, 10]. We can divide the studies by three  
28    categories: the advanced architecture-based method [4, 8, 11, 12], the parameter compression-based  
29    method [13, 14], and the knowledge distillation-based method [6, 15]. The advanced architecture-  
30    based methods manipulate the architecture of ViT to achieve improved training efficiency and  
31    generalized prediction even with the small dataset. On the other hand, the parameter compression-  
32    based methods focus on a low-rank approximation of the transformer encoder in ViT, which results  
33    in the boosted training speed and the suppressed overfitting issue. The knowledge distillation-based  
34    methods utilize the prediction of additional CNN models to avoid the overfitting problem and achieve  
35    rapid training convergence. The previous studies have shown the meaningful development of ViT for  
36    the small dataset and the reduced training time.

37 **However, the previous studies have the remaining limitations where the training datasets must be**  
38 **equivalent for both the student and teacher models.** The architectural manipulation of initial ViT [3]  
39 cannot be applied to different versions of ViTs, hence limiting the handling of new ViT models. Even  
40 though the knowledge distillation-based methods can be employed with the small manipulation of a  
41 model such as the knowledge distillation token, it should be assured to have the initial models trained  
42 by a target dataset, which takes the additional costs for the acquisition of initial models before the  
43 training of main ViT model.

44 In this paper, we propose a novel regularization method to reduce the convergence time and avoid the  
45 overfitting problem on a small dataset, simultaneously. The proposed method utilizes an *attention*  
46 *augmentation module* containing multiple trainable weights that estimate the affinity between the  
47 channel-wise activation map of CNN and the head-wise attention map of ViT. Since the attention  
48 augmentation module is located to regularize the attention of ViT heads, the architecture of ViT can  
49 be perfectly preserved, which lets us enable to employ the proposed algorithm in ViT variants based  
50 on the attention maps. In addition, since the attention map can be obtained even when the task of  
51 the pre-trained CNN is not equivalent to the target task, we can employ our method without the  
52 pre-trained CNN model with the same target dataset. We validate our regularization method by using  
53 ImageNet and CIFAR10 datasets with various scenarios, which show the outperforming accuracy  
54 and the reduction of epochs required for its training convergence. Furthermore, we investigate the  
55 important factors for ViT to avoid the overfitting issue by analyzing the trained weights of the attention  
56 augmentation module, and through the investigation, we present the dissimilarity of the deep layers'  
57 roles between CNN and ViT.

58 We can summarize our contributions as following:

- 59 • We propose a novel regularization method to resolve the issues of overfitting and lengthy  
60 training time of ViT through the trainable attention links between the ViT attention maps  
61 and CNN activation maps.
- 62 • The proposed scheme preserves the original architecture of ViT, which results in its general  
63 employment regardless of the architecture of ViT.
- 64 • Through the proposed algorithm, the performance of ViT can be dramatically improved with  
65 the limited size of dataset, and the training time is reduced without the loss of performance  
66 in various scenarios.
- 67 • The relationship between ViT and CNN is analyzed in terms of attentional regions, which  
68 validates the analysis from the previous studies.

## 69 **2 Related Work**

### 70 **2.1 Transformers in Vision**

71 Transformer models introduced by [16] are neural networks that purely utilize the attention mech-  
72 anism. While they have been used broadly in the field of natural language processing, Vision  
73 Transformer (ViT) [3] adapted them in the domain of computer vision with minimal modification to  
74 its architecture. ViT showed comparable performance to CNN in the condition of large pre-training.  
75 For the advanced optimization of ViT, CaiT [4] used layer normalization in ViT layers and changed  
76 the input location of class token to prevent saturation of performance in deep layers. Swin Trans-  
77 formers [12] adopted a hierarchical transformer that computes shifted windows to make it suitable for the  
78 vision domain. PiT [10] introduced the concept of pooling in ViT from CNN, improving the gener-  
79 alization of ViT. T2T-ViT [8] enhanced sample efficiency by reshaping input tokens and changing  
80 the backbone of networks motivated by several CNN architectures. Raghu et al. [17] measured the  
81 similarity of representations between specific layers of CNN and ViT using centered kernel alignment.  
82 With additional relative positional encoding, Cordonnier et al. [18] proved attention mechanisms in  
83 ViTs can perform as convolution layers in CNNs and showed their functional similarity. From the  
84 investigation, ConViT [7] was motivated to use relative positional encoding to give locality – the  
85 inductive bias of CNNs – to ViT.

86 The research was extended to [19], reparameterizing pre-trained convolutional layers as a format of  
87 ConViT. Refiner [11] tackled the over-smoothing problem between tokens in deep layers of ViT, and  
88 relieve it by projecting attention heads into the higher dimensions and applying convolution directly  
89 to attention maps to learn local relationship among the tokens. Those variants of ViT improved  
90 the optimization and data efficiency of the initial ViT model by modifying the architecture itself.

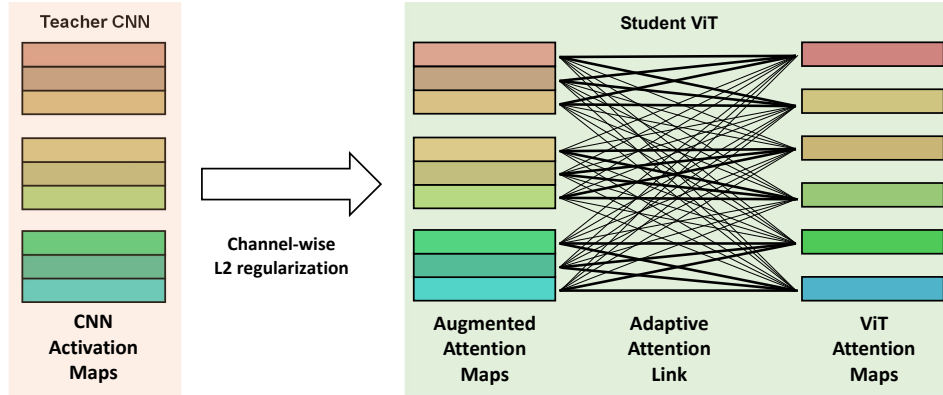


Figure 1: Our distillation procedure. ViT attention maps are augmented through adaptive attention links, which build linear combinations of different original maps. Then, those augmented attention maps mimic CNN activation maps at one-on-one correspondence.

91 However, our method does not touch any part of ViT modules, but only connects attention between  
 92 ViT and CNN, transferring attention through links to give ViT a learning signal from the teacher.

### 93 2.2 Knowledge Distillation:

94 In knowledge distillation, a student model leverages a pre-trained teacher model’s soft prediction  
 95 divided by the same temperature values [20]. The softened predictions can be regarded as label-  
 96 smoothing, and by using them, the student model can achieve the data augmentation effect. Distillation  
 97 between different types of neural architectures has also been proposed, DistilBERT [21] showed the  
 98 effectiveness of a distilled knowledge from BERT [22] into LSTM [23]. DeiT [6] distilled knowledge  
 99 from CNN to ViT, which seems similar to our work. **The knowledge distillation has been also  
 100 employed into the transformer-based model of natural language processing, which results in the  
 101 performance improvement by using the teacher model [24, 25]** However, in contrast to the previous  
 102 study using the prediction for the knowledge distillation, our framework transfers the knowledge  
 103 based on the similarity of the latent feature maps. As a result, we can extend the range of teacher  
 104 models to cover the models in which the prediction vectors differ from the prediction of the task.

105 On the other hand, we can transfer latent representations of teacher models to those of student models.  
 106 FitNets [26] improved the stability of deep network training by guiding the latent layers to the  
 107 teacher’s well-trained latent representation. Zagoruyko et al. [27] considered attention as projected  
 108 activation maps of CNN into a spatial dimension, which could be regarded as spatial attention. They  
 109 showed that spatial attention contains valuable information that is useful to improve the performance  
 110 of the student network. Kim et al. [28] used a paraphraser to extract and pass the teacher’s knowledge  
 111 to the student’s translator to learn its representation. Meanwhile, Heo et al. [5] demonstrated that  
 112 the knowledge transfer based on the neurons’ activation is a more classification-friendly approach  
 113 than the direct transfer using output values. Attention-based feature distillation [29] measured the  
 114 similarities between teacher and student features through attention, which determines the importance  
 115 of knowledge to transfer.

## 116 3 Attention Link-based ViT Regularization

117 In this section, we first explain the backgrounds of the self-attention mechanism and ViT. Then,  
 118 we explain the method to extract the attention maps from ViT, followed by the description of the  
 119 architecture and the training method of the augmented attention module is described. The overall  
 120 framework is depicted in Fig. 1

### 121 3.1 Background of ViT

122 We first explain the self-attention mechanism and the original ViT model referred to by [3]. The  
 123 self-attention mechanism mimics the human cognition system making the attention to the external

124 stimulus, which is designed by a transformer-based model with the attention matrix estimated by  
 125 pairs of key and query.

126 **Self-attention Mechanism:** We define the input sequence by  $\mathbf{X} \in \mathcal{R}^{L \times D_{in}}$  where  $L$  is the length  
 127 of the sequence and  $D_{in}$  means the dimension of one sequential element in the sequence. Then, we  
 128 can estimate the elements of the attention mechanism composed of key, query, and value vectors by  
 129 linearly projecting the input sequence by the corresponding embedding weights  $\mathbf{W}_k$ ,  $\mathbf{W}_q$ , and  $\mathbf{W}_v$ ,  
 130 respectively. Thus, when we define the key, query, and value vectors by  $\mathbf{Q}$ ,  $\mathbf{K}$ , and  $\mathbf{V}$ , respectively,  
 131 we obtain the vectors as following:

$$\mathbf{K} = \mathbf{X}\mathbf{W}_k, \mathbf{Q} = \mathbf{X}\mathbf{W}_q, \mathbf{V} = \mathbf{X}\mathbf{W}_v, \quad (1)$$

132 where  $\mathbf{W} \in \mathcal{R}^{D_{in} \times D_{head}}$  from  $\mathbf{W} \in \{\mathbf{W}_k, \mathbf{W}_q, \mathbf{W}_v\}$  and  $D_{head}$  is the dimension of the head  
 133 embedding.

134 Then, the self-attention of the head can be estimated by:

$$f(\mathbf{X}) = s\left(\mathbf{Q}\mathbf{K}^T / \sqrt{D_{head}}\right)\mathbf{V} \in \mathcal{R}^{L \times D_{head}}, \quad (2)$$

135 where  $s(\mathbf{Z})$  is a function to transfer each row vector of input matrix  $\mathbf{Z}$  by softmax. According to the  
 136 derivation, self-attention can consider the semantic dependency among sequential inputs.

137 Many transformer-based models are based on the architecture stacked by the Multi-Head Self-  
 138 Attention layers (MHSA) containing multiple self-attention heads with independent embedding  
 139 weights. For the given input  $\mathbf{X}$ , we define the output of  $m$ -th self-attention head at  $n$ -th level depth  
 140 by  $f_{(m,n)}(\mathbf{X})$ . Then, we denote the corresponding key, query, and value vectors by  $\mathbf{K}_{(m,n)}$ ,  $\mathbf{Q}_{(m,n)}$ ,  
 141 and  $\mathbf{V}_{(m,n)}$ , respectively.

142 **ViT Framework:** The original ViT model directly employed the conventional transformer-based  
 143 model built for the natural language processing of the visual classification task. We can summarize  
 144 the inference process of the original ViT as following. At first, we divide an input image by  $P^2$   
 145 patches with the same size and sequentially order the patches after their vectorization. Since the  
 146 transformer network is invariant to the order of the sequential data, ViT concatenates positional  
 147 embedding vectors to the input patches to represent the original position of the patch.

148 We define the sequential data obtained from one image by  $\mathbf{X}_0 \in \mathcal{R}^{P^2 \times D_{im}}$ , where  $D_{im}$  is the  
 149 size of the vector linearly projected from the vectorized image patch and the positional embedding  
 150 vector. Before we feed  $\mathbf{X}_0$  into the transformer modules, the trainable class token sized by  $\mathcal{R}^{D_{im}}$  is  
 151 sequentially connected ahead of  $\mathbf{X}'_0$ , which results in  $\mathbf{X}_0 \in \mathcal{R}^{(P^2+1) \times D_{im}}$ .

152 The transformer-based encoders of ViT are the modules containing the series of a layer normalization,  
 153 a self-attention multi-head module, a fully connected layer, and a layer normalization, where every  
 154 normalization layer has a residual connection. We define the serial process by a function of  $\mathbf{X}_{n+1} =$   
 155  $g(\mathbf{X}_n) \in \mathcal{R}^{(P^2+1) \times D}$  where  $D$  is the size of latent vectors. When  $N$  modules are stacked in the  
 156 transformer-based encoder, the class-wise score is estimated by linearly projecting the final output of  
 157 the class token as following:  $p(cls|\mathbf{X}) = softmax(FC(\mathbf{X}_N))$ . For a detailed explanation of ViT,  
 158 you can be referred to [3].

### 159 3.2 Attention Map Extraction

160 We need to compare the ViT attention map and the CNN activation map for our regularization-based  
 161 algorithm. Instead of the relative positional embedding [7] or the attention bias [11], we preserve the  
 162 original architecture of ViT to generalize the usability of our framework to cover the ViT variants.

163 To obtain the attention map from the original ViT, we utilize the attention value between the class  
 164 token and the image patch. The class token takes a key role to determine the final prediction, so  
 165 we can assume that the attention to the class token may represent the importance of image patches  
 166 for the classification result. Thus, when feeding the class token into the transformer module as its  
 167 query vector, we obtain the attention value by estimating the dot product between the embedding  
 168 vectors of the class token and the corresponding image patch. Then, for  $m$ -th head in  $n$ -th multi-head  
 169 self-attention layer, we can estimate the attention value as:

$$\mathbf{A}_{(m,n)} = Rec\left(s\left(\mathbf{Q}_{(m,n)}[0]\mathbf{K}_{(m,n)}[1:]^T\right)\right) \in \mathcal{R}^{P \times P}, \quad (3)$$

170 where  $Rec$  is a function to reconstruct the rectangular matrix of  $\mathcal{P} \times \mathcal{P}$  from its input vector of  $\mathcal{P}^\epsilon$   
 171 according to the order of the sequential patches, and  $\mathbf{Q}_{(m,n)}[0]$  and  $\mathbf{K}_{(m,n)}[1:]$  represent the first  
 172 query vector of the class token and the key vectors of the image patches, respectively.

173 In the case of the CNN activation map, we extract the activation maps after the normalization of every  
 174 convolution block. Instead of integrating the channel-wise activation maps, we consider the separated  
 175 activation maps independently to improve the degree of freedom of our attention augmentation  
 176 module. In contrast to the constant resolution of ViT attention maps, the resolution of the CNN  
 177 activation maps decreases with deep layers by pooling layers and strides of convolution layers. Thus,  
 178 to preserve the resolution of every activation map, we resize all the CNN activation maps to have the  
 179 same size with the ViT attention maps **by using bicubic interpolation**. We define the  $c$ -th resized CNN  
 180 activation map by  $\mathbf{B}_c$ , where  $c \in \{1, \dots, C\}$  and  $C$  is the number of entire CNN activation maps.

### 181 3.3 Attention Augmentation Module

#### 182 3.3.1 Module Architecture

183 Even though both the CNN activation and ViT attention maps represent the key parts of the target  
 184 object for the prediction, their distribution such as a center point and a variance would be different  
 185 from each other due to the dissimilarities of their operations. For example, while the ViT attention map  
 186 is distributed between 0 and 1 because of the softmax estimation, the values in the CNN activation  
 187 map are normalized by a batch normalization, which can contain negative values. Furthermore, in  
 188 general, the number of CNN activation maps is much larger than the number of ViT attention maps  
 189 due to the large channel-wise depth of CNNs. As a result, it is impossible to directly compare each of  
 190 the CNN activation maps with the ViT attention maps.

191 The attention augmentation module is designed to solve the problems of different distributions and a  
 192 varying number of maps. We design the attention augmentation module to contain multiple attention  
 193 links which are the trainable weight parameters to scale the ViT attention maps. By estimating the  
 194 weighted summation of ViT attention maps with the attention links, we can obtain the augmented  
 195 attention maps where the number is equivalent to the number of CNN activation maps. Thus, we can  
 196 estimate the augmented attention maps as following:

$$197 \mathbf{A}_c^+ = \sum_{m,n=1}^{M,N} w_c^{(m,n)} \mathbf{A}_{(m,n)} + b_c, \quad (4)$$

197 where  $w_c^{(m,n)}$  and  $b_c$  are the attention link and a trainable bias for  $c$ -th augmented attention map  
 198 ( $c \in \{1, \dots, C\}$ ), respectively.  **$M$  is the number of self-attention heads in one level depth and  $N$**   
 199 **presents the maximum level depth**. Note that the weight of attention link  $w_c^{(m,n)}$  is used to analyze  
 200 the strength of connectivity for each CNN/ViT layer in section 4.1.

201 We implement the augmented attention module by a  $1 \times 1$  convolution layer generating  $C$  augmented  
 202 attention maps from a tensor of  $\mathcal{R}^{P^2 \times MN}$  where the ViT attention maps are stacked. Because we  
 203 only use the augmented attention maps only for the training loss, the attention augmentation module  
 204 has no role in the inference, which can be removed after the training of ViT.

#### 205 3.3.2 Module Training

206 By using the augmented attention module, we can obtain the same number of augmented attention  
 207 maps  $\mathbf{A}_c^+$  with the CNN attention maps  $\mathbf{B}_c$ . To ignore the remaining scale gap between the two  
 208 maps, we first apply the  $l_2$  normalization, and then the mean squared error is estimated to build the  
 209 attention-based regularization loss as:

$$209 \mathcal{L}_{att} = \|\mathbf{A}_c^+ / \|\mathbf{A}_c^+\|_2 - \mathbf{B}_c / \|\mathbf{B}_c\|_2\|_2. \quad (5)$$

210 Then, we integrate the attention-based regularization loss  $\mathcal{L}_{att}$  with the cross-entropy loss  $\mathcal{L}_{CE}$  of  
 211 original ViT as:

$$211 \mathcal{L} = \mathcal{L}_{CE} + \lambda \mathcal{L}_{att}, \quad (6)$$

212 where  $\lambda$  is a scaling factor to control the effect of our regularization. Since the regularization loss can  
 213 work as an obstacle to ignoring the inductive bias, referred by [7], we suppress the value of  $\lambda$  at the  
 214 specified epochs to increase the effectiveness of the cross-entropy loss. We exponentially decay the  
 215 value of  $\lambda$  by multiplying a decay constant between 1 and 0 at every epoch.

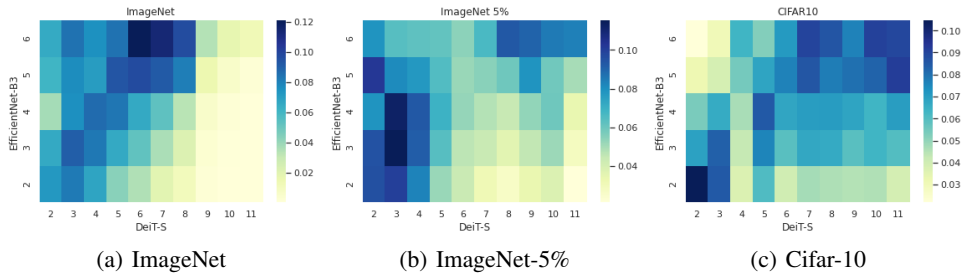


Figure 2: Relations between CNN Activation Maps and ViT attention maps. The x-axis and y-axis indicates the depth of ViT layer and CNN block, respectively. We obtain the heatmap by averaging the magnitudes of attention links.

CNN Block Level		1	2	3	4	5	6	7
ViT Layer Level	$\alpha$ -link	1 ~ 3	1 ~ 5	3 ~ 5	4 ~ 6	4 ~ 6	6 ~ 11	7 ~ 12
	$\beta$ -link	1 ~ 2	1 ~ 5	3 ~ 5	4 ~ 6	4 ~ 7	6 ~ 10	7 ~ 10

Table 1: Selective Link Configuration.

## 216 4 Link Selection for Advanced Regularization

217 In this section, we build the advanced architecture of the attention augmentation module based on the  
 218 analysis of the fully-trained attention links. After showing the resultant attention links, we explain the  
 219 advanced link designed by considering the relation between CNN activation and ViT attention maps.

### 220 4.1 Analysis of Resultant Attention Links

221 We visualize heat maps to show the scale distribution of the attention links after their training.  
 222 To compare the relationship between the ViT attention and the CNN activation maps, we only  
 223 consider the magnitude of the weight parameters in the attention links. As shown in Fig. 2, we obtain  
 224 multiple heat maps by using three datasets including *ImageNet* [30], a 5% subset of *ImageNet*, and  
 225 *Cifar-10* [31].

226 As analyzed in many previous studies [18, 17], the ViT attention maps are highly related to the CNN  
 227 activation maps located at a similar level. The results validate that the self-attention multi-heads of  
 228 ViT can train the hierarchical information by the stacked architecture, which is similar to the training  
 229 mechanism of CNN. Thus, the CNN activation maps would be helpful to regularize the ViT attention  
 230 maps if we can find their level similarity.

231 In addition to the well-known hypothesis, the heat maps from the attention links show the interesting  
 232 characteristic where the high-level heads are not regularized when the large dataset is given. **When a  
 233 large dataset is provided, we can observe the suppressed magnitudes of the attention links at high-level  
 234 heads in comparison to small datasets.** The discovery exploits the high-level heads should not be  
 235 trained by the regularization of the high-level CNN layers, which verifies the high-level heads can  
 236 represent more complicated semantic information than CNN layers. In other words, the representation  
 237 can be seen out of the inductive bias of CNNs, so we can show **semantic information that overwhelms**  
 238 the inductive bias is hard to be trained without a large dataset.

239 DeiT [6] and Swin Transformers [12] showed that the employment of the inductive bias of CNN is  
 240 effective for the training of ViT. At the end of the training, we observed that high-level heads are  
 241 disconnected from augmented attention maps, which means they are no more regularized by CNN  
 242 activation maps. This indicates that high-level heads would escape the inductive bias of locality, and  
 243 they are trained by a long-range dependency that cannot be acquired by CNNs.

244 We can summarize the two hypotheses from the analysis as following:

- 245 • The ViT attention and the CNN activation maps are similar to each other at a similar level.
- 246 • The high-level ViT heads can present the semantic information that cannot be represented  
 247 by the CNN layers, but training the semantic information requires a large dataset.

Train Size	Top-1				Top-5			
	DeiT	ConViT	AAL (Ours)	Gap	DeiT	ConViT	AAL(Ours)	Gap
5%	34.8%	47.8%	<b>51.7%</b>	49%/8%	57.8%	70.7%	<b>75.9%</b>	31%/7%
10%	48.0%	59.6%	<b>64.7%</b>	35%/9%	71.5%	80.3%	<b>85.8%</b>	20%/7%
30%	66.1%	73.7%	<b>76.1%</b>	15%/4%	86.0%	90.7%	<b>93.0%</b>	8%/3%
50%	74.6%	78.2%	<b>78.9%</b>	6%/1%	91.8%	93.8%	<b>94.5%</b>	3%/1%
100%	79.9%	<b>81.4%</b>	81.0%	1%/0%	95.0%	<b>95.8%</b>	95.5%	1%/0%

Table 2: ImageNet test accuracy with different sampling ratios. The Gap columns represent the relative performance improvement of the AAL over DeiT and ConViT, respectively.

## 248 4.2 Selective Attention Link

249 Based on the analysis, we additionally propose the selective attention link to improve the training  
 250 efficiency of the attention augmentation module. Instead of the full link in the original attention  
 251 augmentation module, only a part of the attention links are utilized to obtain the augmented attention  
 252 maps. Based on the two hypothesis from our analysis, the augmented attention map is generated  
 253 only by the ViT attention maps with the similar levels, and no link is connected to the high-level ViT  
 254 attention maps when the training dataset is sufficiently large.

255 Accordingly, we build two types of selective attention link, which are denoted by  $\alpha$ -link and  $\beta$ -link.  
 256  $\alpha$ -link connects the ViT attention maps to only the augmented attention maps at a similar level.  $\beta$ -link  
 257 is similar to the  $\alpha$ -link but the links to the high-level heads are entirely disconnected. The detailed  
 258 connections are given in Table 1.

## 259 5 Experiments

260 In experiments, we showed that transferring attention from pre-trained CNN models to ViTs can inject  
 261 CNN’s inductive bias (i.e locality) naturally in standard self-attention layers, without the necessity  
 262 of additional modules extending the self-attention network. Thus, we examine how efficiently the  
 263 method helps ViT to be converged for optimal performance, especially showing a large gap in a small  
 264 data regime.

### 265 5.1 Experimental Settings

266 **Implementation Details:** The computing resource used in our experiments is Nvidia A100.  
 267 If not mentioned otherwise, the student ViT model used for experiments is DeiT-S (distilled  
 268 version) and used EfficientNet-B3 [32] as the teacher CNN model. We set  $\lambda$  to 2000 and the decay constant for  $\lambda$  is set to 0.99 for the first 200  
 270 epochs and 0.98 for the last 100 epochs. For a fair comparison, we preserve the values of the remaining  
 271 hyperparameters and the training strategies from our baseline model of DeiT [6].  
 272  
 273

Models	DeiT-B	ConViT	AAL
Top-1	97.5%	95.4%	97.5%

Table 3: CIFAR10 Top-1 test accuracy

274 **Comparisons and Dataset:** For comparison, we consider two previous studies, which include DeiT  
 275 and ConViT. DeiT utilizes the knowledge distillation method to improve the ViT-based models, and  
 276 ConViT shows the state-of-the-art performance when the training data is given sufficiently even  
 277 without using the knowledge distillation methods. To show the generality of our algorithm, we utilize  
 278 four classification datasets: ImageNet, CIFAR10, Caltech-UCSD Birds-200-2011 (CUB-200), and  
 279 Oxford 102 Flowers (Flower-102). In the case of ImageNet, we extract the subsets randomly sampled  
 280 with the various ratios (5%, 10%, 30%, 50%), maintaining class balance, to show the validity of the  
 281 proposed algorithm when the insufficient data is given for the training.

### 282 5.2 Quantitative Results

283 We first perform the comparisons with the various subsets of ImageNet. As shown in Table 2, the  
 284 proposed algorithm shows the state-of-the-art performance when the ImageNet subsets are used  
 285 to train the model. The performance of our framework is similar to that of ConViT when the  
 286 entire dataset is considered for training. However, the performance gap between our framework and  
 287 ConViT becomes enlarged with the insufficient training data. Furthermore, we should notice that

Train Size	Top-1			Top-5		
	Full-link	Selective-link	Gap	Full-link	Selective-link	Gap
5%	48.9%	<b>51.7%</b>	5.7%	73.6%	<b>75.9%</b>	3.1%
10%	63.0%	<b>64.7%</b>	2.6%	84.6%	<b>85.8%</b>	1.4%
30%	75.2%	<b>76.1%</b>	1.2%	92.4%	<b>93.0%</b>	0.6%
50%	78.5%	<b>78.9%</b>	0.5%	94.3%	<b>95.0%</b>	0.7%
100%	<b>81.0%</b>	80.9%	-0.1%	95.5%	<b>95.5%</b>	0.0%

Table 4: ImageNet test accuracy with different sampling ratios. The Gap columns represent the relative performance improvement of Selective-link.

Teacher Model	Student Model	Teacher Model Top-1	Student Model Top-1
ResNet34	DeiT-S w/ distill	73.3%	79.4%
EfficientNet-B3	DeiT-B w/ distill	81.1%	82.8%

Table 5: ImageNet test accuracy with various teacher and student models

288 EfficientNet-B3 that is our teacher model needs only 12.2M parameters, which is much smaller  
289 than 86.6M parameters of RegNetY-16GF [33] used in DeiT [6]. Thus, we can validate that our  
290 proposed framework can overwhelm DeiT-B even by using the light teacher model. In addition, while  
291 ConViT-S needs 5M more parameters than ours or DeiT, our method outperforms both of DeiT and  
292 ConViT-S, which validates the efficiency of our framework. The quantitative results for CUB and  
293 Flower datasets are represented in the supplementary material.

294 Table 3 shows the experimental results with CIFAR10 dataset. The results verifies that the proposed  
295 algorithm can increase the robustness to the insufficient size of training data. In addition, in the DeiT  
296 paper, 7200 training epochs were needed to achieve 97.5% top-1 test accuracy when training from  
297 scratch using the DeiT-B model which has more attention heads than DeiT-S. On the other hand, our  
298 method only needed to train 300 training epochs to reach the same test accuracy while using the  
299 DeiT-S model, which validates its training efficiency.

### 300 5.3 Analysis

301 In addition to the following analysis, we present the additional experiments to show the validity of  
302 our framework in the supplementary material. The additional experiments include the performance of  
303 weakly supervised object localization, the qualitative results for attention maps, the learning curve,  
304 and epoch-wise qualitative changes of attention links.

305 **Effectiveness of Selective Links:** To show that our selective attention link-based transfer efficiently  
306 matches ViT attention maps with CNN activation maps, we compared two different settings on the  
307 attention augmentation module. *Full-link* fully connects each ViT attention map to produce augmented  
308 attention maps that match CNN activation maps as one-to-one channel-wise correspondence. In the  
309 case of the full ImageNet dataset,  $\beta$ -link was used for the selective link, while we utilized  $\alpha$ -link for  
310 the other small datasets. As shown in Table 4, the attention transfer with a fully connected attention  
311 link shows superior performance to the accuracy of DeiT and ConViT (Table 2) in low data regime,  
312 and the selective attention links show further improvement from its results.

313 **Robustness to Variety of Models:** We add results with the variants composing of different teacher  
314 and student models to show the generality of our method upon various environments. As shown in  
315 Table 5, the proposed framework successfully improves the performance of its teacher model even  
316 with the different teacher and student models. Interestingly, when we use a light teacher model, we  
317 can achieve the large performance gap between the teacher and student models.

318 **Data and Model Efficiency:** Our additional trainable module, which is the attention augmentation  
319 module, includes only a single 1x1 Conv layer which augments the attention maps of the student ViT.  
320 In our default settings, the number of the parameter is 0.068M, which is quite small compared to  
321 DeiT-S of 22M parameters. As we mentioned, the lengthy training time of ViT is a critical drawback  
322 especially when the computational resources are limited. The reduced training time of our method  
323 can be validated through the learning curve represented in the supplementary material.



Strong Data Aug. (Default)			Weak Data Aug.		
Methods	Top-1	Top-5	Methods	Top-1	Top-5
Cross Entropy (CE)	91.3%	99.6%	Cross Entropy (CE)	84.2%	98.7%
CE + AAL	97.4%	99.9%	CE + AAL	92.5%	99.7%
CE + Soft Distillation	91.0%	99.6%	CE + Soft Distillation	84.0%	98.9%
CE + Hard Distillation	92.0%	99.8%	CE + Hard Distillation	85.1%	99.0%
CE + AAL + Hard Dist.	96.5%	99.9%	CE + AAL + Hard Dist.	94.1%	99.7%

Table 6: Ablation Test with different settings of data augmentation and distillation methods

324 **Ablation studies:** For additional verification of our knowledge transfer method, we trained on  
325 CIFAR10 with different scenarios. To sternly check the performance difference of knowledge dis-  
326 tillation effect from each method, we used much weaker data augmentation than the setting used  
327 in other experiments with only simple techniques such as random crops and horizontal flips. This  
328 allows us to confirm the data efficiency in a low data regime. In addition, we compared our method  
329 to other knowledge distillation methods introduced by DeiT with a teacher model pre-trained on  
330 the CIFAR10 dataset. As shown in Table 6, for both soft label distillation and hard label distillation,  
331 our method outperforms the class prediction-based distillation method. From this result, we could  
332 infer that directly transferring attention gives a better learning signal than giving the teacher model’s  
333 output predictions. Furthermore, we could confirm that knowledge earned from a large dataset can  
334 give a good learning direction. **Compared to the result of Table 3 where the teacher model pre-trained  
335 by ImageNet was used, the performance drops due to the lack of information in the teacher model  
336 pre-trained by CIFAR10.** This could be another advantage since teacher models in DeiT are restricted  
337 to be trained on the target dataset to give proper output prediction.

338 **Various Baseline:** In Table 7, we show that applying our method is not only limited to standard ViT.  
339 In the experiments, we employ our method to Pooling based ViT (PiT-S) [10], **and** we observed the  
340 sample efficiency of the model increased by a large margin using our method.

341 **Robustness to random initialization:** We per-  
342 form several trials with different random seeds  
343 as shown in Table 8. Our algorithm shows con-  
344 sistency even with the various initial parameters.  
345 Due to our limited computation, we run 240  
346 epochs of training in contrast to 300 epochs of  
347 training in the default setting.

	PiT-S	PiT-S + AAL
Top-1	12.2%	44.0%
Top-5	25.2%	67.3%

Table 7: Our method with PiT (ImageNet 5%)

348 **Prediction-based distillation and Fine-tuning:** Our method can be jointly applied with class  
349 prediction-based distillation. Thus, both Table 2 and Table 3 verify that our method can show synergy  
350 with the distillation method proposed in [6]. In addition, we perform the additional comparison to  
351 fine-tuning algorithms with self-supervised learning (SSL) for ViT [34], SSL with linear classifier,  
352 and SSL with k-NN classifier respectively show 77% and 74.5% for ImageNet top-1 test accuracy.  
353 **Every accuracy is lower than our performance of 81.0% with the same ViT model, which validates  
354 the synergy of our method with the self-supervised fine-tuning mechanisms.**

## 355 6 Conclusion

356 In this paper, we have introduced a novel method  
357 of transferring knowledge from CNN to ViT.  
358 By accessing attention of CNNs and adaptively  
359 adopting them, student ViT was able to earn high  
360 quality of learning signal with CNN’s inductive  
361 bias. By applying our method, we could train  
362 ViT in less training epochs without overfitting  
363 even with the small dataset or limited labeled  
364 data. Also, we revealed relations between inter-  
365 mediate representations from those different types of neural networks, which varied due to the training  
366 dataset. Furthermore, by analyzing those relationship with trained attention links, we could take  
367 advantage of more efficient connection between networks. We leave wider application of our methods  
368 to new ViT architectures in future works.

	Trial I	Trial II	Trial III
Top-1	47.3%	46.5%	47.2%
Top-5	71.9%	71.3%	72.0%

Table 8: Repeated trials (ImageNet 5%)

369 **References**

- 370 [1] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional  
371 neural networks. In *NIPS*, 2012.
- 372 [2] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.  
373 In *CVPR*, 2016.
- 374 [3] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas  
375 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and  
376 Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*,  
377 2021.
- 378 [4] Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Hervé Jégou. Going deeper  
379 with image transformers. *ArXiv*, abs/2103.17239, 2021.
- 380 [5] Byeongho Heo, Minsik Lee, Sangdoon Yun, and Jin Young Choi. Knowledge transfer via distillation of  
381 activation boundaries formed by hidden neurons. In *AAAI*, 2019.
- 382 [6] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé  
383 Jégou. Training data-efficient image transformers & distillation through attention. In *ICML*, 2021.
- 384 [7] Stéphane d’Ascoli, Hugo Touvron, Matthew L. Leavitt, Ari S. Morcos, Giulio Biroli, and Levent Sagun.  
385 Convit: Improving vision transformers with soft convolutional inductive biases. In *ICML*, 2021.
- 386 [8] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zi-Hang Jiang, Francis E.H. Tay, Jiashi Feng,  
387 and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. In *ICCV*,  
388 2021.
- 389 [9] Andreas Steiner, Alexander Kolesnikov, Xiaohua Zhai, Ross Wightman, Jakob Uszkoreit, and Lucas Beyer.  
390 How to train your vit? data, augmentation, and regularization in vision transformers, 2021.
- 391 [10] Byeongho Heo, Sangdoon Yun, Dongyoon Han, Sanghyuk Chun, Junsuk Choe, and Seong Joon Oh.  
392 Rethinking spatial dimensions of vision transformers. *ArXiv*, abs/2103.16302, 2021.
- 393 [11] Daquan Zhou, Yujun Shi, Bingyi Kang, Weihao Yu, Zihang Jiang, Yuan Li, Xiaojie Jin, Qibin Hou, and  
394 Jiashi Feng. Refiner: Refining self-attention for vision transformers. *ArXiv*, abs/2106.03714, 2021.
- 395 [12] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin  
396 transformer: Hierarchical vision transformer using shifted windows. In *2021 IEEE/CVF International  
397 Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*, pages 9992–  
398 10002. IEEE, 2021.
- 399 [13] Krzysztof Marcin Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane,  
400 Tamás Szepesvári, Peter Hawkins, Jared Quincy Davis, Afroz Mohiuddin, Lukasz Kaiser, David Benjamin  
401 Belanger, Lucy J. Colwell, and Adrian Weller. Rethinking attention with performers. In *ICLR*, 2021.
- 402 [14] Sinong Wang, Belinda Z. Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with  
403 linear complexity. *ArXiv*, abs/2006.04768, 2020.
- 404 [15] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand  
405 Joulin. Emerging properties in self-supervised vision transformers. *ArXiv*, abs/2104.14294, 2021.
- 406 [16] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz  
407 Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, 2017.
- 408 [17] Maithra Raghu, Thomas Unterthiner, Simon Kornblith, Chiyuan Zhang, and Alexey Dosovitskiy. Do vision  
409 transformers see like convolutional neural networks? *ArXiv*, abs/2108.08810, 2021.
- 410 [18] Jean-Baptiste Cordonnier, Andreas Loukas, and Martin Jaggi. On the relationship between self-attention  
411 and convolutional layers. In *ICLR*, 2020.
- 412 [19] Stéphane d’Ascoli, Levent Sagun, Giulio Biroli, and Ari Morcos. Transformed cnns: recasting pre-trained  
413 convolutional layers with self-attention. *ArXiv*, abs/2106.05795, 2021.
- 414 [20] Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network. *ArXiv*,  
415 abs/1503.02531, 2015.
- 416 [21] Raphael Tang, Yao Lu, Linqing Liu, Lili Mou, Olga Vechtomova, and Jimmy Lin. Distilling task-specific  
417 knowledge from BERT into simple neural networks. *ArXiv*, abs/1903.12136, 2019.

- 418 [22] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep  
419 bidirectional transformers for language understanding. In *NAACL-HLT*, 2019.
- 420 [23] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Comput.*, 9(8):1735–1780,  
421 1997.
- 422 [24] Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. Minilm: Deep self-attention  
423 distillation for task-agnostic compression of pre-trained transformers. In *NeurIPS*, 2020.
- 424 [25] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu.  
425 Tinybert: Distilling bert for natural language understanding. In *EMNLP*, 2020.
- 426 [26] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua  
427 Bengio. Fitnets: Hints for thin deep nets. In *ICLR*, 2015.
- 428 [27] Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the performance  
429 of convolutional neural networks via attention transfer. In *ICLR*, 2017.
- 430 [28] Jangho Kim, Seonguk Park, and Nojun Kwak. Paraphrasing complex network: Network compression via  
431 factor transfer. In *NIPS*, 2018.
- 432 [29] Mingi Ji, Byeongho Heo, and Sungrae Park. Show, attend and distill: Knowledge distillation via attention-  
433 based feature matching. In *AAAI*, 2021.
- 434 [30] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical  
435 image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255,  
436 2009.
- 437 [31] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 (canadian institute for advanced research).
- 438 [32] Mingxing Tan and Quoc V. Le. Efficientnet: Rethinking model scaling for convolutional neural networks.  
439 In *ICML*, 2019.
- 440 [33] Ilija Radosavovic, Raj Prateek Kosaraju, Ross B. Girshick, Kaiming He, and Piotr Dollár. Designing  
441 network design spaces. *CoRR*, abs/2003.13678, 2020.
- 442 [34] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand  
443 Joulin. Emerging properties in self-supervised vision transformers. In *2021 IEEE/CVF International  
444 Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*, pages 9630–  
445 9640. IEEE, 2021.

## 446 Checklist

447 The checklist follows the references. Please read the checklist guidelines carefully for information on  
448 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or  
449 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing  
450 the appropriate section of your paper or providing a brief inline description. For example:

- 451 • Did you include the license to the code and datasets? **[Yes]** See Section ??.
- 452 • Did you include the license to the code and datasets? **[No]** The code and the data are  
453 proprietary.
- 454 • Did you include the license to the code and datasets? **[N/A]**

455 Please do not modify the questions and only use the provided macros for your answers. Note that the  
456 Checklist section does not count towards the page limit. In your paper, please delete this instructions  
457 block and only keep the Checklist section heading above along with the questions/answers below.

- 458 1. For all authors...
- 459 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
460 contributions and scope? **[Yes]** See Section 1
- 461 (b) Did you describe the limitations of your work? **[Yes]** See Section 6
- 462 (c) Did you discuss any potential negative societal impacts of your work? **[No]** Our research  
463 do not have any potential negative societal impacts

- 464 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
465 them? [Yes]
- 466 2. If you are including theoretical results...
- 467 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Section 3  
468 (b) Did you include complete proofs of all theoretical results? [Yes] See Section 4
- 469 3. If you ran experiments...
- 470 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
471 mental results (either in the supplemental material or as a URL)? [Yes] See Section  
472 4
- 473 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
474 were chosen)? [Yes] See Section 5
- 475 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
476 ments multiple times)? [Yes] See Section 5
- 477 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
478 of GPUs, internal cluster, or cloud provider)? [Yes] See Section 5
- 479 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 480 (a) If your work uses existing assets, did you cite the creators? [Yes] See Section 5  
481 (b) Did you mention the license of the assets? [Yes] See Section 5  
482 (c) Did you include any new assets either in the supplemental material or as a URL? [No]  
483 There are no new assets
- 484 (d) Did you discuss whether and how consent was obtained from people whose data you're  
485 using/curating? [No] I only use open sources
- 486 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
487 information or offensive content? [No]
- 488 5. If you used crowdsourcing or conducted research with human subjects...
- 489 (a) Did you include the full text of instructions given to participants and screenshots, if  
490 applicable? [No]
- 491 (b) Did you describe any potential participant risks, with links to Institutional Review  
492 Board (IRB) approvals, if applicable? [No]
- 493 (c) Did you include the estimated hourly wage paid to participants and the total amount  
494 spent on participant compensation? [No]