
ScaleKD: Strong Vision Transformers Could Be Excellent Teachers

Anonymous Author(s)

Affiliation
Address
email

Abstract

1 In this paper, we question if well pre-trained vision transformer (ViT) models could
2 be used as teachers that exhibit scalable properties to advance cross architecture
3 knowledge distillation research, in the context of adopting mainstream large-scale
4 visual recognition datasets for evaluation. To make this possible, our analysis under-
5 lines the importance of seeking effective strategies to align (1) feature computing
6 paradigm differences, (2) model scale differences, and (3) knowledge density dif-
7 ferences. By combining three closely coupled components namely *cross attention*
8 *projector*, *dual-view feature mimicking* and *teacher parameter perception* tailored
9 to address the alignment problems stated above, we present a simple and effective
10 knowledge distillation method, called *ScaleKD*. Our method can train student
11 backbones that span across a variety of convolutional neural network (CNN), multi-
12 layer perceptron (MLP), and ViT architectures on image classification datasets
13 with significantly improved performance, achieving state-of-the-art distillation
14 performance. For instance, taking a well-trained Swin-L as the teacher model,
15 our method gets 75.15%|82.03%|84.16%|78.63%|81.96%|83.93%|83.80%|85.53%
16 top-1 accuracies for MobileNet-V1|ResNet-50|ConvNeXt-T|Mixer-S/16|Mixer-
17 B/16|ViT-S/16|Swin-T|ViT-B/16 models trained on ImageNet-1K dataset from
18 scratch, showing 3.05%|3.39%|2.02%|4.61%|5.52%|4.03%|2.62%|3.73% absolute
19 gains to the individually trained counterparts with the same experimental settings.
20 Intriguingly, when scaling up the size of teacher models or their pre-training
21 datasets, our method showcases larger gains to student models. Empirically, the
22 student backbones trained by our method transfer well on downstream MS-COCO
23 and ADE20K datasets. Moreover, our method shows the potential to be an efficient
24 substitution for the time-intensive pre-training of any target student on large-scale
25 datasets if a strong pre-trained ViT is available. The code will be released soon.

26

1 Introduction

27 The great success of deep learning in computer vision (CV) has been driven by an explosion of neural
28 network architectures among which convolutional neural networks (CNNs) [1–3], vision transformers
29 (ViTs) [4, 5] and multi-layer perceptrons (MLPs) [6–8] are three major model categories. While
30 CNNs were the de facto models for about a decade, recent progress shows that large ViT models have
31 attained state-of-the-art performance on many visual recognition tasks such as image classification,
32 image segmentation, and object detection. In principle, ViTs extend the philosophy of predominant
33 transformer architectures [9] in natural language processing (NLP) to vision tasks. They convert an
34 image into a sequence of equal-sized patches treated as tokens resembling words in NLP, then apply
35 the dot-product self-attention mechanism over the sequence of image patches. ViTs designed in this
36 way couple with a powerful data-hungry learning paradigm: models are first pre-trained on massive
37 datasets (with supervised or self-supervised [10, 11] or cross-modality learning [12, 13]) and then
38 fine-tuned on target datasets (with supervised learning). As the size of ViT models and pre-training
39 datasets increases, the pre-trained models tend to have improved generalization performance. Despite

40 this notable model performance scalability, the pre-training process of ViTs leads to significantly huge
41 expenses. Besides, in industry, CNNs and MLPs are still widely used, due to the wider availability
42 of effective implementations and optimizations compared to ViTs. Furthermore, large models
43 including ViTs, CNNs, and MLPs are memory-hungry and computationally intensive, prohibiting the
44 deployment of them in many resource-constrained application scenarios.

45 In parallel, knowledge distillation (KD) has proven to be a promising model compression pathway
46 and has attracted lots of research interests. It relies on a teacher-student framework that transfers the
47 knowledge learned by a large teacher model to a compact student model, aiming to make the student
48 model can have improved performance to substitute the teacher model in deployment. However,
49 most existing KD methods [14–35] focus on CNN architectures, and usually perform evaluation on
50 small datasets with non-mainstream student models for industrial applications, lagging far behind the
51 evolution of neural network architectures. In this paper, *we attempt to connect knowledge distillation*
52 *research with well pre-trained ViT models that stand out for their remarkable scalability, via a new*
53 *perspective*. Specifically, *we question whether well pre-trained ViT models could be used as teachers*
54 *that effectively transfer scalable properties to target student models with different typed architectures*
55 *such as CNNs and MLPs or heterogeneous ViT structures (in this work, we refer ‘cross architecture*
56 *KD’ to such a more generalized formulation), in the context of using mainstream large-scale visual*
57 *recognition benchmarks. Note there have been few recent efforts on cross architecture KD research,*
58 *but they explored narrow focuses such as logits distillation [36] and feature transform [37, 38],*
59 *following the paths previously explored in CNNs based KD methods. To answer the question in our*
60 *motivation, we think the knowledge transfer difficulties are rooted in the following three aspects of*
61 *differences: (1) Differences in feature computing paradigm*. In terms of semantic units, ViTs operate
62 *on a sequence of equal-sized image patches added with positional embeddings, whereas CNNs*
63 *operate on regular grids of pixels. In terms of core operations, ViTs rely on self-attention operations*
64 *to model global feature dependencies, whereas CNNs rely on convolution operations to model local*
65 *features. Although MLPs also use a patchify stem as ViTs, they rely on fully connected operations*
66 *instead of self-attention operations and do not use positional embeddings, showing inferior feature*
67 *learning ability. These differences in feature computing paradigm pose the first knowledge transfer*
68 *barrier to overcome. (2) Differences in model scale*. On the micro scale, model scale differences
69 *among ViTs, CNNs, and MLPs lie in network width, network depth, building blocks, etc. On the*
70 *macro scale, model scale differences come from the capability of scaling the model size for ViTs,*
71 *CNNs and MLPs towards better performance and generalization ability. As a result, these differences*
72 *in model scale make the capacity of different network architectures typically vary significantly,*
73 *emerging as the second knowledge transfer barrier to address. (3) Differences in knowledge density*.
74 Under the prevalent pre-training and fine-tuning paradigm, when scaling up pre-training datasets,
75 large ViTs usually exhibit obviously superior performance scalability than top-performing CNNs
76 and MLPs in terms of fine-tuning on both upstream image classification tasks and downstream dense
77 prediction tasks [5, 39, 40]. As for knowledge distillation in this work, we assume that pre-training
78 datasets are no longer accessible and only well pre-trained ViT teacher models are available, avoiding
79 expensive pre-training costs and making the setting suited for real applications. When training student
80 models on upstream image classification datasets like ImageNet-1K, the knowledge density between
81 teacher and student models is different, which appears as the third barrier to handle.

82 From the above analysis, we can conclude that the design of effective schemes to align (1) feature
83 computing paradigm differences, (2) model scale differences, and (3) knowledge density differences
84 between the pre-trained ViT teacher and target student models, plays the key role to attain our goal.
85 Accordingly, we present scalable knowledge distillation (ScaleKD), a simple and effective cross
86 architecture KD method, which addresses them in a progressive manner. Regarding the first alignment
87 problem, we notice that some previous works leverage cross attention mechanisms to align different
88 modalities [41–43]. Extending this philosophy to cross architecture distillation, we propose our *cross*
89 *attention projector (CAP)* (shown in Figure 1(a)), bridging the gaps across different semantic units
90 and core operations. For one, CAP utilizes positional embedding and patchify stem to transform the
91 semantic units of CNN and MLP into transformer-like tokens. For another, by employing trainable
92 queries that share the same attributes as the teacher’s features, the student’s features can be globally
93 queried, modeling interdependencies based on positional information. In this way, CAP aligns
94 feature computing differences in form and serves as the base component of the other two designs.
95 Regarding the second and third alignment problems, as the high model capacity and the pre-training
96 are inextricably bound up with each other for ViTs, we think they co-act on models’ feature space
97 and parameter space. Surprisingly, we observe two interesting phenomena related to the two spaces:

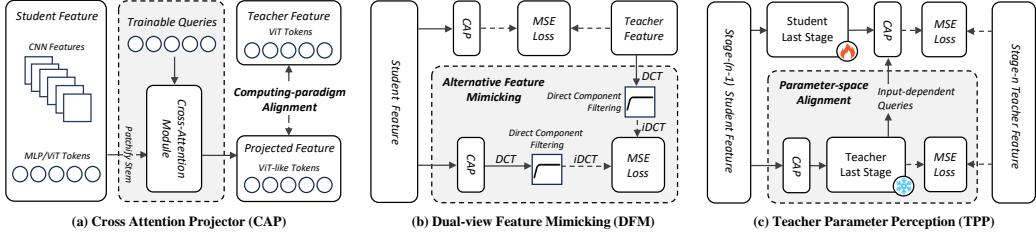


Figure 1: Overview of three core components in our ScaleKD, which are cross attention projector (left), dual-view feature mimicking (middle), and teacher parameter perception (right). Notably, the teacher’s parameters are frozen in the training stage,

- As shown in Figure 2, the frequency distributions of the features of the pre-trained ViTs are extremely imbalanced, where the direct component (zero frequency) response is dominant among all frequencies. Conducting feature distillation under such an imbalanced distribution may neglect the features of all other alternative components.
- As the parameters in the fine-tuning stage are slightly changed, the pre-training knowledge remains in the parameter space of the pre-trained ViTs. Although the pre-training datasets are invisible, the student still has the opportunity to obtain the knowledge by aligning its parameter space to the teacher’s.

Inspired by these two insightful observations, we try to tackle the problems from two new perspectives. Based on the first observation, we design *dual-view feature mimicking (DFM)* (shown in Figure 1(b)), whose key insight is to complement the neglected alternative features in KD process. Specifically, DFM employs CAP as the feature projector and incorporates two feature mimicking paths. In the first path, DFM conducts feature mimicking in the original space to learn the teacher’s global features. In the second path, by removing the direct component in frequency space, DFM highlights the subtle alternative responses in feature mimicking, thus avoiding the neglect of these features. As a result, the two paths are complementary to each other, jointly promoting the feature space alignment. Based on the second observation, we propose *teacher parameter perception (TPP)* (shown in Figure 1(c)), whose target is to align the teacher’s and student’s parameter spaces. Thanks to the aligned feature computing paradigm by CAP, TPP could bridge the student’s early stages and the teacher’s later stages and form a proxy feature processing path. By applying feature distillation in this path, the student’s parameter space would be gradually aligned to the teacher’s. In this way, the pre-training knowledge could be transferred from the teacher to the student. Notably, the distillation learning processes in feature space and parameter space are the two sides of the same coin, indicating that DFM and TPP are mutually promoted by each other. Thanks to progressive designs, CAP, DFM, and TPP can be seamlessly integrated into a neat and effective cross-architecture knowledge distillation method, called *ScaleKD*, which addresses the above three problems simultaneously. Although ScaleKD has multiple feature mimicking paths, they only exist in the training stage, which means that ScaleKD does not alter the student’s structure and introduces no additional cost in the inference stage. By conducting systematic experiments on several large mainstream benchmarks, we validate the effectiveness and generalization ability of our method.

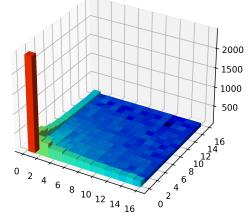
2 Scalable Knowledge Distillation

Given a pre-trained ViT teacher having m stages and a target student (CNN or MLP or ViT) having n stages, let F^{s_i} and F^{t_j} denote features from i -th stage of the student and j -th of the teacher, respectively. In what follows, we formulate all components in the form of performing feature distillation, for better clarifying the relationships of three coupled components in our method.

2.1 Cross Attention Projector

As shown in Figure 1(a), CAP adopts the structure of the standard transformer decoder block, consisting of a transformer decoder layer and an MLP layer, but incorporates three critical modifications: i) patchifying regular grids of pixels in CNN; ii) adding

Figure 2: Feature distribution of BEiT-L [40] in the frequency domain, where the direct component response is dominant. Details on drawing this figure are shown in Figure 5.



140 positional embedding; iii) setting queries in the transformer decoder block as trainable variables that
 141 share the same resolution with the teacher's features. The first two modifications intend to narrow
 142 the discrepancy between different semantic units, while the last modification grants the employed
 143 transformer decoder block great flexibility in aligning feature semantics and spatial resolution. Based
 144 on these modifications, the cross attention operation further provides global dependencies to the
 145 projected features. With the projected features by CAP, the feature distillation loss is defined as:¹:

$$\mathcal{L}_{CAP} = \alpha L(F^t, f_p(F^s; q)) = \alpha \|F^t - f_p(F^s; q)\|_2^2, \quad (1)$$

146 where $f_p, q, \alpha (\geq 0)$ and $L(\cdot)$ denote the CAP and the trainable queries, the loss weight, and the
 147 L_2 -normed distance, respectively.

148 2.2 Dual-view Feature Mimicking

149 As shown in Figure 1(b), building upon CAP, DFM contains two feature mimicking paths. As
 150 we stated in the Introduction section, the first path aims to learn the teacher's global features and
 151 the second path aims to mimic the neglected alternative features. Specifically, in the first path,
 152 DFM conducts feature mimicking in the teacher's original feature space, which is formulated
 153 as: $\mathcal{L}_{ori} = \alpha L(F^t, f_{p1}(F^s))$, where f_{p1} denotes the CAP in the first path. In the second path,
 154 the dominant direct component should be removed. To achieve this goal, we employ Discrete
 155 Cosine Transform (DCT), which maps the feature from the spatial domain to the frequency domain:
 156 $DCT : \mathcal{X} \rightarrow \mathcal{Z}$. We then define an operator ϕ that removes direct component response from the
 157 features:

$$\phi(x) = DCT^{-1}(\sigma(DCT(x))) \quad s.t. \quad \sigma(z) = \begin{cases} 0, & z = 0 \\ z, & z \neq 0 \end{cases}. \quad (2)$$

158 Next, feature mimicking in the second path is formulated as: $\mathcal{L}_{alt} = \alpha L(\phi(F^t), \phi(f_{p2}(F^s)))$, where
 159 f_{p2} denotes the CAP in the second path. Now, the feature distillation loss of DFM is formulated as:

$$\mathcal{L}_{DFM} = \beta \mathcal{L}_{ori} + (1 - \beta) \mathcal{L}_{alt}, \quad (3)$$

160 where $\beta \in [0, 1]$ denotes the balancing weight.

161 2.3 Teacher Parameter Perception

162 As we stated in the Introduction section, TPP establishes a proxy feature processing path by connecting
 163 the student's early stages with the teacher's later stages by a CAP. In our implementation, this proxy
 164 path consists of the student's first $n-1$ stages and the teacher's last stage, as illustrated in Figure
 165 1(c). By conducting feature mimicking in this path, the parameters of the student part are gradually
 166 aligned with the parameters of the teacher part, thus enabling the transfer of the teacher's pre-training
 167 knowledge. The output features of this path are $F^{st} = g_{t_m}(f_p^{st}(F^{s_{n-1}}; q))$, where g_{t_m} and f_p^{st}
 168 denote the teacher's last stage and CAP in the proxy path, respectively. The feature mimicking in
 169 the proxy path is formulated as: $\mathcal{L}^{st} = \alpha L(F^t, F^{st})$. We further introduce F^{st} as input-dependent
 170 queries for the CAP in the original path. This design aims to enhance the capability of CAP, as such
 171 queries contain more teacher-related information. Under such condition, the feature mimicking loss is
 172 changed to $\mathcal{L}^s = \alpha L(F^t, f_p^s(F^{s_n}; F^{st}))$. As the two paths are treated equally, the feature distillation
 173 loss of TPP is defined as:

$$\mathcal{L}^{TPP} = \mathcal{L}^s + \mathcal{L}^{st}. \quad (4)$$

174 2.4 Formulation of ScaleKD

175 From a general perspective, the progressive designs of our above three components are naturally
 176 cohesive. As CAP serves as the basic component in DFM and TPP, we further introduce how to apply
 177 DFM in TPP and thus get a neat formulation of our method, ScaleKD. Specifically, if treating DFM as
 178 an improved version of traditional feature mimicking, it can substitute the original feature mimicking
 179 in each path of TPP. In this way, we present our overall design - ScaleKD, which is formulated as:

$$\mathcal{L}_{ScaleKD} = \mathcal{L}_{task} + \underbrace{\beta \mathcal{L}_{ori}^s + (1 - \beta) \mathcal{L}_{alt}^s}_{DFM \text{ for TPP Student Path}} + \underbrace{\beta \mathcal{L}_{ori}^{st} + (1 - \beta) \mathcal{L}_{alt}^{st}}_{DFM \text{ for TPP Teacher Path}} + \mathcal{L}_{kd}, \quad (5)$$

¹We normalize each token before calculating distillation loss.

Table 1: Pilot experiments on cross architecture distillation with ScaleKD and FD. s_i denotes the distillation is conducted on stage-i. To clearly show the performance gain, experiments in this table are conducted without L_{kd} .

Teacher	Student	Method	Top-1(%)	Δ Top-1(%)
Swin-S (83.02)	ResNet-50	Baseline	76.55	-
		FD (s_4)	77.43	+0.88
		FD (s_3, s_4)	77.74	+1.19
		ScaleKD	79.30	+2.75
	Mixer-S	Baseline	74.02	-
		FD (s_4)	74.88	+0.86
		FD (s_3, s_4)	75.32	+1.30
		ScaleKD	77.24	+3.22

Table 2: Pilot experiments on scaling up the teachers. The advanced training strategy has more advanced data augmentation, optimizer, and longer training epochs, as shown in Table 10.

Teacher	Student	Ratio T/S Params	Top-1(%)	Δ Top-1(%)
<i>Traditional Training Strategy</i>				
Baseline		-	76.55	-
Swin-S (83.02)	ResNet-50	1.94 \times	79.62	+3.07
Swin-B (85.16)		3.43 \times	79.80	+3.25
Swin-L (86.24)		7.68 \times	80.10	+3.55
<i>Advanced Training Strategy</i>				
Baseline		-	78.64	-
Swin-S (83.02)	ResNet-50	1.94 \times	81.43	+2.79
Swin-B (85.16)		3.43 \times	81.77	+3.13
Swin-L (86.24)		7.68 \times	82.03	+3.39

180 where $\beta \in [0, 1]$ indicates the balancing weight and \mathcal{L}_{task} is the cross-entropy loss. Besides,
 181 following previous KD research, we use the vanilla logits-based KD loss [14]. As the features are
 182 standardized, we set $\alpha = 1$ as the default. Thus, we only have one hyper-parameter β in our method.

183 3 Experiments

184 3.1 Datasets and General Experimental Setups

185 We systematically validate the efficacy of our method under different setups: i) training various
 186 student backbones by our ScaleKD on ImageNet-1K [44] (IN-1K) and showing performance gain
 187 compared to from-scratch counterparts; ii) conducting transfer learning on downstream tasks on
 188 MS-COCO [45] and ADE20K [46] to examine whether the performance gains from our method could
 189 be well preserved; iii) showcase linear probing on CIFAR-100 [44] to explore the generalization
 190 ability of our method across different data distributions. *More details are shown in Appendix A and B.*

191 3.2 Experiments on Image Classification Task

192 To obtain pre-training knowledge from large pre-trained ViTs, we design a series of experiments to
 193 extensively study the effectiveness and potentials of our method step by step.

194 Step 1: Pilot Experiments on Our Basic Settings

195 **Cross Architecture Knowledge Distillation.** To illustrate the difficulty of cross architecture feature
 196 distillation and validate the efficacy of ScaleKD under this setting, we compare our ScaleKD with
 197 traditional feature distillation (FD) [15] on two different cross architecture teacher-student network
 198 pairs. From the results shown in Table 1, we can observe: i) Due to architecture gaps between the
 199 teacher and the student, traditional FD shows limited gains; ii) Comparatively, our ScaleKD achieves
 200 significantly better performance, bringing 2.75%|3.22% absolute top-1 gain for ResNet-50|Mixer-S.
 201 By now, we preliminarily verify that ScaleKD could effectively handle cross architecture feature
 202 distillation, which is difficult for traditional FD.

203 **Large Pre-trained ViTs as Teachers.** After showcasing the effectiveness of ScaleKD in cross
 204 architecture KD, we further examine the rationality of selecting large pre-trained ViT as teachers.
 205 Specifically, we gradually scale up the teacher’s model size and capability under two different training
 206 strategies. From the results shown in Table 2, we can observe: *i) ScaleKD can adapt large-scale*
 207 *teachers.* It can be seen that the performance gain is constantly increased under both training strategies
 208 when scaling up the teacher’s model size; *ii) ScaleKD can adapt teachers’ training strategies.* It
 209 can be seen that our ScaleKD still brings comparable performance gain under the advanced training
 210 strategy, although the baseline performance is relatively high.

211 **Conclusion.** According to the above pilot experiments, *our ScaleKD shows basic capabilities on*
 212 *handling cross architecture distillation from large pre-trained ViTs to CNN and MLP students.*

213 Step 2: Main Experiments on Various Teacher-Student Network Pairs

214 Based on the conclusion of the first step, we move forward and perform extensive experiments
 215 to examine the generalization ability of ScaleKD, considering more teacher-student network pairs.
 216 Specifically, we design **eleven** teacher-student network pairs, by choosing two large teachers and **ten**
 217 popular models for students, covering mainstream architectures across ViT, MLP, and CNN.

Table 3: Main results of ScaleKD on eleven teacher-student network pairs. \dagger denotes the model pre-trained on IN-22K [44] and \ddagger denotes the model pre-trained by EVA [40], an MIM distillation on IN-22K from CLIP-based ViT. Underlined results are for the experiments performed by us.

Teacher	Student	Params (M)		FLOPs (G)		Accuracy (%)	
		T	S	T	S	Top-1	Δ Top-1
Swin-L \dagger (86.24)	MobileNet-V1 (72.10)		4.23		0.58	75.15	+3.05
	ResNet-50 (78.64)	196.53	25.56	34.04	4.12	82.03	+3.39
	ConvNeXt-T (82.14)		28.59		4.46	84.16	+2.02
	Mixer-S/16 (74.02)	196.53	18.53	3.78	78.63	+4.61	
	Mixer-B/16 (76.44)		59.88	12.61	81.96	+5.52	
	ViT-S/16 (79.90)		22.05		4.61	83.93	+4.03
	Swin-T (81.18)	196.53	28.29	34.04	4.36	83.80	+2.62
	ViT-B/16 (81.80)		86.57		17.58	85.53	+3.73
	ResNet-50 (78.64)		25.56		4.12	82.34	+3.70
BEiT-L/14 \ddagger (88.58)	Mixer-B/14 (76.62)	304.14	59.88	81.06	16.45	82.89	+6.27
	ViT-B/14 (82.02)			86.57	23.09	86.43	+4.41

Table 4: Performance comparisons with recent top KD methods. Following their settings, the students are trained under the advanced training strategy. Best results are **bolded**.

Model	Method	Teacher	# Epochs	Top-1(%)
Swin-T	From Scratch	-	300	81.18
	DIST [47]	Swin-L (86.30)	300	82.30
	DiffKD [48]	Swin-L (86.30)	300	82.50
	ScaleKD	Swin-L (86.24)	300	83.80
ResNet-50	From Scratch	-	300	78.60
	DIST[47]	Swin-L (86.30)	450	80.20
	DiffKD [48]	Swin-L (86.30)	450	80.50
	OFA [36]	ViT-B (86.53)	300	81.33
	ScaleKD	Swin-L (86.24)	300	82.03
	FunMatch [49]	BiT-Res152x2 (N/A)	1200	81.54
	FunMatch [49]	BiT-Res152x2 (N/A)	9600	82.31
	ScaleKD	Swin-L (86.24)	600	82.55

Table 5: Performance comparisons with various models. More comparisons are shown in Table 15 in the Appendix.

Model	Params (M)	FLOPs (G)	Top-1(%)
<i>CNN-based Architecture</i>			
ResNet-50 [2]	22.56	4.12	78.64
ResNet-50 + ScaleKD	22.56	4.12	82.55
<i>MLP-based Architecture</i>			
Mixer-B/16 [6]	59.88	12.61	76.44
Mixer-B/16 + ScaleKD	59.88	12.61	81.96
gMLP-B [8]	73.00	15.80	81.60
ResMLP-B24 [7]	115.7	23.00	81.00
<i>ViT-based Architecture</i>			
ViT-S/16 [4]	22.05	4.61	79.90
ViT-S/16 + ScaleKD	22.05	4.61	83.93
Swin-T [5]	73.00	15.80	81.18
Swin-B [5]	87.77	15.14	83.50

218 **Experimental Results.** From the results shown in Table 3, we observe: i) In general, our ScaleKD
219 shows great generalization ability across various teacher-student network pairs, no matter for CNN,
220 MLP and ViT students. Over eleven teacher-student pairs, the mean top-1 accuracy improvement
221 reaches 3.94%, and the maximal gain is 6.27%; ii) Considering the acceleration, with Swin-L
222 as the teacher, ResNet-50|Mixer-S/16|ViT-S/16 trained by ScaleKD even outperforms individually
223 trained ResNet-152|Mixer-B/16|ViT-B/16 by a margin of 0.28%|2.19%|2.13%, achieving over
224 $2.35 \times | 3.23 \times | 3.83 \times$ compression on model parameter size; iii) The absolute performance gain could
225 still be increased, when choosing stronger teachers. For instance, ScaleKD bring 0.32% additional
226 gain to ResNet-50 when changing the teacher from Swin-L to BEiT-L/14.

227 Step 3: Performance Comparisons with Three Different Families of Methods

228 In Step 2, we have verified our ScaleKD could help various student models obtain large performance
229 gains. To further evaluate the value of these gains, we make three kinds of comparisons.

230 **Comparisons with KD Methods.** We first compare our ScaleKD with various top KD methods,
231 to demonstrate the superiority of our method. From the results shown in Table 4, we can see that:
232 i) Compared to DIST, DiffKD and OFA, although our teacher is not the best and the number of
233 training epochs is not the largest, our ScaleKD still outperforms all these methods by clear margins
234 (0.70%|1.30% on ResNet-50|Swin-T); ii) Comparing to FunMatch, our cross architecture distillation
235 is even superior, outperforming FunMatch by a margin of 0.24% but only using less than 1/10 training
236 epochs.

Table 6: Performance comparisons with various pre-training methods. \Rightarrow denotes transfer learning and * denotes the model is trained and tested with 384×384 resolution samples.

Model	Method	Training Dataset	Dataset Samples \times Epochs (M)	Viewed Samples (M)	Top-1(%)
<i>Supervised pre-training</i>					
ViT-B/16	Pre-training [4]	IN-22K \Rightarrow IN-1K JFT-300 \Rightarrow IN-1K	13.7 \times 90 + 1.28 \times 32 300 \times 7 + 1.28 \times 32	1274 2141	83.97 84.15
	ScaleKD	IN-1K	1.28 \times 300	384	85.53
<i>Self-supervised pre-training</i>					
ViT-B/16	BEiT [39] iBOT [11]	IN-22K \Rightarrow IN-1K IN-22K \Rightarrow IN-1K	13.7 \times 150 + 1.28 \times 100 13.7 \times 320 + 1.28 \times 100	2183 4512	83.70 84.40
	ScaleKD	IN-1K	1.28 \times 300	384	85.64
<i>Cross-modal pre-training</i>					
ViT-B/16	CLIP [13]	LAION-2B \Rightarrow IN-1K LAION-2B \Rightarrow IN-12K \Rightarrow IN-1K	2320 \times 32 + 1.28 \times 50 2320 \times 32 + 12.1 \times 60 + 1.28 \times 50	74304 75030	85.47 86.17
	ScaleKD	IN-1K	1.28 \times 300	384	86.19
<i>EVA Hybrid pre-training (MIM distillation from the cross-modal pre-trained teacher)</i>					
EVA02-S/14*	EVA-02 [51]	IN-22K \Rightarrow IN-1K	13.7 \times 240 + 1.28 \times 50	3352	85.80
	ScaleKD	IN-1K	1.28 \times 300	384	86.22

Table 7: Transfer learning performance on MS-COCO.

Framework	Backbone	Pre-training	Classification (IN-1K) Top-1(%)	Object Detection (COCO)				Instance Segmentation (COCO)			
				AP	AP_S	AP_M	AP_L	AP	AP_S	AP_M	AP_L
Mask R-CNN	ResNet-50	Baseline Ours	78.64 82.03 (+3.39)	40.2 42.3	23.0 25.5	44.3 46.5	52.5 54.6	37.1 39.1	18.0 19.3	40.1 42.5	54.9 57.1
	Swin-T	Baseline Ours	81.18 83.80 (+2.62)	42.7 44.4	26.5 28.7	45.9 47.9	56.6 58.6	39.3 40.8	20.5 21.8	41.8 43.7	57.8 59.8

237 **Comparisons with Model Designs.** To illustrate the value of our method to model performance, we
238 examine whether our ScaleKD could bring comparable performance gains as network architecture
239 engineering. We apply ScaleKD to the basic design of each architecture and make comparisons with
240 recent advanced designs. From the results shown in Table 5, we observe that: i) The performance
241 gains from our ScaleKD could help the basic designs reach similar or even higher performance than
242 advanced models. ii) For some basic models that suffer from the data-hungry problem (Mixer-B) or
243 low convergence speed (ViT-S), our ScaleKD could even help them exceed larger advanced models.

244 **Comparisons with Pre-training Methods.** We finally examine the value of our method to the
245 training regime. We compare the performance between models trained by ScaleKD and models
246 pre-trained by various methods on upstream datasets, to figure out whether our ScaleKD could be
247 an efficient substitution to pre-training of any target student on large-scale datasets if we have a
248 strong pre-trained ViT. For making comprehensive comparisons, we select the methods across various
249 pre-training paradigms, such as supervised pre-training, self-supervised pre-training, cross-modal
250 pre-training, and hybrid pre-training. From the results shown in Table 6, we can summarize: i) Our
251 ScaleKD has superior performance than various methods across all kinds of pre-training; ii) Our
252 ScaleKD is a more efficient approach, viewing less than $5.58 \times 11.75 \times 195.39 \times 8.73 \times$ samples than
253 these counterpart methods based on different pre-training paradigms.

254 3.3 Experiments on Downstream Tasks

255 To further examine whether the performance
256 gains from our method could be well pre-
257 served in transfer learning, we conduct com-
258 parative experiments on MS-COCO for ob-
259 ject detection and instance segmentation, and
260 on ADE20K for semantic segmentation.

Table 8: Transfer learning performance on ADE20K

Framework	Backbone	Pre-training	IN-1K (Top-1)	ADE20K (mIOU)
UpNet	ResNet-50	Baseline Ours	78.64 82.03 (+3.39)	42.37 44.50 (+2.13)
	Swin-T	Baseline Ours	81.18 83.80 (+2.62)	44.41 46.33 (+1.92)
	ViT-B/16	Baseline Ours	81.80 85.53 (+3.73)	46.75 50.84 (+4.09)

261 **Experimental Results.** The results on MS-COCO and ADE20K are shown in Table 7 and Table 8,
262 respectively, from which we can observe that: i) Overall, our pre-trained models outperform baselines
263 by significant margins across three tasks and different architectures, demonstrating the performance
264 gains from ScaleKD could be transferred to downstream tasks; ii) For semantic segmentation tasks,

Table 9: Ablation studies. Experiments in (b)-(d) are performed on Swin-S→ResNet-50. As DFM and TPP are designed based on CAP, CAP is added by default when choosing DFM and TPP in (a). Because of this, we treat CAP as another baseline method, when analyzing DFM and TPP in (c)-(d).

(a) Ablation on the overall design			(b) Ablation on CAP			(c) Ablation on DFM		
Teacher	Student		CAP	DFM	TPP	KD	Top-1	
Swin-S	ResNet-50						76.55	
		✓					77.87	
		✓	✓				78.51	
		✓		✓			78.62	
		✓	✓	✓			79.30	
		✓	✓	✓	✓	✓	79.62	
Swin-S	Mixer-S						74.02	
		✓					75.03	
		✓	✓				76.42	
		✓		✓			76.28	
		✓	✓	✓			77.24	
		✓	✓	✓	✓	✓	77.59	
(d) Ablation on TPP								
Method	TPP Design			Accuracy (%)			Top-1	Δ Top-1
	Teacher's Param	Provide Queries		Top-1				
Baseline	-	-	-	76.55	-			
CAP	-	-	-	77.87	+1.32			
TPP				78.50	+1.95		78.62	+2.07
	✓	✓	✓	✓				

265 ViT-B/16 achieves the highest 4.09% absolute performance gain across three backbones, even higher
266 than its gain on IN-1K; iii) For object detection and instance segmentation on MS-COCO, ResNet-
267 50|Swin-T pre-trained by ScaleKD outperforms baseline on AP by a margin of 2.1%|1.7% and
268 2.0%|1.5%.

269 3.4 Ablation Study

270 **Ablation Study on the Overall Design.** To exploit the gains from different designs and further
271 examine whether our three components are complementary to each other, we perform this study.
272 From the results shown in Table 9a, we notice: i) When gradually applying more of the designs, the
273 performance of ResNet-50 and Mixer-S show similar increasing trends, showing that each component
274 of ScaleKD is not designed for specific student architecture; ii) Although CAP can bring the two
275 students promising performance gains, DFM and TPP can further bring ResNet-50|Mixer-S extra
276 performance improvements, 0.64%|0.75% and 1.25%|1.39%, respectively. This observation verifies
277 that DFM and TPP are complementary to CAP; iii) When using DFM and TPP together, both ResNet-
278 50 and Mixer-S can obtain additional performance boosts, which indicates that DFM and TPP are
279 complementary with each other.

280 **Ablation Study on CAP, DFM and TPP.** i) We compare CAP with two popular projectors, *Linear-*
281 *projector* and *Conv-projector*, to verify the superiority of CAP. The former projector consists of a
282 linear layer and the latter projector consists of two 3×3 convolutional layers with ReLU. From the
283 results shown in Table 9b, we can notice that our CAP outperforms them, which partially validates
284 that our CAP can boost the distillation performance by aligning the computing paradigm in form.
285 ii) We then investigate the effectiveness of two designs in DFM. From the results shown in Table
286 9c, we find that plain dual-path feature mimicking brings an extra 0.25% performance gain to CAP,
287 and direct component filtering further brings 0.39% performance improvement based on dual-path
288 feature mimicking. This experiment validates that the two designs are essential and complementary
289 to each other. iii) To explore the necessities of teacher parameter reusing and input-dependent queries
290 in TPP, we conduct the experiment in Table 9d. The results show that reusing teacher parameters can
291 provide students with performance improvement and providing input-dependent queries can further
292 enhance the effectiveness of TPP.

293 *More ablation studies on hyper-parameters and the training efficiency are provided in Appendix D.*

294 3.5 Verification Experiments and Discussion

295 Finally, we conduct experiments to study: i) Could ScaleKD help the student mimic the teacher's high-
296 frequency features? ii) Could ScaleKD help the student learn the teacher's pre-training knowledge?

297 **High-Frequency Feature Mimicking.** To clearly illustrate that our ScaleKD could help the student
298 capture the teacher's high-frequency features, we make comparisons between the method with DFM
299 and without DFM. Specifically, we measure the distances between the student's features and the

300 teacher’s features of alternative components in the spatial domain, as shown in Figure 3. Clearly,
301 methods with DFM can reduce high-frequency feature distances, answering the first question.

302 **Pre-training Knowledge Learning.** The pre-training knowl-
303 edge of teachers can be generally summarized into two parts: the
304 generalization ability within training data distribution and the
305 generalization ability across different data distributions. In Table
306 6, we have verified the capability of ScaleKD under the former
307 setting. To further examine the generalization across different
308 distributions, we conduct several linear probing experiments on
309 CIFAR-100, based on models in Table 6. From the results shown
310 in Table 12 in the Appendix, we can observe that: i) models
311 pre-trained by CLIP greatly improve the backbone’s generaliza-
312 tion ability across datasets; ii) our ScaleKD helps the student
313 model reach similar performance as CLIP-based pre-training,
314 without viewing pre-training data. The two observations prove
315 the gain from ScaleKD in this experiment is from the knowledge
316 of upstream datasets, thus answering the second question.

317 4 Related Work

318 **Knowledge Distillation.** Traditional Knowledge distillation methods [14–35] generally focus
319 on CNN-based teacher-student network pairs with small model scale gaps. Some recent works
320 [52, 53, 47] further study how to conduct knowledge distillation on larger teachers. As vision trans-
321 formers suffer from low convergence speeds, some recent works [54–56] explore leveraging CNNs
322 to accelerate the training of vision transformers. Meanwhile, [36–38] discuss how to bridge the
323 architecture gap when the teacher and the student are in different model categories.

324 **Frequency-based Knowledge Distillation.** As traditional feature distillation only focus on pixel-to-
325 pixel differences, FAM [57] designs frequency-based attention and conducts KD between attention
326 maps. FreeKD [58] explores how to eliminate unfavorable information in the frequency domain to
327 enhance the distillation performance on dense prediction tasks. Different from our ScaleKD, they
328 consider feature distillation on CNN-based network pairs and have different formulations.

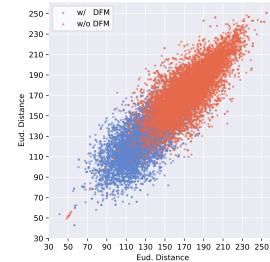
329 **Teacher Parameter Reuse.** Some previous KD methods also leverage the teacher’s parameter for
330 reusing a better classifier [32] and initializing the student’s neck and head [59, 60] or dismissing the
331 shortcuts in residual architectures [61]. Unlike our ScaleKD, the motivation of these works focuses on
332 parameter reuse or equivalent substitution, rather than aligning the two parameter spaces for pursuing
333 the teacher’s pre-training knowledge.

334 5 Conclusion

335 In this paper, we present ScaleKD, a new cross architecture KD approach for transferring the scalable
336 properties from pre-trained vision transformers to various CNNs and MLPs or heterogeneous ViTs.
337 Our method presents principled designs to address the difficulties of feature computing paradigm
338 differences, model scale differences, and knowledge density differences between teacher and student.
339 By conducting systematic experiments on various large mainstream benchmarks, we broadly validate
340 the effectiveness and generalization ability of our method.

341 **Limitations and Broader Impacts.** Restricted by our computational resources, we do not conduct
342 experiments on very large teachers, such as ViT-22B [62], or on large students, such as ViT-L [4].
343 Furthermore, with the increasing model scale of teachers, the training costs of ScaleKD increase,
344 which is a common limitation to KD research. According to the analysis in Appendix D, the
345 extra training costs of ScaleKD is acceptable in most cases, because its performance gain shows
346 the potential to replace time-intensive pre-training of students on large-scale datasets, thus being
347 significantly more efficient. For broader impacts, we think our research could be beneficial to many
348 real applications by promoting their model performance and training efficiency, which could also be
349 a reference for the research on KD with large vision models.

Figure 3: Feature distances of alterna-tional components in the spa-tial domain. Details on the figure drawing are in Figure 6



350 **References**

- 351 [1] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep
352 convolutional neural networks. In *NIPS*, 2012.
- 353 [2] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
354 recognition. In *CVPR*, 2016.
- 355 [3] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining
356 Xie. A convnet for the 2020s. In *CVPR*, 2022.
- 357 [4] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
358 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al.
359 An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2020.
- 360 [5] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining
361 Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, 2021.
- 362 [6] Ilya O Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas
363 Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, et al. Mlp-mixer:
364 An all-mlp architecture for vision. In *NeurIPS*, 2021.
- 365 [7] Hugo Touvron, Piotr Bojanowski, Mathilde Caron, Matthieu Cord, Alaaeldin El-Nouby,
366 Edouard Grave, Gautier Izacard, Armand Joulin, Gabriel Synnaeve, Jakob Verbeek, et al.
367 Resmlp: Feedforward networks for image classification with data-efficient training. *TPAMI*,
368 2022.
- 369 [8] Hanxiao Liu, Zihang Dai, David So, and Quoc V Le. Pay attention to mlps. In *NeurIPS*, 2021.
- 370 [9] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
371 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, 2017.
- 372 [10] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked
373 autoencoders are scalable vision learners. In *CVPR*, 2022.
- 374 [11] Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong.
375 Image bert pre-training with online tokenizer. In *ICLR*, 2021.
- 376 [12] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
377 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
378 models from natural language supervision. In *ICML*, 2021.
- 379 [13] Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade
380 Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws
381 for contrastive language-image learning. In *CVPR*, 2023.
- 382 [14] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network.
383 *arXiv preprint arXiv:1503.02531*, 2015.
- 384 [15] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and
385 Yoshua Bengio. Fitnets: Hints for thin deep nets. In *ICLR*, 2015.
- 386 [16] Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the
387 performance of convolutional neural networks via attention transfer. In *ICLR*, 2016.
- 388 [17] Sungsoo Ahn, Shell Xu Hu, Andreas Damianou, Neil D Lawrence, and Zhenwen Dai. Varia-
389 tional information distillation for knowledge transfer. In *CVPR*, 2019.
- 390 [18] Frederick Tung and Greg Mori. Similarity-preserving knowledge distillation. In *ICCV*, 2019.
- 391 [19] Baoyun Peng, Xiao Jin, Jiaheng Liu, Dongsheng Li, Yichao Wu, Yu Liu, Shunfeng Zhou, and
392 Zhaoning Zhang. Correlation congruence for knowledge distillation. In *ICCV*, 2019.
- 393 [20] Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. Relational knowledge distillation. In
394 *CVPR*, 2019.

- 395 [21] Nikolaos Passalis and Anastasios Tefas. Learning deep representations with probabilistic
396 knowledge transfer. In *ECCV*, 2018.
- 397 [22] Byeongho Heo, Minsik Lee, Sangdoo Yun, and Jin Young Choi. Knowledge transfer via
398 distillation of activation boundaries formed by hidden neurons. In *AAAI*, 2019.
- 399 [23] Jangho Kim, SeongUk Park, and Nojun Kwak. Paraphrasing complex network: Network
400 compression via factor transfer. In *NeurIPS*, 2018.
- 401 [24] Junho Yim, Donggyu Joo, Jihoon Bae, and Junmo Kim. A gift from knowledge distillation:
402 Fast optimization, network minimization and transfer learning. In *CVPR*, 2017.
- 403 [25] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive representation distillation. In
404 *ICLR*, 2020.
- 405 [26] Byeongho Heo, Jeesoo Kim, Sangdoo Yun, Hyojin Park, Nojun Kwak, and Jin Young Choi. A
406 comprehensive overhaul of feature distillation. In *ICCV*, 2019.
- 407 [27] Guodong Xu, Ziwei Liu, Xiaoxiao Li, and Chen Change Loy. Knowledge distillation meets
408 self-supervision. In *ECCV*, 2020.
- 409 [28] Xiatian Zhu, Shaogang Gong, et al. Knowledge distillation by on-the-fly native ensemble. In
410 *NeurIPS*, 2018.
- 411 [29] Guile Wu and Shaogang Gong. Peer collaborative learning for online knowledge distillation. In
412 *AAAI*, 2021.
- 413 [30] Jing Yang, Brais Marinez, Bulat Adrian, and Tzimiropoulos Georgios. Knowledge distillation
414 via softmax regression representation learning. In *ICLR*, 2021.
- 415 [31] Defang Chen, Jian-Ping Mei, Yuan Zhang, Can Wang, Zhe Wang, Yan Feng, and Chun Chen.
416 Cross-layer distillation with semantic calibration. In *AAAI*, 2021.
- 417 [32] Defang Chen, Jian-Ping Mei, Hailin Zhang, Can Wang, Yan Feng, and Chun Chen. Knowledge
418 distillation with the reused teacher classifier. In *CVPR*, 2022.
- 419 [33] Xueqing Deng, Dawei Sun, Shawn Newsam, and Peng Wang. Distpro: Searching a fast
420 knowledge distillation process via meta optimization. In *ECCV*, 2022.
- 421 [34] Xiaolong Liu, LUJUN LI, Chao Li, and Anbang Yao. Norm: Knowledge distillation via
422 n-to-one representation matching. In *ICLR*, 2022.
- 423 [35] Borui Zhao, Quan Cui, Renjie Song, Yiyu Qiu, and Jiajun Liang. Decoupled knowledge
424 distillation. In *CVPR*, 2022.
- 425 [36] Zhiwei Hao, Jianyuan Guo, Kai Han, Yehui Tang, Han Hu, Yunhe Wang, and Chang Xu.
426 One-for-all: Bridge the gap between heterogeneous architectures in knowledge distillation. In
427 *NeurIPS*, 2023.
- 428 [37] Yufan Liu, Jiajiong Cao, Bing Li, Weiming Hu, Jingting Ding, and Liang Li. Cross-architecture
429 knowledge distillation. In *ACCV*, 2022.
- 430 [38] Yixuan Wei, Han Hu, Zhenda Xie, Zheng Zhang, Yue Cao, Jianmin Bao, Dong Chen, and
431 Baining Guo. Contrastive learning rivals masked image modeling in fine-tuning via feature
432 distillation. *arXiv preprint arXiv:2205.14141*, 2022.
- 433 [39] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image trans-
434 formers. In *ICLR*, 2021.
- 435 [40] Yuxin Fang, Wen Wang, Binhui Xie, Quan Sun, Ledell Wu, Xinggang Wang, Tiejun Huang,
436 Xinlong Wang, and Yue Cao. Eva: Exploring the limits of masked visual representation learning
437 at scale. In *CVPR*, 2023.
- 438 [41] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image
439 diffusion models. In *ICCV*, 2023.

- 440 [42] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson,
441 Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In
442 *CVPR*, 2023.
- 443 [43] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer.
444 High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022.
- 445 [44] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
446 hierarchical image database. In *CVPR*, 2009.
- 447 [45] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan,
448 Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*,
449 2014.
- 450 [46] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio
451 Torralba. Semantic understanding of scenes through the ade20k dataset. *IJCV*, 2019.
- 452 [47] Tao Huang, Shan You, Fei Wang, Chen Qian, and Chang Xu. Knowledge distillation from a
453 stronger teacher. In *NeurIPS*, 2022.
- 454 [48] Tao Huang, Yuan Zhang, Mingkai Zheng, Shan You, Fei Wang, Chen Qian, and Chang Xu.
455 Knowledge diffusion for distillation. In *NeurIPS*, 2023.
- 456 [49] Lucas Beyer, Xiaohua Zhai, Amélie Royer, Larisa Markeeva, Rohan Anil, and Alexander
457 Kolesnikov. Knowledge distillation: A good teacher is patient and consistent. In *CVPR*, 2022.
- 458 [50] Ao Wang, Hui Chen, Zijia Lin, Hengjun Pu, and Guiguang Ding. Repvit: Revisiting mobile
459 cnn from vit perspective. *arXiv preprint arXiv:2307.09283*, 2023.
- 460 [51] Yuxin Fang, Quan Sun, Xinggang Wang, Tiejun Huang, Xinlong Wang, and Yue Cao. Eva-02:
461 A visual representation for neon genesis. *arXiv preprint arXiv:2303.11331*, 2023.
- 462 [52] Seyed Iman Mirzadeh, Mehrdad Farajtabar, Ang Li, Nir Levine, Akihiro Matsukawa, and
463 Hassan Ghasemzadeh. Improved knowledge distillation via teacher assistant. In *AAAI*, 2020.
- 464 [53] Wonchul Son, Jaemin Na, Junyoung Choi, and Wonjun Hwang. Densely guided knowledge
465 distillation using multiple teacher assistants. In *ICCV*, 2021.
- 466 [54] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and
467 Hervé Jégou. Training data-efficient image transformers & distillation through attention. In
468 *International conference on machine learning*, 2021.
- 469 [55] Xianing Chen, Qiong Cao, Yujie Zhong, Jing Zhang, Shenghua Gao, and Dacheng Tao. Dearkd:
470 Data-efficient early knowledge distillation for vision transformers. In *CVPR*, 2022.
- 471 [56] Borui Zhao, Renjie Song, and Jiajun Liang. Cumulative spatial knowledge distillation for vision
472 transformers. In *ICCV*, 2023.
- 473 [57] Cuong Pham, Van-Anh Nguyen, Trung Le, Dinh Phung, Gustavo Carneiro, and Thanh-Toan Do.
474 Frequency attention for knowledge distillation. In *WACV*, 2024.
- 475 [58] Yuan Zhang, Tao Huang, Jiaming Liu, Tao Jiang, Kuan Cheng, and Shanghang Zhang. Freekd:
476 Knowledge distillation via semantic frequency prompt. In *CVPR*, 2024.
- 477 [59] Zijian Kang, Peizhen Zhang, Xiangyu Zhang, Jian Sun, and Nanning Zheng. Instance-
478 conditional knowledge distillation for object detection. In *NeurIPS*, 2021.
- 479 [60] Zhendong Yang, Zhe Li, Xiaohu Jiang, Yuan Gong, Zehuan Yuan, Danpei Zhao, and Chun
480 Yuan. Focal and global knowledge distillation for detectors. In *CVPR*, 2022.
- 481 [61] Guilin Li, Junlei Zhang, Yunhe Wang, Chuanjian Liu, Matthias Tan, Yunfeng Lin, Wei Zhang,
482 Jiashi Feng, and Tong Zhang. Residual distillation: Towards portable deep neural networks
483 without shortcuts. In *NeurIPS*, 2020.

- 484 [62] Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin
 485 Gilmer, Andreas Peter Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, et al.
 486 Scaling vision transformers to 22 billion parameters. In *ICML*, 2023.
- 487 [63] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
 488 *Tech Report*, 2009.
- 489 [64] MMPreTrain Contributors. Openmmlab’s pre-training toolbox and benchmark. <https://github.com/open-mmlab/mmpretrain>, 2023.
- 490 [65] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias
 491 Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural
 492 networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- 493 [66] Pengguang Chen, Shu Liu, Hengshuang Zhao, and Jiaya Jia. Distilling knowledge via knowledge
 494 review. In *CVPR*, 2021.
- 495 [67] Zhendong Yang, Zhe Li, Mingqi Shao, Dachuan Shi, Zehuan Yuan, and Chun Yuan. Masked
 496 generative distillation. In *ECCV*, 2022.
- 497 [68] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun,
 498 Wansen Feng, Ziwei Liu, Jiarui Xu, Zheng Zhang, Dazhi Cheng, Chenchen Zhu, Tianheng
 499 Cheng, Qijie Zhao, Buyu Li, Xin Lu, Rui Zhu, Yue Wu, Jifeng Dai, Jingdong Wang, Jianping
 500 Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. MMDetection: Open mmlab detection
 501 toolbox and benchmark. *arXiv preprint arXiv:1906.07155*, 2019.
- 502 [69] MMSegmentation Contributors. MMSegmentation: Openmmlab semantic segmentation toolbox
 503 and benchmark. <https://github.com/open-mmlab/mmsegmentation>, 2020.
- 504 [70] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *ICCV*, 2017.
- 505 [71] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing
 506 for scene understanding. In *ECCV*, 2018.
- 507 [72] Xiaohan Ding, Yiyuan Zhang, Yixiao Ge, Sijie Zhao, Lin Song, Xiangyu Yue, and Ying
 508 Shan. Unireplknet: A universal perception large-kernel convnet for audio, video, point cloud,
 509 time-series and image recognition. *arXiv preprint arXiv:2311.15599*, 2023.
- 510 [73] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural
 511 networks. In *ICML*, 2019.
- 512 [74] Kan Wu, Jinnian Zhang, Houwen Peng, Mengchen Liu, Bin Xiao, Jianlong Fu, and Lu Yuan.
 513 Tinyvit: Fast pretraining distillation for small vision transformers. In *ECCV*, 2022.
- 514 [75] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zi-Hang Jiang, Francis EH Tay,
 515 Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch
 516 on imagenet. In *ICCV*, 2021.
- 517 [76] Mingyu Ding, Bin Xiao, Noel Codella, Ping Luo, Jingdong Wang, and Lu Yuan. Davit: Dual
 518 attention vision transformers. In *ECCV*, 2022.
- 519 [77] Jiahao Wang, Wenqi Shao, Mengzhao Chen, Chengyue Wu, Yong Liu, Kaipeng Zhang,
 520 Songyang Zhang, Kai Chen, and Ping Luo. Adapting llama decoder to vision transformer. *arXiv
 521 preprint arXiv:2404.06773*, 2024.
- 522

Appendix for "ScaleKD: Strong Vision Transformers Could Be Excellent Teachers"

Anonymous Author(s)

Affiliation
Address
email

523 **A Dataset**

524 **ImageNet-1K** [44] is a well-known large-scale classification dataset, comprising over 1.2 million
525 training images and 50,000 validation images with 1,000 object categories.

526 **MS-COCO** [45] is a large-scale dataset for object detection and instance segmentation, which
527 contains 118,000 training images and 5,000 validation images with 80 object categories.

528 **ADE20K** [46] is a challenging semantic segmentation dataset, containing 20,210 training samples,
529 2,000 validation samples, and 3,352 testing samples with 150 categories.

530 **CIFAR-100** [63] is a popular classification dataset, which contains 50,000 training images with 500
531 images per class and 10,000 test images.

532 **B Experimental Setups**

533 **B.1 Experimental Setups on IN-1K**

534 **Training Strategy.** We conduct our experiments with two training strategies: traditional training
535 strategy and advanced training strategy. The traditional training strategy is commonly used in previous
536 KD approaches (shown in Table 10a) and the advanced training strategy is adopted in training recently
537 proposed CNNs, MLPs, and ViTs (shown in Table 10b).

538 **Hardware.** The experiments using the traditional training strategy are conducted on $8 \times$ NVIDIA
539 Tesla-V100 GPUs, while the experiments using the advanced training strategy are conducted on $32 \times$
540 NVIDIA Tesla-V100 GPUs.

541 **Implementation Codebase.** We implement our method based on MMClassification [64].

542 **Hyper-parameter Settings.** The overall loss for our ScaleKD is defined as Formulation 5. Thanks to
543 the simplicity of our design, we have only one hyper-parameter β in our ScaleKD. From the ablation
544 study in the Appendix D, we find that the best choice is $\beta = 0.6$ and we use it as the default setting
545 throughout all experiments.

546 **Selection of Teacher-Student Network Pairs.** Overall, we choose eleven teacher-student network
547 pairs, which consist of two strong ViTs and ten students, covering mainstream architectures of ViT,
548 MLP, and CNN. Specifically, for the teacher, we choose two different types of strong ViTs: supervised
549 pre-trained Swin-L [5] with the hierarchical architecture and hybrid pre-trained BEiT-L [39] with the
550 typical ViT architecture. Moreover, compared to Swin-L, BEiT-L is much larger in terms of model
551 size and stronger in terms of model performance. For the student selections, we first choose the basic
552 design in each architecture, such as ResNet-50 [2], Mixer-S/16 [6], and ViT-S/16 [4]. Then, we also
553 select some popular models, such as MobileNet-V1 [65], ConvNeXt-T [3], and Swin-S. Next, we
554 expand the basic designs to larger ones, like Mixer-B/16, Mixer-B/14, ViT-B/16 and ViT-B/14. After

Table 10: Detailed settings of traditional training strategy and advanced training strategy on IN-1K.

(a) Traditional training strategy		(b) Advanced training strategy	
Configuration	CNN	Config	CNN / MLP / ViT
Batch Size	256	Batch Size	2048 / 1536 / 1024
Learning Rate	0.1	Learning Rate	5e-3 / 7e-4 / 1e-3
Learning Rate Schedule	Stepwise Decay/Cosine Decay	Learning Rate Schedule	Cosine Decay
Optimizer	SGD	Optimizer	Lamb / AdamW / AdamW
Optimizer Hyper-Parameters	momentum= 0.9	Optimizer Hyper-Parameters	$\beta_1, \beta_2, \epsilon = 0.9, 0.009, 1e-8$
Weight Decay	1e-4	Weight Decay	0.02 / 0.07 / 0.05
Training Epochs	100	Training Epochs	300
Warmup Epochs	x	Warmup Epochs	5 / 20 / 20
Drop Path	x	Drop Path	0.05 / 0.1 / 0.1
Label Smoothing	x	Label Smoothing	0.1
Random Flip	0.5	Random Flip	0.5
Random Resize Crop	(0.08,1)	Random Resize Crop	(0.08,1)
Random Augmentation	x	Random Augmentation	(7,0.5) / (9,0.5) / (9,0.5)
Random Erasing	x	Random Erasing	0.25

Table 11: Detailed settings of transfer learning strategies on MS-COCO and ADE20K.

(a) MS-COCO		(b) ADE20K	
Configuration	ResNet-50 / Swin-S	Configuration	ResNet-50 / Swin-S / ViT-B
Weight Initialization	Pre-trained Checkpoint	Weight Initialization	Pre-trained Checkpoint
Batch Size	16	Batch Size	16
Learning Rate	1e-4	Learning Rate Decay	1e-4 / 1e-4 / 2e-4
Learning Rate Decay	Stage (0.7)	LR decay	Stage (0.9) / Stage (0.9) / Layer (0.6)
Learning Rate Schedule	Cosine Decay	Learning Rate Schedule	Cosine Decay
Optimizer	AdamW	Optimizer	AdamW
Optimizer Hyper-Parameters	$\beta_1, \beta_2, \epsilon = 0.9, 0.009, 1e-8$	Optimizer Hyper-Parameters	$\beta_1, \beta_2, \epsilon = 0.9, 0.009, 1e-8$
Weight Decay	0.05	Weight Decay	0.05
Training Epochs	8	Training Iterations	160000
Crop Size	(1333, 800)	Crop Size	(512, 512)
Drop Path	0.0 / 0.2	Drop Path	0.0 / 0.3 / 0.2

555 separately selecting the teachers and students, we finally organize them into eleven teacher-student
556 network pairs for experiments.

557 **Counterpart Knowledge Distillation Methods.** In the main paper, we make comparisons with
558 many recent KD methods, such as DIST [47], DiffKD [48], OFA [36] and FuncMatch [49]. In the
559 appendix, we further compare with CNN-based methods, such as KD [14], AT [16], OFD [26], RKD
560 [20], CRD [25], DKD [35], SRRL [30], ReviewKD [66], DistPro [33] and MGD [67].

561 **Counterpart Model Designs.** To explicitly evaluate the significance of the performance gains,
562 we apply ScaleKD on the basic designs of each architecture and make comparisons with various
563 advanced counterparts. Driven by this target, we mainly select the next-generation model based on
564 the selected designs. For ResNet-50, we select ConvNeXt-T [3] and RepViT-2.3M [50] for making
565 comparisons. The former one is the typical design of the new-era CNN and the latter one is a
566 popular model for deployment. For Mixer-B, we select gMLP-B [8] and ResMLP-B24 [7], which are
567 optimized based on the weaknesses of MLP-Mixer. For ViT-S, we choose Swin-S [5] and Swin-B as
568 compared counterparts, to validate whether our ScaleKD could outperform larger advanced designs.

569 **Counterpart Pre-training Methods.** We select state-of-the-art methods in each paradigm. For
570 supervised pre-training, we directly compare with the results in the original paper [4]. For self-
571 supervised pre-training, we choose BEiT [39] and iBoT [11]. For cross-modal pre-training and hybrid
572 pre-training, we choose CLIP [13] and EVA-02 [51], respectively.

Table 12: Performance comparisons of ViT-B/16 between ScaleKD and CLIP for linear probing.

Model	Method	Pre-training Dataset	IN-1K (Training)	CIFAR-100 (Linear Probing)				
				1 shot	5 shot	10 shot	25 shot	Full
ViT-B/16	From-scratch	IN-1K	81.80	33.86	60.30	66.77	72.65	81.76
	CLIP	LAION-300M	-	-	-	71.96	77.21	84.07
		LAION-2B, IN-1K	85.49	41.10	69.00	72.34	78.64	85.51
		LAION-2B, IN-12K, IN-1K	86.17	44.90	70.19	76.77	81.43	88.88
		CLIP OpenAI, IN-12K, IN-1K	85.99	47.40	70.85	77.37	81.52	88.92
Teacher: BEiT-L/14	Ours	IN-1K	86.19	48.14	70.91	77.52	81.50	89.11
	EVA	CLIP OpenAI, IN-22K, IN-1K	88.58	63.74	85.20	87.39	89.27	93.36

Table 13: Experiments on the training efficiency of ScaleKD. The student model in all experiments is ResNet-50. In (a), we compare our ScaleKD with traditional FD on three teachers with different model scales. In (b), we conduct the experiments based on Swin-S→ResNet-50 teacher-student network pair to illustrate the training costs (memory and time) introduced by each component. Experiments are conducted on $8 \times$ NVIDIA Tesla-V100 GPUs

(a) Training costs comparisons with FD					(b) Training costs of each design in ScaleKD							
Teacher	Method	Top-1 (%)	GPU Memory (G)	T_{train} (d)	Method	CAP	DFM	Designs TPP	KD	Top-1 (%)	GPU Memory (G)	T_{train} (d)
Swin-S	FD	77.43	3.66	1.67	FD	-	-	-	-	77.43	3.66	1.67
	ScaleKD	79.62	6.65	2.10		✓	✓	✓	✓	77.87	3.77	1.70
Swin-B	FD	77.76	3.66	1.83	ScaleKD	✓	✓	✓	✓	78.51	4.02	1.77
	ScaleKD	79.80	7.26	2.53		✓	✓	✓	✓	78.62	5.13	1.84
Swin-L	FD	77.72	4.72	2.24		✓	✓	✓	✓	79.30	6.60	2.08
	ScaleKD	80.10	9.11	3.51		✓	✓	✓	✓	79.62	6.65	2.10

573 B.2 Experimental Setups on MS-COCO and ADE20K

574 **Training Strategy and Hyper-parameter Settings.** For the experiments on MS-COCO, we adopt
575 the settings shown in Table 11a, while for experiments on ADE20K, we adopt the settings shown in
576 Table 11b.

577 **Hardware.** All experiments on MS-COCO and ADE20K are conducted on $8 \times$ NVIDIA Tesla-V100
578 GPUs.

579 **Implementation Codebase.** We conduct experiments based on MMDetection [68] and MM Segmen-
580 tation [69].

581 **Framework and Backbone Selections.** For task frameworks, we choose Mask R-CNN [70] for
582 object detection and instance segmentation, and UperNet [71] for semantic segmentation. As for
583 backbones, we select ResNet-50, Swin-T, and ViT-B/16.

584 C More Experimental Results

585 **Performance Comparisons with More KD Methods.** In our main paper, we only compare our
586 ScaleKD with some related cross architecture KD approaches in Table 4, as few previous works
587 conduct experiments based on medium-sized students, such as ResNet-50. To make comparisons with
588 more CNN-based KD methods, we conduct experiments on a traditional student network, MobileNet-
589 V1, using the same training strategy as them. From the results shown in Table 14, we can observe
590 that by utilizing Swin-L as the teacher, our ScaleKD could help MobileNet-V1 reach 74.21% top-1
591 accuracy, outperforming previous methods which utilize ResNet-50 as the teacher by clear margins.

592 **Performance Comparisons with More Model Designs.** In our main paper, we apply our ScaleKD
593 to the basic design of each architecture and make comparisons with more recent variant architectures.
594 As illustrated in Table 15, we choose more designs to make more comprehensive comparisons.

595 D More Ablation Studies

596 **Ablation on Training Efficiency of ScaleKD.** As the TPP in our ScaleKD leverages the teacher’s last
597 stage, it will introduce additional training costs. To clearly study its training efficiency, we conduct

Table 14: Top-1 accuracy (%) comparison on IN-1K with more CNN-based KD methods. In the experiment, we adopt the same traditional training strategy as these methods.

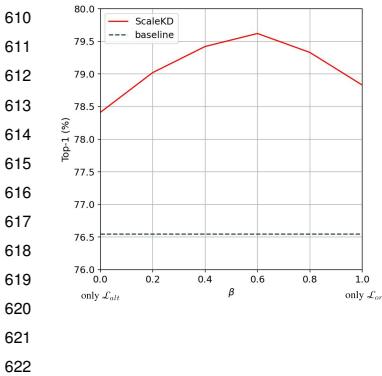
Model	Teacher	Method	Top-1(%)
MobileNet-V1	ResNet-50 (76.16)	From Scratch	69.63
		KD [14]	70.68
		AT [16]	70.72
		OFD [26]	71.25
		RKD [20]	71.23
		CRD [25]	71.40
		DKD [35]	72.05
		SRRL [30]	72.49
		ReviewKD [66]	72.56
		DIST [47]	73.24
		DistPro [33]	73.26
		MDG [67]	73.35
		DiffKD [48]	73.62
		Swin-L (86.24)	ScaleKD
			74.21

Table 15: Performance comparisons on IN-1K with various models. We conduct ScaleKD on the simplest design of each architecture and then make performance comparisons with various designs. Ours*, Ours†, Ours‡ denote choosing ViT-B (training from scratch), Swin-L (with IN-22K pre-training) and BEiT-L (with EVA pre-training) as the teacher, respectively.

Model	Training Dataset	Sample Resolution	Params (M)	FLOPs (G)	Top-1(%)
<i>CNN-based Architecture</i>					
ConvNext-T [3]	IN-1K	224 ²	28.59	4.46	82.14
ConvNext-T + Ours [†]	IN-1K	224 ²	28.59	4.46	84.16
ConvNext-T [3]	IN-22K \Rightarrow IN-1K	224 ²	28.59	4.46	82.90
ConvNext-B [3]	IN-1K	224 ²	87.77	15.14	83.80
UniRepLKNet-T [72]	IN-1K	224 ²	31.00	4.90	83.20
EfficientNet-B5 [73]	IN-1K	456 ²	30.00	9.90	83.60
RepViT-M2.3 [50]	IN-1K	224 ²	22.90	-	83.70
<i>MLP-based Architecture</i>					
Mixer-B/16 [6]	IN-1K	224 ²	59.88	12.61	76.44
Mixer-B/16 + Ours*	IN-1K	224 ²	59.88	12.61	81.62
Mixer-B/16 + Ours†	IN-1K	224 ²	59.88	12.61	81.96
Mixer-B/14 + Ours‡	IN-1K	224 ²	59.88	16.45	82.89
Mixer-B/16 [6]	IN-22K \Rightarrow IN-1K	224 ²	59.88	12.61	80.64
Mixer-L/16 [6]	IN-22K \Rightarrow IN-1K	224 ²	208.2	44.57	82.89
ResMLP-B24 [7]	IN-1K	224 ²	115.7	23.0	81.00
gMLP-B [8]	IN-1K	224 ²	73.00	15.80	81.60
<i>Transformer-based Architecture</i>					
ViT-S/16 [4]	IN-1K	224 ²	22.05	4.61	79.90
ViT-S/16 + Ours [†]	IN-1K	224 ²	22.05	4.61	83.93
ViT-S/16 [4, 74]	IN-22K \Rightarrow IN-1K	224 ²	22.05	4.61	80.50
Swin-T [5, 74]	IN-22K \Rightarrow IN-1K	224 ²	28.29	4.36	81.90
T2T-ViT _t -14 [75]	IN-1K	224 ²	21.47	4.34	81.83
DaViT-T [76]	IN-1K	224 ²	28.36	4.54	82.24
iLLaMA-S [77]	IN-1K	224 ²	21.90	-	79.90
EVA-02-S/14 [51]	IN-1K	336 ²	22.13	15.51	81.12
EVA-02-S/14 + Ours‡	IN-1K	336 ²	22.13	15.51	86.22
ViT-B/16 [4]	IN-1K	224 ²	86.57	17.58	81.80
Swin-B [5]	IN-1K	224 ²	87.77	15.14	83.50
T2T-ViT _t -24 [75]	IN-1K	224 ²	64.00	12.69	82.71
DaViT-B [76]	IN-1K	224 ²	87.95	15.51	84.09
ViT-B/16 [4]	IN-22K \Rightarrow IN-1K	224 ²	86.57	17.58	83.97
Swin-B [5]	IN-22K \Rightarrow IN-1K	224 ²	87.77	15.14	85.20
iLLaMA-B [77]	IN-22K \Rightarrow IN-1K	224 ²	86.30	-	85.00

598 ablative experiments in this section: i) As shown in Table 13a, we first compare the training efficiency
 599 of ScaleKD with traditional FD on three network pairs with increasing teacher’s model scales; ii) Then,
 600 as shown in Table 13b, we conduct the experiments on each design in ScaleKD. The experimental
 601 results show that: i) Using large teachers would induce more GPU memory occupation and longer
 602 training time. ii) Comparatively, TPP is the most resource-consuming component, especially after
 603 combining it with DFM. In summary, our ScaleKD needs additional training resources compared to
 604 traditional FD. However, if considering the significant performance gain it brings, these additional
 605 costs are acceptable in most cases.

606 Figure 4: Ablation Study on the
 607 hyper-parameter β .



Ablation Study on Hyper-parameter β . We conduct the ablation study on Swin-S → ResNet-50 network pair. According to Formulation 5 in the main paper, our method only has one hyper-parameter β , which is the balancing weight of the two views in DFM. To be specific, $\beta = 1.0$ indicates that only the first feature mimicking path exists, while $\beta = 0$ indicates that only the second feature mimicking path exists. And we select the β uniformly from 0 to 1 to determine the optimal β in ScaleKD. As shown in Figure 4, we can observe: i) In general, ScaleKD outperforms the baseline by significant margins at each setting, validating the efficacy of ScaleKD; ii) When $\beta = 0.6$, our ScaleKD achieves the best performance; iii) Though the second feature mimicking path could be used individually, it is inferior to the first path, indicating that the direct component is essential in feature mimicking; iv) When the two paths work collaboratively, they perform better than individually applied, which suggests that the two views are complementary with each other.

623 **Ablation Study on Pre-training and Distillation.** In
 624 this study, we explore the originality of the gains. We
 625 compare models trained by ScaleKD with upstream pre-
 626 trained models and upstream pre-training models with
 627 KD. From the results shown in Table 16, we notice: i)
 628 Compared to individual pre-training, applying KD under
 629 this stage can significantly boost model performance;
 630 ii) ViT-S/16 trained by ScaleKD significantly outper-
 631 forms the models trained with KD on IN-22K. The two
 632 observations illustrate that: i) Small students are diffi-
 633 cult to capture pre-training knowledge; ii) Our ScaleKD
 634 could effectively help the student to learn useful pre-
 635 training knowledge from the teacher without viewing
 636 the pre-training dataset.

637 E More Visualization Results

638 In this section, we provide more visualization results. In Figure 5, we provide the frequency
 639 distributions of three large pre-trained ViT models. We can observe the consistency in pre-trained
 640 ViTs’ unbalanced frequency distributions, that the direct responses are salient and obviously stronger
 641 than any other responses.

Table 16: Ablation study on pre-training and distillation.

Model	Method	Top-1(%)
ViT-S/16	Training from scratch on IN-1K	79.90
	Training from scratch on IN-1K w/ KD	81.42
	Pre-training on IN-22K	80.05
	Pre-training on IN-22K w/ KD	82.00
Training from scratch on IN-1K w/ ScaleKD		83.93

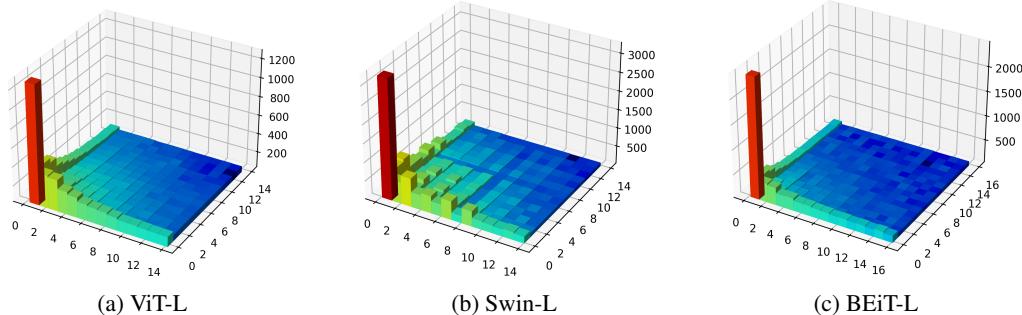


Figure 5: More illustrative feature distributions of large pre-trained ViTs in the frequency domain. We first collect the last stage output feature maps of 1600 samples from IN-1K, then conduct DCT on each channel, and finally take the average value across these samples, after turning all responses into absolute values.

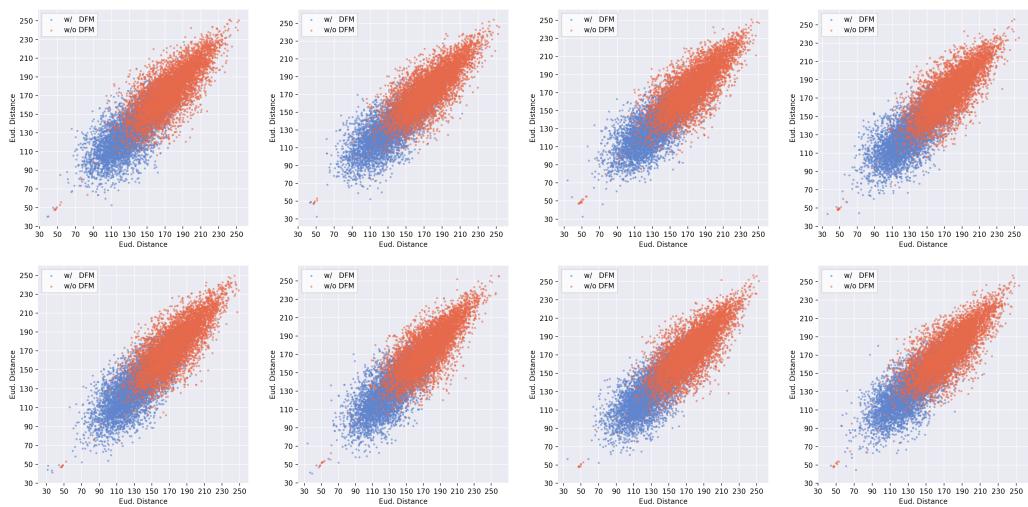


Figure 6: Feature distance distribution of alternative components for the last stage features between teacher and student on IN-1K. We obtain 64,000 feature pairs from 64,000 samples. After calculating the distance between teacher and student, we project the high-dimension distances into two-dimension space for illustration. Finally, we randomly select 6,400 data points for 8 times to draw the scatters. Blue points denote the distances without DFM, while orange points denote the distances with DFM.

642 **NeurIPS Paper Checklist**

643 **1. Claims**

644 Question: Do the main claims made in the abstract and introduction accurately reflect the
645 paper's contributions and scope?

646 Answer: [Yes]

647 Justification: The abstract and introduction contain three parts: i) motivation, ii) methodology,
648 and iii) experimental results, which are explained in Section 1, Section 2, and Section 3.
649 We discuss the related works in Section 4.

650 Guidelines:

- 651 • The answer NA means that the abstract and introduction do not include the claims
652 made in the paper.
- 653 • The abstract and/or introduction should clearly state the claims made, including the
654 contributions made in the paper and important assumptions and limitations. A No or
655 NA answer to this question will not be perceived well by the reviewers.
- 656 • The claims made should match theoretical and experimental results, and reflect how
657 much the results can be expected to generalize to other settings.
- 658 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
659 are not attained by the paper.

660 **2. Limitations**

661 Question: Does the paper discuss the limitations of the work performed by the authors?

662 Answer: [Yes]

663 Justification: We discuss the limitations in the Conclusion part in our main paper

664 Guidelines:

- 665 • The answer NA means that the paper has no limitation while the answer No means that
666 the paper has limitations, but those are not discussed in the paper.
- 667 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 668 • The paper should point out any strong assumptions and how robust the results are to
669 violations of these assumptions (e.g., independence assumptions, noiseless settings,
670 model well-specification, asymptotic approximations only holding locally). The authors
671 should reflect on how these assumptions might be violated in practice and what the
672 implications would be.
- 673 • The authors should reflect on the scope of the claims made, e.g., if the approach was
674 only tested on a few datasets or with a few runs. In general, empirical results often
675 depend on implicit assumptions, which should be articulated.
- 676 • The authors should reflect on the factors that influence the performance of the approach.
677 For example, a facial recognition algorithm may perform poorly when image resolution
678 is low or images are taken in low lighting. Or a speech-to-text system might not be
679 used reliably to provide closed captions for online lectures because it fails to handle
680 technical jargon.
- 681 • The authors should discuss the computational efficiency of the proposed algorithms
682 and how they scale with dataset size.
- 683 • If applicable, the authors should discuss possible limitations of their approach to
684 address problems of privacy and fairness.
- 685 • While the authors might fear that complete honesty about limitations might be used by
686 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
687 limitations that aren't acknowledged in the paper. The authors should use their best
688 judgment and recognize that individual actions in favor of transparency play an impor-
689 tant role in developing norms that preserve the integrity of the community. Reviewers
690 will be specifically instructed to not penalize honesty concerning limitations.

691 **3. Theory Assumptions and Proofs**

692 Question: For each theoretical result, does the paper provide the full set of assumptions and
693 a complete (and correct) proof?

694 Answer: [NA]

695 Justification: The paper does not include theoretical results

696 Guidelines:

- 697 • The answer NA means that the paper does not include theoretical results.
- 698 • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- 699 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 700 • The proofs can either appear in the main paper or the supplemental material, but if 701 they appear in the supplemental material, the authors are encouraged to provide a short 702 proof sketch to provide intuition.
- 703 • Inversely, any informal proof provided in the core of the paper should be complemented 704 by formal proofs provided in appendix or supplemental material.
- 705 • Theorems and Lemmas that the proof relies upon should be properly referenced.

706 **4. Experimental Result Reproducibility**

708 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
709 perimental results of the paper to the extent that it affects the main claims and/or conclusions
710 of the paper (regardless of whether the code and data are provided or not)?

711 Answer: [Yes]

712 Justification: We provide all training setups and implementation details in Appendix B

713 Guidelines:

- 714 • The answer NA means that the paper does not include experiments.
- 715 • If the paper includes experiments, a No answer to this question will not be perceived 716 well by the reviewers: Making the paper reproducible is important, regardless of 717 whether the code and data are provided or not.
- 718 • If the contribution is a dataset and/or model, the authors should describe the steps taken 719 to make their results reproducible or verifiable.
- 720 • Depending on the contribution, reproducibility can be accomplished in various ways.
721 For example, if the contribution is a novel architecture, describing the architecture fully
722 might suffice, or if the contribution is a specific model and empirical evaluation, it may
723 be necessary to either make it possible for others to replicate the model with the same
724 dataset, or provide access to the model. In general, releasing code and data is often
725 one good way to accomplish this, but reproducibility can also be provided via detailed
726 instructions for how to replicate the results, access to a hosted model (e.g., in the case
727 of a large language model), releasing of a model checkpoint, or other means that are
728 appropriate to the research performed.
- 729 • While NeurIPS does not require releasing code, the conference does require all submis-
730 sions to provide some reasonable avenue for reproducibility, which may depend on the
731 nature of the contribution. For example
 - 732 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
733 to reproduce that algorithm.
 - 734 (b) If the contribution is primarily a new model architecture, the paper should describe
735 the architecture clearly and fully.
 - 736 (c) If the contribution is a new model (e.g., a large language model), then there should
737 either be a way to access this model for reproducing the results or a way to reproduce
738 the model (e.g., with an open-source dataset or instructions for how to construct
739 the dataset).
 - 740 (d) We recognize that reproducibility may be tricky in some cases, in which case
741 authors are welcome to describe the particular way they provide for reproducibility.
742 In the case of closed-source models, it may be that access to the model is limited in
743 some way (e.g., to registered users), but it should be possible for other researchers
744 to have some path to reproducing or verifying the results.

745 **5. Open access to data and code**

746 Question: Does the paper provide open access to the data and code, with sufficient instruc-
747 tions to faithfully reproduce the main experimental results, as described in supplemental
748 material?

749 Answer: [No]

750 Justification: The code will be made publicly available. While it is currently under internal
751 review process, we decide not to release the code at this moment.

752 Guidelines:

- 753 • The answer NA means that paper does not include experiments requiring code.
- 754 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 755 • While we encourage the release of code and data, we understand that this might not be
756 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
757 including code, unless this is central to the contribution (e.g., for a new open-source
758 benchmark).
- 759 • The instructions should contain the exact command and environment needed to run to
760 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 761 • The authors should provide instructions on data access and preparation, including how
762 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 763 • The authors should provide scripts to reproduce all experimental results for the new
764 proposed method and baselines. If only a subset of experiments are reproducible, they
765 should state which ones are omitted from the script and why.
- 766 • At submission time, to preserve anonymity, the authors should release anonymized
767 versions (if applicable).
- 768 • Providing as much information as possible in supplemental material (appended to the
769 paper) is recommended, but including URLs to data and code is permitted.

772 6. Experimental Setting/Details

773 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
774 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
775 results?

776 Answer: [Yes]

777 Justification: We detailed explain all experimental settings in Appendix B

778 Guidelines:

- 779 • The answer NA means that the paper does not include experiments.
- 780 • The experimental setting should be presented in the core of the paper to a level of detail
781 that is necessary to appreciate the results and make sense of them.
- 782 • The full details can be provided either with the code, in appendix, or as supplemental
783 material.

784 7. Experiment Statistical Significance

785 Question: Does the paper report error bars suitably and correctly defined or other appropriate
786 information about the statistical significance of the experiments?

787 Answer: [No]

788 Justification: The main experiments are conducted on large-scale datasets, whose results are
789 stable.

790 Guidelines:

- 791 • The answer NA means that the paper does not include experiments.
- 792 • The authors should answer “Yes” if the results are accompanied by error bars, confi-
793 dence intervals, or statistical significance tests, at least for the experiments that support
794 the main claims of the paper.
- 795 • The factors of variability that the error bars are capturing should be clearly stated (for
796 example, train/test split, initialization, random drawing of some parameter, or overall
797 run with given experimental conditions).

- 798 • The method for calculating the error bars should be explained (closed form formula,
 799 call to a library function, bootstrap, etc.)
 800 • The assumptions made should be given (e.g., Normally distributed errors).
 801 • It should be clear whether the error bar is the standard deviation or the standard error
 802 of the mean.
 803 • It is OK to report 1-sigma error bars, but one should state it. The authors should
 804 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
 805 of Normality of errors is not verified.
 806 • For asymmetric distributions, the authors should be careful not to show in tables or
 807 figures symmetric error bars that would yield results that are out of range (e.g. negative
 808 error rates).
 809 • If error bars are reported in tables or plots, The authors should explain in the text how
 810 they were calculated and reference the corresponding figures or tables in the text.

811 **8. Experiments Compute Resources**

812 Question: For each experiment, does the paper provide sufficient information on the com-
 813 puter resources (type of compute workers, memory, time of execution) needed to reproduce
 814 the experiments?

815 Answer: [Yes]

816 Justification: We give the demand compute resources in Appendix B

817 Guidelines:

- 818 • The answer NA means that the paper does not include experiments.
 819 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
 820 or cloud provider, including relevant memory and storage.
 821 • The paper should provide the amount of compute required for each of the individual
 822 experimental runs as well as estimate the total compute.
 823 • The paper should disclose whether the full research project required more compute
 824 than the experiments reported in the paper (e.g., preliminary or failed experiments that
 825 didn't make it into the paper).

826 **9. Code Of Ethics**

827 Question: Does the research conducted in the paper conform, in every respect, with the
 828 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

829 Answer: [Yes]

830 Justification: We have reviewed the code of ethics.

831 Guidelines:

- 832 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
 833 • If the authors answer No, they should explain the special circumstances that require a
 834 deviation from the Code of Ethics.
 835 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
 836 eration due to laws or regulations in their jurisdiction).

837 **10. Broader Impacts**

838 Question: Does the paper discuss both potential positive societal impacts and negative
 839 societal impacts of the work performed?

840 Answer: [Yes]

841 Justification: We discuss the broader impacts in the Conclusion Part of the main paper.

842 Guidelines:

- 843 • The answer NA means that there is no societal impact of the work performed.
 844 • If the authors answer NA or No, they should explain why their work has no societal
 845 impact or why the paper does not address societal impact.
 846 • Examples of negative societal impacts include potential malicious or unintended uses
 847 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
 848 (e.g., deployment of technologies that could make decisions that unfairly impact specific
 849 groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pre-trained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We clearly illustrate the dataset and codebase we use in the Appendix A and the Appendix B, respectively.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

- 903 • If this information is not available online, the authors are encouraged to reach out to
904 the asset's creators.

905 **13. New Assets**

906 Question: Are new assets introduced in the paper well documented and is the documentation
907 provided alongside the assets?

908 Answer: [Yes]

909 Justification: We obtained some high-performance models in this paper.

910 Guidelines:

- 911 • The answer NA means that the paper does not release new assets.
912 • Researchers should communicate the details of the dataset/code/model as part of their
913 submissions via structured templates. This includes details about training, license,
914 limitations, etc.
915 • The paper should discuss whether and how consent was obtained from people whose
916 asset is used.
917 • At submission time, remember to anonymize your assets (if applicable). You can either
918 create an anonymized URL or include an anonymized zip file.

919 **14. Crowdsourcing and Research with Human Subjects**

920 Question: For crowdsourcing experiments and research with human subjects, does the paper
921 include the full text of instructions given to participants and screenshots, if applicable, as
922 well as details about compensation (if any)?

923 Answer: [NA]

924 Justification: The paper does not involve crowdsourcing nor research with human subjects.

925 Guidelines:

- 926 • The answer NA means that the paper does not involve crowdsourcing nor research with
927 human subjects.
928 • Including this information in the supplemental material is fine, but if the main contribu-
929 tion of the paper involves human subjects, then as much detail as possible should be
930 included in the main paper.
931 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
932 or other labor should be paid at least the minimum wage in the country of the data
933 collector.

934 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
935 Subjects**

936 Question: Does the paper describe potential risks incurred by study participants, whether
937 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
938 approvals (or an equivalent approval/review based on the requirements of your country or
939 institution) were obtained?

940 Answer: [NA]

941 Justification: The paper does not involve crowdsourcing nor research with human subjects.

942 Guidelines:

- 943 • The answer NA means that the paper does not involve crowdsourcing nor research with
944 human subjects.
945 • Depending on the country in which research is conducted, IRB approval (or equivalent)
946 may be required for any human subjects research. If you obtained IRB approval, you
947 should clearly state this in the paper.
948 • We recognize that the procedures for this may vary significantly between institutions
949 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
950 guidelines for their institution.
951 • For initial submissions, do not include any information that would break anonymity (if
952 applicable), such as the institution conducting the review.