

---

# Hybrid Distillation: Connecting Masked Autoencoders with Contrastive Learners

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Representation learning has been evolving from traditional supervised training to  
2 Contrastive Learning (CL) and Masked Image Modeling (MIM). Previous works  
3 have demonstrated their pros and cons in specific scenarios, *i.e.*, CL and supervised  
4 pre-training excel at capturing longer-range global patterns and enabling better  
5 feature discrimination, while MIM can introduce more local and diverse attention  
6 across all transformer layers. In this paper, we explore how to obtain a model  
7 that combines their strengths. We start by examining previous feature distillation  
8 and mask feature reconstruction methods and identify their limitations. We find  
9 that their increasing diversity mainly derives from the asymmetric designs, but  
10 these designs may in turn compromise the discrimination ability. In order to  
11 better obtain both discrimination and diversity, we propose a simple but effective  
12 Hybrid Distillation strategy, which utilizes both the supervised/CL teacher and the  
13 MIM teacher to jointly guide the student model. Hybrid Distill imitates the token  
14 relations of the MIM teacher to alleviate attention collapse, as well as distills the  
15 feature maps of the supervised/CL teacher to enable discrimination. Furthermore, a  
16 progressive redundant token masking strategy is also utilized to reduce the distilling  
17 costs and avoid falling into local optima. Experiment results prove that Hybrid  
18 Distill can achieve superior performance on different benchmarks.

## 19 1 Introduction

20 Pre-training followed by fine-tuning has been a common paradigm for computer vision tasks since  
21 the advent of deep learning. In the past decade, supervised image classification [16, 10, 24] over the  
22 widely used ImageNet [32] has dominated the pretraining mode. Recently, self-supervised learning  
23 has emerged as a promising alternative, particularly with two approaches: Contrastive Learning (CL)  
24 and Masked Image Modeling (MIM). The former one, typical representatives are MoCo [14] and  
25 SimCLR [4], learns invariant representation for positive views, which are usually defined as different  
26 augmentations of the same image. Furthermore, CLIP [30] extends CL to a multi-modal manner,  
27 which utilizes the corresponding text description of the given image as positive pairs. While the  
28 latter, including MAE [13] and SimMIM [44], aims to reconstruct the masked image patches and has  
29 become mainstream due to its efficiency brought by mask operations.

30 The different pre-training paradigms of CL and MIM facilitate a series of studies [43, 27, 38] that  
31 aim at understanding their respective properties. These studies point out that CL pre-training behaves  
32 more similar to supervised pre-training, *i.e.*, it provides models with longer-range global patterns  
33 targeting object shape, particularly in the last few layers [27], and enables feature representation with  
34 better **discrimination**. However, as shown in Fig. 1(a), CL pre-training causes self-attention in the  
35 last few layers to collapse into homogeneity, with attention distances located within a very small  
36 distance range. In contrast, MIM pre-training can bring more diverse attention and evenly distributed  
37 representations to all layers [43, 27], and this **diversity** contributes to its better generalization on

38 downstream fine-tuning. Nevertheless, MIM pre-training is slower to converge and underperforms in  
39 linear probing, mainly due to its lack of discrimination ability.

40 Since discrimination and diversity are both crucial for downstream adaptation, previous methods  
41 [41, 11, 23, 40, 29] propose to utilize feature distillation to combine the benefits of CL and MIM.  
42 Among them, dBOT [23] replaces the reconstructing objective of MAE with the feature maps of  
43 different pre-trained teachers. It finds that feature distillation can bring diverse attention no matter  
44 what the teacher model is, and the performance is comparable across different teachers, even with  
45 the randomly initialized ones, after multi-stage distillation. Also observing that distillation can yield  
46 diversity benefits, FD [41] directly distills feature maps from supervised/CL teachers to relieve the  
47 attention collapse and achieves considerable downstream performance gains. Although interesting  
48 and important, we argue that their findings are incomplete.

49 This paper re-examines these findings and reconsiders the importance of diversity and discrimination.  
50 Our study reveals the following observations: (i) **The increase in diversity derives from the  
51 asymmetric architecture designs, rather than feature distillation itself.** (Section 2.2) After  
52 removing the asymmetric attention in [41] and encoder-decoder designs in [23] and keeping the same  
53 teacher and student structures, we observe a negligible increase (or even a decrease) in attention  
54 diversity. (ii) **The asymmetric decoder de facto harm the discrimination over the encoder  
55 side, for it migrates the semantic information of the teacher model.** (Section 2.3) Due to the  
56 decomposition of the encoding and decoding functions, student encoders tend to summarize more  
57 general information, thus gradually losing the semantics obtained from teachers and yielding similar  
58 results after multi-stage distillation [23]. (iii) **Mask reconstruction of high-level semantics does  
59 not help improve diversity.** (Section 2.4) The phenomenon of reconstructing high-level information  
60 [29, 11, 40] is similar to direct feature distillation and lacks the diversity found in MIM, which  
61 implies that the attention diversity of MIM mainly comes from low-level reconstruction objectives.

62 Based on the above observations, we argue that a better distillation strategy is needed to help student  
63 models inherit both diversity and discrimination. To this end, we propose a simple but effective  
64 feature distillation method, termed as **Hybrid Distill**, to fully exploit the pre-trained model. Unlike  
65 previous works, Hybrid Distill aims to distill knowledge from both the supervised/CL and MIM  
66 teacher, allowing the student model to benefit from their respective advantages. To realize this, Hybrid  
67 Distill makes careful designs for the distilling target and location. Specifically, we find that **the  
68 relational modeling ability of MIM is crucial for preserving token diversity, while the feature  
69 maps of supervised/CL teachers are beneficial for discrimination.** Accordingly, we set the token  
70 relations of the MIM teacher and the feature maps of the supervised/CL teacher as the distilling  
71 objectives of Hybrid Distill. The token relations are distilled in layers preceding the final layer where  
72 attention collapse tends to occur, while the feature maps are distilled in the final layer to preserve  
73 semantics. Additionally, Hybrid Distill proposes a progressive redundant token masking strategy  
74 to reduce distilling costs and prevent falling into local optima. Experiment results show that the  
75 above distilling strategy works surprisingly well even when using MAE and CLIP teachers, *i.e.*, MAE  
76 pretrained with only 1.28M ImageNet images can also boost the large-scale (400M) pretrained CLIP  
77 teacher on different downstream tasks.

78 In a nutshell, this paper makes the following distribution:

- 79 • We re-examine the findings of previous feature distilling methods and point out that their increas-  
80 ing diversity mainly arises from the use of asymmetric designs, while these designs may in turn  
81 compromise the discrimination.
- 82 • We further propose a Hybrid Distill framework that utilized both supervised/CL and MIM teacher  
83 to provide the student with higher-quality discrimination and diversity. Distilling targets and locations  
84 are carefully designed in Hybrid Distill to fully exploit the strengths of both teachers.
- 85 • We conduct property analysis to demonstrate that the representations exhibit both discrimination  
86 and diversity in our Hybrid Distill. Experiments on various downstream tasks, including classification,  
87 detection, and segmentation, also showcase its superiority.

## 88 2 Model Evaluation: Diversity and Discrimination

89 This section re-examines the findings of previous feature distillation or mask feature reconstruction  
90 works illustrated in Sec. 1 and highlights their limitations in incorporating diversity and discrimination.

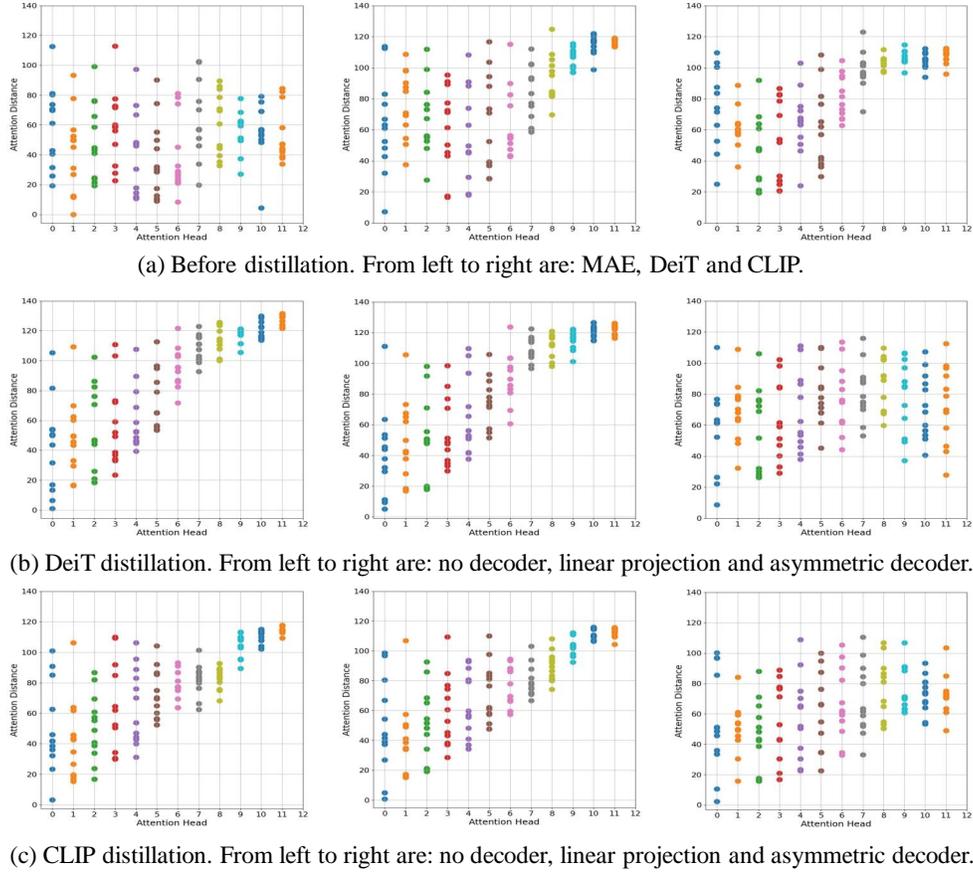


Figure 1: Average head distance after feature distillation with various decoders. (a) are the baselines. (b) use the supervised DeiT model as the teacher. (c) use the CL-based CLIP model as the teacher.

## 91 2.1 Preliminary

92 We first introduce the definitions of diversity and discrimination and the evaluation strategies we used.  
 93 **Discrimination** means that the representations contain more global patterns tailored to object shapes,  
 94 which is beneficial for recognizing objects and distinguishing images. **Diversity** is a relative concept,  
 95 which means that the model pays more attention to local information and can achieve more evenly  
 96 distributed representations, particularly in the last few layers.

97 We measure these properties by **average head distance** [41, 10] and **normalized mutual information**  
 98 **(NMI)** [33]. The former calculates the average distance between the query tokens and the key tokens  
 99 based on their attention weights, providing insight into whether the attention is global or local. The  
 100 latter measures whether the attention is attending to different tokens or similar ones and is calculated  
 101 following [27]. Specifically, let a uniform distribution  $p(q) = \frac{1}{N}$  represent the distribution of query  
 102 tokens, where  $N$  is the total token number. The joint distribution of query and key is then computed  
 103 as  $p(q, k) = \pi(k|q)p(q)$ , where  $\pi(k|q)$  is the normalized self-attention matrix. Thus, NMI can be  
 104 calculated by  $\frac{I(q,k)}{\sqrt{H(q)H(k)}}$  where  $I(\cdot, \cdot)$  is the mutual information and  $H(\cdot)$  is the marginal entropy.

## 105 2.2 The Increase in Diversity Derives from the Asymmetric Designs

106 Fig. 1 measures the average head distance after feature distillation with a consistent encoder structure  
 107 (vanilla Vision Transformer (ViT) [10]) for both the teacher and student models, along with various  
 108 decoders only for the student. It can be seen that when the encoder is kept the same, using no decoder  
 109 or linear projection decoder leads to a negligible increase (or even decrease) in attention diversity,  
 110 reflecting that feature distilling itself cannot bring benefits to diversity. Adding some extra attention  
 111 layers to the decoder can make the student encoder more diverse, but it hinders discrimination since  
 112 the last layer no longer captures long-range patterns. Fig. 2(a) further compares NMI using the DeiT  
 113 teacher and the results are in line with the attention visualization, *i.e.*, without asymmetric designs,

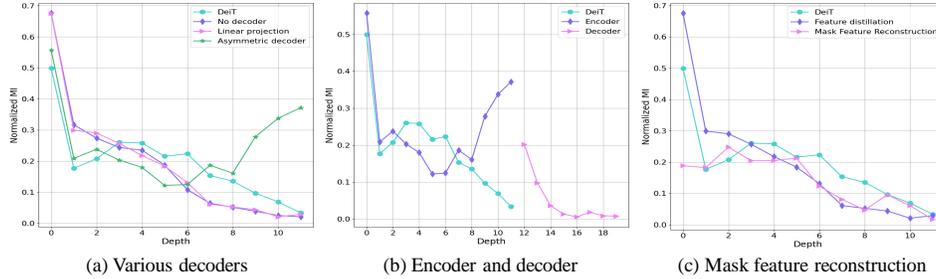


Figure 2: The normalized mutual information (NMI) of (a) various decoders, (b) encoder and decoder, and (c) mask feature reconstruction.

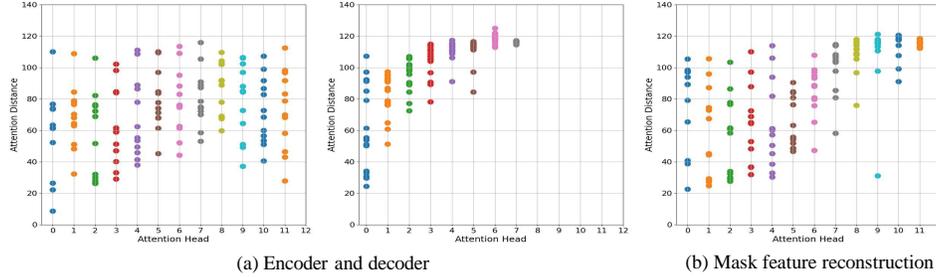


Figure 3: Average head distance of (a) encoder and decoder, and (b) mask feature reconstruction.

114 the student collapses into homogeneity and pays attention to similar tokens in the last few layers.  
 115 Conversely, the use of asymmetric decoders greatly reduces discrimination.

116 The above discussions focus on varying decoders, while FD [41] introduces asymmetric designs to  
 117 the encoder by adding additional learnable parameters and relative position bias to the attention layers  
 118 of the student. In the appendix, we demonstrate that the increase in diversity observed in FD also  
 119 arises from these designs and the diversity brought by them is not always significant.

### 120 2.3 The Asymmetric Decoder Harms the Encoder Discrimination

121 Fig. 3(a) and Fig. 2(b) further measure the average head distance and NMI of the asymmetric  
 122 decoder. Our findings suggest that the decoder has transferred the discrimination of the teacher, as its  
 123 behavior is similar to that of the last few layers of the teacher model where attention collapse occurs.  
 124 Reducing the number of decoder layers does not eliminate this transfer, as further demonstrated  
 125 in the appendix. Since only the student encoder is retained and applied to downstream tasks after  
 126 distillation, the semantic information that the model maintained is weakened, which explains why in  
 127 dBOT, different teachers tend to yield similarly-behaving models after multi-stage distilling. Note  
 128 that dBOT conducts feature distilling in a mask reconstruction way, while we demonstrate in both  
 129 Sec. 2.4 and the visualization in the appendix that it behaves similarly to directly distilling features.

### 130 2.4 Mask Reconstruction of High-Level Semantics Does not Help Improve Diversity

131 Fig. 3(b) and Fig. 2(c) examine the influence of mask reconstructing high-level information. To  
 132 eliminate the effect of the asymmetric decoder, we feed both the masks and tokens into the encoder  
 133 simultaneously and use only linear projection as the decoder. The overall process is thus similar  
 134 to SimMIM [44], except that we use the high-level information obtained from the supervised/CL  
 135 teacher as the distilling objective. Fig. 3(b) proves that reconstructing high-level information brings  
 136 no diversity gains towards directly distilling features, which is consistent with the finding of [45], *i.e.*,  
 137 reconstruction is unnecessary for MIM with semantic-rich teachers. This phenomenon also implies  
 138 that the diversity of MIM mainly arises from the low-level reconstructing objective rather than from  
 139 the reconstruction itself, since diversity is absent in high-level reconstruction.

## 140 3 Hybrid Distillation

141 From the above discussion, we conclude that existing distillation pipelines have limitations in  
 142 providing discrimination and diversity. Thus, we further propose a novel hybrid distillation framework  
 143 to ensure these important properties, and this section elaborates on its details.

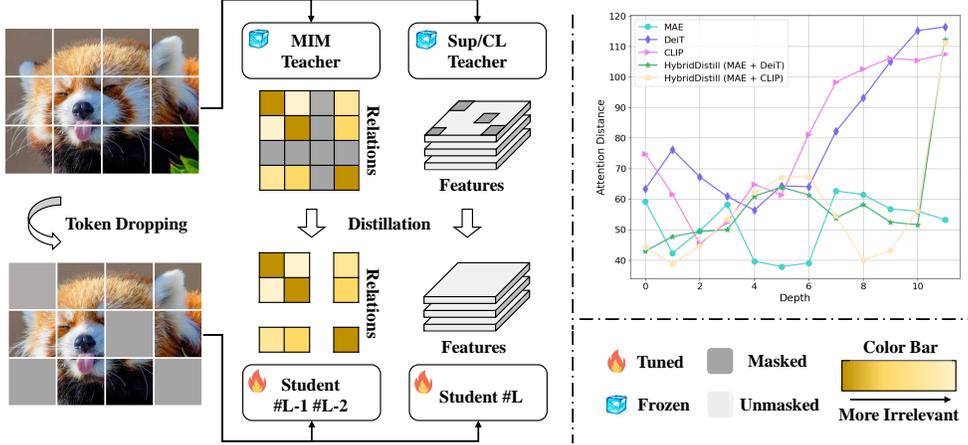


Figure 4: Hybrid Distill pipeline and its effectiveness in ensuring discrimination and diversity.

### 144 3.1 Overview

145 Given a supervised/CL pre-trained model  $T_c$ , and a MIM pre-trained model  $T_m$ , Hybrid Distill  
 146 simultaneously distills knowledge from these two different types of pre-trained teachers, aims at  
 147 combining their respective advantages to enhance the new representations in a randomly initialized  
 148 student model  $S_\theta$  where  $\theta$  is its learnable parameters. ViT [10] is adopted for all the models in Hybrid  
 149 Distill, and  $T_m$  is provided by MAE [13] while  $T_c$  is provided by DeiT [36] or CLIP [30].  
 150 Specifically, the Hybrid Distill framework is shown in Fig. 4 and its overall objective is:

$$\begin{aligned} \max_{\theta} \mathbb{E}_{x \sim \mathcal{X}} \mathcal{D} \{T_c(x) \odot M, S_\theta(M \odot x)\} \\ + \alpha \mathcal{D} \{T'_m(x) \odot M, S'_\theta(M \odot x)\}, \end{aligned} \quad (1)$$

151 where  $\odot$  is an element-wise product operation.  $M$  is a mask provided by the teacher model using  
 152 the strategy described in Sec. 3.2 and  $M \odot x$  denotes the unmasked patches.  $\mathcal{D}(\cdot, \cdot)$  is the distance  
 153 measurement, and we use smooth L1 distance in our experiment.  $\alpha$  is the hyperparameter that controls  
 154 the contribution of the two teacher models. Note that we do not distill the final output features  $T_m(x)$   
 155 for the MIM pre-trained model but instead use the token relations in the previous ViT layers, denote  
 156 as  $T'_m(x)$ , as the learning objective. Details are illustrated in Sec. 3.2.

### 157 3.2 Distilling Strategies

158 **What to distill?** Different from previous works [41, 11, 45] that directly distill the features of  
 159 teacher models, we analyze that the diversity of MIM pre-trained models arises from their superior  
 160 token-level relationship modeling, while supervised/CL pre-trained models excel at image-level  
 161 discrimination. Hence, we apply different distilling targets to  $T_c$  and  $T_m$  to fully utilize their respective  
 162 advantages. Specifically, taking  $T_m$  as an example, we decompose  $T_m$  into  $T_m^1 \circ T_m^2 \circ \dots \circ T_m^L$ ,  
 163 where  $T_m^i$  is the  $i^{th}$  layer of  $T_m$  and is composed of a multi-head self-attention (MSA) layer and an  
 164 MLP layer. Given  $x_m^i$  as the input of the  $i^{th}$  layer, the calculation in  $T_m^i$  can be represented as:

$$\begin{aligned} R_m^i(x_m^i) &= Q_m^i(x_m^i)K_m^i(x_m^i)^T, \\ \text{MSA}_m^i(x_m^i) &= \text{Softmax} \left( R_m^i(x_m^i) / \sqrt{d} \right) V_m^i(x_m^i), \\ T_m^i(x_m^i) &= x_m^i + \text{MLP}(x_m^i + \text{MSA}_m^i(x_m^i)), \end{aligned} \quad (2)$$

165 where  $Q_m^i$ ,  $K_m^i$ , and  $V_m^i$  denotes the linear mappings for  $x_m^i$  and  $d$  equals to the dimension of  $x_m^i$ .  
 166 Then, for MIM pre-trained model  $T_m$ , we set the token relation  $R_m^i(x_m^i)$  as the distilling target, while  
 167 for supervised/CL pretrained model  $T_c$ , we set the output features  $T_c^i(x_c^i)$  as the target.

168 **Where to distill?** As shown in Fig. 1(a), supervised and CL models tend to collapse into ho-  
 169 mogeneity in the last few layers, so Hybrid Distill chooses to distill token relations from  $T_m$  in  
 170 these layers to address this collapse and improve diversity. While for the last layer of  $S$  which is

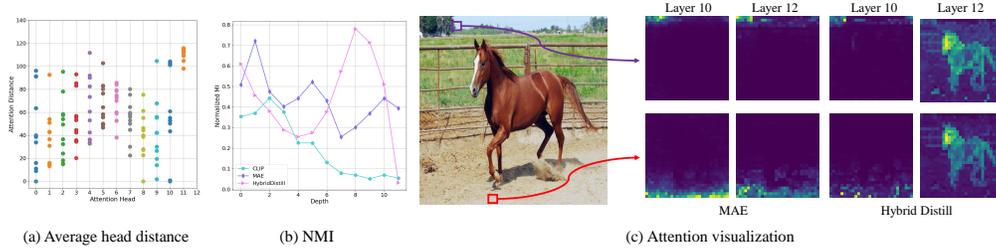


Figure 5: The (a) average head distance, (b) NMI, and (c) attention visualization of the student model obtained from Hybrid Distill with MAE and CLIP teachers.

171 crucial for discrimination, Hybrid Distill directly distills knowledge from  $T_c$  using the output features.  
 172 Specifically, we distill token relations from  $T_m$  at the  $L - 1$  and  $L - 2$  layers and distill features from  
 173  $T_c$  at the  $L$  layer of ViT. Accordingly, the learning objective  $T_c(x)$  and  $T'_m(x)$  in Eq. 1 become:

$$\begin{aligned} T_c(x) &= T_c^L(x), \\ T'_m(x) &= [R_m^{L-1}(x), R_m^{L-2}(x)]. \end{aligned} \quad (3)$$

174 **Distillation acceleration via redundant token dropping.** Suppose the input is divided into  $N$   
 175 tokens, *i.e.*,  $x \in \mathbb{R}^{N \times d}$ , Hybrid Distill can directly distill token relations and features using all the  $N$   
 176 tokens. However, since some tokens in the image may be redundant, it is promising to mask these  
 177 tokens for the student model  $S$  to reduce memory and time costs. Furthermore, removing redundant  
 178 tokens can play a regulatory role, helping the model avoid local optima during the distillation process.

179 Specifically, we use the MIM pre-trained teacher  $T_m$  to guide the identification of redundant tokens  
 180 and provide the token mask. Inspired by [20], we propose a progressive redundant token masking  
 181 strategy, which generates token masks at different layers of  $T_m$  in a progressive manner. Given  $x_m^i$   
 182 and the mask  $M_m^{i-1}$  provided by the previous layer, we define the tokens in  $x_m^i \odot M_m^{i-1}$  and are  
 183 top  $K\%$  similar to their average token as redundant tokens in the  $i^{th}$  layer and generate a redundant  
 184 token mask for them. The above process is denoted as  $T(x_m^i \odot M_m^{i-1}, K)$ . Next, we update  $M_m^i$   
 185 using  $T(x_m^i \odot M_m^{i-1}, K)$  and  $M_m^{i-1}$  as follows:

$$M_m^i = \begin{cases} M_m^{i-1} - T(x_m^i \odot M_m^{i-1}, K), & \text{if } i \in I, \\ M_m^{i-1} & \text{if } i \notin I. \end{cases} \quad (4)$$

186 where  $I$  is the set of layers required to update the token mask. For  $M_m^0$ , all elements are set to 1.  
 187 Finally, we set the mask  $M$  for the student model as  $M = M_m^L$ .

### 188 3.3 Property Analysis

189 **Average head distance.** Fig. 5(a) visualizes the average head distance of the student model with CLIP  
 190 and MAE as teachers, while the visualization of CLIP and MAE teachers themselves are included in  
 191 Fig. 1(a). These visualizations demonstrate that Hybrid Distill enhances the discrimination ability of  
 192 the student model, compensating for the semantic lacking problem of the MAE teacher. Moreover,  
 193 Hybrid Distill avoids succeeding attention collapse from the CLIP teacher and generates more diverse  
 194 representations in the last few layers.

195 **Normalized mutual information.** Fig. 5(b) further inspects the NMI. The results demonstrate  
 196 that the mutual information between tokens is significantly enhanced in the layers where the MAE  
 197 token relationships are distilled. Besides, this enhancement does not compromise the discrimination  
 198 obtained from CLIP, as evidenced by attention in the final layers still attending to similar tokens.

199 **Attention visualization.** Fig. 5(c) further visualizes the attention between a given query and other  
 200 keys at different layers to examine behaviors. Compared to MAE, Hybrid Distill exhibits better  
 201 discrimination ability, *i.e.*, the query tokens of the last layer have global attention towards the main  
 202 object of the images, regardless of their location. Besides, Hybrid Distill also improves the locality of  
 203 the model in the  $10^{th}$  layer, where attention collapse is known to occur in the CLIP teacher.

### 204 3.4 Discussion with Other Distillation Methods

205 Compared to previous distillation methods [41, 11, 23, 40, 29], Hybrid Distill stands out by not being  
 206 restricted to using a single teacher network. In addition to addressing the limitations of single-teacher

Table 1: Main results on ImageNet-1k classification, COCO detection and instance segmentation, and ADE20K semantic segmentation. \*: using MAE+DeiT teachers. †: using MAE+CLIP teachers.

Method	Backbone	Distill.	IN-1K	COCO		ADE20K
				AP <sup>box</sup>	AP <sup>Mask</sup>	
DeiT [36]	ViT-B		81.8	46.9	41.5	47.0
MoCo v3 [7]			83.2	45.5	40.5	47.1
DINO [2]			83.3	46.8	41.5	47.2
MAE [13]			83.6	48.4	42.6	48.1
CAE [5]			83.3	48.0	42.3	47.7
SdAE [8]			84.1	48.9	43.0	48.6
CLIP [30]			83.6	47.6	42.3	49.6
MAE [13]		ViT-L		85.9	54.0	47.1
CLIP [30]			86.1	52.7	46.2	54.2
Distill-DeiT	ViT-B		82.0	47.7	42.1	47.3
Distill-MAE		✓	83.7	49.1	43.1	47.8
Distill-CLIP			84.8	49.5	43.5	50.3
Hybrid Distill*	ViT-B	✓	83.7	50.3	44.2	49.1
Hybrid Distill†			<b>85.1</b>	<b>50.6</b>	<b>44.4</b>	<b>51.5</b>
Hybrid Distill†	ViT-L	✓	<b>88.0</b>	<b>54.6</b>	<b>47.6</b>	<b>56.3</b>

Table 2: Classification results on CIFAR100, Cars and INaturalist19. \*: using MAE+DeiT teachers. †: using MAE+CLIP teachers.

Method	Backbone	CIFAR100	Cars	INaturalist19	Mean
DeiT [36]	ViT-B	91.4	92.0	77.3	86.9
MAE [13]	ViT-B	89.6	89.5	75.2	84.8
Distill-DeiT	ViT-B	91.2	92.5	78.3	87.3
Distill-MAE	ViT-B	90.3	93.1	79.0	87.5
Distill-CLIP	ViT-B	91.6	94.3	81.6	89.2
Hybrid Distill*	ViT-B	91.7	94.1	80.2	88.7
Hybrid Distill†	ViT-B	<b>92.0</b>	<b>94.5</b>	<b>81.9</b>	<b>89.5</b>
Hybrid Distill†	ViT-L	<b>94.5</b>	<b>95.6</b>	<b>85.3</b>	<b>91.8</b>

207 distillation in enriching diversity (as discussed in Sec. 2), a more direct factor is that single-teacher  
 208 distillation cannot create new knowledge, *e.g.*, creating additional discrimination for the student  
 209 model when using the MIM teacher. Therefore, we believe that combining and utilizing existing  
 210 knowledge from various teachers is more effective and convenient. Furthermore, with the growing  
 211 availability of large-scale pre-trained models within the community, it becomes increasingly valuable  
 212 to explore new ways to utilize these models and combine their strengths. This further enhances the  
 213 practical value of our Hybrid Distill, and we hope our work would shed light on new directions.

## 214 4 Experiments

### 215 4.1 Implementation Details

216 Hybrid Distill is conducted on 8 V100 GPUs and is built on the codebase of dBOT [23], so most of its  
 217 settings are in line with dBOT. Specifically, the batch size, learning rate, and weight decay are set to  
 218 1024 and 6e-4, and 0.05, respectively. AdamW [26] optimizer and cosine decay [25] schedule is used.  
 219 The input size is 224<sup>2</sup>. For ViT-B, the distillation is based on ImageNet-1K and the epoch is 300  
 220 for main results and 100 for ablation studies. For ViT-L, the distillation is based on ImageNet-21K  
 221 and the epoch is 40. The hyperparameter  $\alpha$  is set to 1.0 and the redundant token masking set  $I$  is set  
 222 to  $[0, L/3, 2L/3]$  following [20]. The performances are tested on different downstream tasks. For  
 223 classification, we report results on ImageNet-1K, CIFAR100 [19], Cars [18], and iNaturalist19 [37].  
 224 For object detection and instance segmentation, we fine-tune the student model on COCO [22] using  
 225 Mask-RCNN [15] following [5]. For semantic segmentation, the evaluation is conducted on ADE20K  
 226 [47] using the ViT with UperNet [42] following [5, 8]. More details are included in the appendix.

### 227 4.2 Main Results

228 This section presents benchmark results of Hybrid Distill on different downstream. We also list results  
 229 for supervised and self-supervised pre-trained models, as well as 300-epoch uni-distillation baselines

Table 3: Different combinations of two teacher models.  $T_c(x)$ : DeiT,  $T_m(x)$ : MAE.

Targets	AP <sup>box</sup>	AP <sup>mask</sup>
$T_c(x)$	47.5	41.8
$T_m(x)$	48.9	43.1
$T_c(x) + T'_c(x)$	46.8	41.5
$T_m(x) + T'_m(x)$	48.9	43.2
$T_c(x) + T'_m(x)$	<b>50.0</b>	<b>43.9</b>

Table 5: The distilling targets of  $T'_m(x)$ .  $T_c(x)$ : DeiT,  $T_m(x)$ : MAE.  $\star$  means distilling MAE and DeiT features at the last layer.

Targets	AP <sup>box</sup>	AP <sup>mask</sup>
$T_m^i \star$	47.7	42.1
$T_m^i$	49.6	43.5
MSA <sup>i</sup> <sub>m</sub>	49.8	43.7
R <sup>i</sup> <sub>m</sub>	<b>50.0</b>	<b>43.9</b>

Table 4: Different combinations of two teacher models.  $T_c(x)$ : CLIP,  $T_m(x)$ : MAE.  $\star$ : using the ImageNet-100 pretrained weights.

Targets	AP <sup>box</sup>	AP <sup>mask</sup>
$T_c(x)$	49.1	43.1
$T_m(x)$	48.9	43.1
$T_c(x) + T'_c(x)$	49.1	43.2
$T_c(x) + T'_m(x)$	<b>50.4</b>	<b>44.1</b>
$T_c(x) + T'_m(x) \star$	49.5	43.5

Table 6: The distilling targets of  $T'_m(x)$ .  $T_c(x)$ : CLIP,  $T_m(x)$ : MAE.

Targets	AP <sup>box</sup>	AP <sup>mask</sup>
$T_m^i$	49.9	44.0
MSA <sup>i</sup> <sub>m</sub>	50.1	44.0
R <sup>i</sup> <sub>m</sub>	<b>50.4</b>	<b>44.1</b>

230 which use the same symmetrical structures as Hybrid Distill, for comparison. As shown in Tab. 1,  
 231 Hybrid Distill achieves performance gains on all downstream tasks, especially for the dense-level  
 232 ones. Specifically, although the performance of DeiT is suboptimal, its strength can be complementary  
 233 to MAE and brings considerable benefits, *i.e.*, when using DeiT and MAE teachers, Hybrid Distill  
 234 achieves 50.3 AP<sup>box</sup> and 44.2 AP<sup>mask</sup> on COCO, as well as 49.1 mIoU on ADE20K, surpassing  
 235 Distill-MAE by 1.2, 1.1, and 1.3, respectively. Similarly, Hybrid Distill achieves 50.6 AP<sup>box</sup> and  
 236 44.4 AP<sup>mask</sup> on COCO, as well as 51.5 mIoU on ADE20K when using CLIP and MAE teachers,  
 237 outperforming Distill-CLIP by 1.1, 0.9, and 1.2, respectively. When using the ViT-L backbone, the  
 238 performance can be further boosted to 54.6 AP<sup>box</sup>, 47.6 AP<sup>mask</sup> and 56.3 mIoU on respective tasks.  
 239 The improvement on ImageNet-1k is not significant, probably because the distillation is performed on  
 240 the same dataset, thus increasing diversity fails to bring further gains. In Tab. 2, we further evaluate  
 241 Hybrid Distill on several small-scale classification datasets and observe more significant gains.

### 242 4.3 Ablation Study

243 This section ablates different variants of Hybrid Distill. The results are reported on dense-level COCO  
 244 detection and segmentation tasks, as diversity has a stronger influence on these dense-level tasks [27].

245 **Different combinations of two teachers.** We first evaluate the benefits of combining two teachers  
 246 for distillation. As shown in Tab. 3, adding additional MAE attention regularization can bring  
 247 noticeable improvements (2.5 on AP<sup>box</sup> and 2.1 on AP<sup>mask</sup>) compared to directly distilling from the  
 248 DeiT teacher. Moreover, the additional attention regularization cannot bring benefits when only using  
 249 a single DeiT teacher, which suggests that the benefits come from the introduction of MAE teacher.  
 250 The above conclusions are consistent when using CLIP and MAE teachers, as illustrated in Tab. 4.  
 251 We also try a much weaker version of MAE teacher which is only pre-trained on ImageNet-100 for  
 252 100 epochs in Tab. 4. We lower the weight of this teacher to avoid its impact on discrimination. The  
 253 results are still positive, which reflects the power of the MIM pre-training in modeling diversity.

254 **Distilling target of the MIM teacher.** We then examine the distilling target of the MIM teacher.  
 255 As shown in Tab. 5, distilling the relation R<sup>i</sup><sub>m</sub> brings the best detection performance (50.0AP<sup>box</sup>).  
 256 Distilling MSA<sup>i</sup><sub>m</sub> achieves a close performance (49.8AP<sup>box</sup>) since its essential is also distilling  
 257 relationships, while directly distilling the feature maps T<sup>i</sup><sub>m</sub> brings the worst performance (49.6AP<sup>box</sup>).  
 258 Nevertheless, all these schemes outperform the DeiT distillation baseline, and the trends are consistent  
 259 when using CLIP and MAE teachers, as shown in Tab. 6. Besides, we also evaluate a basic setting  
 260 that directly distills the features of both the MAE and DeiT teachers at the last layer. The result is far  
 261 from satisfactory, which highlights the effectiveness of the designs in Hybrid Distill.

262 **Distilling position of the MIM teacher.** Tab. 7 inspect the distilling position of the MIM teacher.  
 263 We first experiment with distilling MAE relations at the front, middle, and back layers. Distilling at  
 264 the back layers achieves better results, *i.e.*, 1.5AP<sup>box</sup> and 2.4AP<sup>box</sup> gains towards distilling at the

Table 7: The distilling position of  $T_m$ .

Distilling layers	AP <sup>box</sup>	AP <sup>mask</sup>
1-11	48.8	43.0
1,2,3	47.4	41.9
5,6,7	48.3	42.7
9,10,11	49.8	43.7
10,11	<b>50.0</b>	<b>43.9</b>
11	49.2	43.3

Table 8: The token masking strategy.

Strategy	Ratio	AP <sup>box</sup>	AP <sup>mask</sup>
No	100%	<b>50.0</b>	<b>43.9</b>
Random	35%	49.2	43.3
Direct	35%	49.6	43.7
Progressive	13%(50% <sup>3</sup> )	48.4	42.8
Progressive	34%(70% <sup>3</sup> )	49.9	43.8
Progressive	73%(90% <sup>3</sup> )	49.9	43.8

265 front and middle, respectively. The results are consistent with the fact that attention collapse tends to  
 266 occur in these back layers. We then ablate the number of distilling layers and find that distilling at the  
 267 two layers preceding the final layer (*i.e.*, 10,11) contributes to the best results.

268 **Token masking strategy.** Tab. 8 studies different masking strategies for the student model. Since  
 269 we progressive drop the redundant tokens three times, the actual tokens used in the student model are  
 270  $(1 - K)^3\%$ . We observe that when dropping 30% tokens at a time, Hybrid Distill achieves very close  
 271 performance (49.9AP<sup>box</sup> and 43.8AP<sup>mask</sup>) to the no masking results and outperforms the random  
 272 masking strategy and the direct masking strategy which only generates token mask at the last layer. In  
 273 addition, we notice that our token masking strategy also has a regularizing effect, which can prevent  
 274 the model from falling into a locally optimal when training for longer epochs. Details about this  
 275 effect are included in the appendix.

## 276 5 Related Work

277 **Representation learning.** Pre-training on large-scale datasets (e.g., ImageNet [32], JFT [34], Kinetics  
 278 [3], etc.) is typically utilized for downstream initialization. Except for the common supervised pre-  
 279 training [16, 10, 24], contrastive learning (CL) [4, 14, 6, 12] and masked image modeling (MIM)  
 280 [1, 44, 13] dominate the recent research. The former is achieved by pulling close the features of two  
 281 different augment views of the input image. While the latter, inspired by masked language modeling  
 282 [17, 46] in NLP, is realized by reconstructing the mask part of the input image. Recently multi-model  
 283 extensions [30, 9, 21] of the CL pre-training have also been proposed by utilizing the paired text  
 284 description of the given image. These different types of pre-training frameworks are proven to have  
 285 different properties [27, 43], and this paper aims to combine their respective excellent properties to  
 286 boost a student model.

287 **Knowledge distillation.** Knowledge distillation [28, 35, 31] utilizes a well-trained teacher to guide  
 288 the feature learning of the student model, thus transferring its ability to the student. Beyond its  
 289 success in supervised learning, some recent works [41, 11, 39, 40, 29] utilize it to extend existing  
 290 pretrained models or paradigms. Feature distillation (FD) [41] finds that distilling the feature map  
 291 of the supervised/CL pretrained teacher can bring diverse representation to the student and make it  
 292 more friendly for downstream fine-tuning. dBOT [23], MVP [40], and BEiT v2 [29] change the mask  
 293 reconstruction object of MIM to the knowledge of the teacher model to boost MIM pre-training with  
 294 semantic information. In this paper, we analyze their properties and propose a new hybrid distillation  
 295 framework to deal with their deficiencies.

## 296 6 Conclusion

297 This paper proposed a hybrid distillation framework that simultaneously distills knowledge from  
 298 both the supervised/CL pre-trained teacher and MIM pre-trained teacher to enhance the diversity and  
 299 discrimination of the student. The framework addresses the limitations of single-teacher distillation,  
 300 where increasing diversity through the use of asymmetric designs may harm discrimination. Specifi-  
 301 cally, Hybrid Distill carefully designs the distilling target and location, *i.e.*, distilling relations from  
 302 MIM in layers where attention collapse tends to occur and distilling features from supervised/CL  
 303 in the last layer to preserve discrimination. A progressive redundant token masking strategy is also  
 304 proposed for reducing the distilling costs. Experiments prove that Hybrid Distill can acquire better  
 305 properties and achieve promising results on various downstream. We hope our research would shed  
 306 light on a new direction for applying existing large-scale pre-trained models.

307 **References**

- 308 [1] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. BEit: BERT pre-training of image  
309 transformers. In *International Conference on Learning Representations*, 2022.
- 310 [2] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski,  
311 and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings*  
312 *of the International Conference on Computer Vision (ICCV)*, 2021.
- 313 [3] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the  
314 kinetics dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern*  
315 *Recognition (CVPR)*, 2017.
- 316 [4] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework  
317 for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*, 2020.
- 318 [5] Xiaokang Chen, Mingyu Ding, Xiaodi Wang, Ying Xin, Shentong Mo, Yunhao Wang, Shumin  
319 Han, Ping Luo, Gang Zeng, and Jingdong Wang. Context autoencoder for self-supervised  
320 representation learning. *arXiv preprint arXiv:2202.03026*, 2022.
- 321 [6] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum  
322 contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020.
- 323 [7] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised  
324 vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer*  
325 *Vision*, pages 9640–9649, 2021.
- 326 [8] Yabo Chen, Yuchen Liu, Dongsheng Jiang, Xiaopeng Zhang, Wenrui Dai, Hongkai Xiong, and  
327 Qi Tian. Sdae: Self-distillated masked autoencoder. In *ECCV*, 2022.
- 328 [9] Yufeng Cui, Lichen Zhao, Feng Liang, Yangguang Li, and Jing Shao. Democratizing contrastive  
329 language-image pre-training: A clip benchmark of data, model, and supervision, 2022.
- 330 [10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,  
331 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al.  
332 An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint*  
333 *arXiv:2010.11929*, 2020.
- 334 [11] Yuxin Fang, Wen Wang, Binhui Xie, Quan Sun, Ledell Wu, Xinggang Wang, Tiejun Huang,  
335 Xinlong Wang, and Yue Cao. Eva: Exploring the limits of masked visual representation learning  
336 at scale. *arXiv preprint arXiv:2211.07636*, 2022.
- 337 [12] Jean-Bastien Grill, Florian Strub, Florent Alché, Corentin Tallec, Pierre H Richemond, Elena  
338 Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi  
339 Azar, et al. Bootstrap your own latent: A new approach to self-supervised learning. *arXiv*  
340 *preprint arXiv:2006.07733*, 2020.
- 341 [13] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked  
342 autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on*  
343 *Computer Vision and Pattern Recognition*, pages 16000–16009, 2022.
- 344 [14] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for  
345 unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on*  
346 *Computer Vision and Pattern Recognition (CVPR)*, pages 9729–9738, 2020.
- 347 [15] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *ICCV*, pages  
348 2961–2969, 2017.
- 349 [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
350 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*  
351 *(CVPR)*, pages 770–778, 2016.
- 352 [17] Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. Bert: Pre-training of deep  
353 bidirectional transformers for language understanding. In *NAACL-HLT*, pages 4171–4186,  
354 2019.

- 355 [18] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-  
356 grained categorization. In *Proceedings of the IEEE international conference on computer vision*  
357 *workshops*, pages 554–561, 2013.
- 358 [19] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.  
359 2009.
- 360 [20] Jin Li, Yaoming Wang, XIAOPENG ZHANG, Yabo Chen, Dongsheng Jiang, Wenrui Dai,  
361 Chenglin Li, Hongkai Xiong, and Qi Tian. Progressively compressed auto-encoder for self-  
362 supervised representation learning. In *The Eleventh International Conference on Learning*  
363 *Representations*, 2023.
- 364 [21] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image  
365 pre-training for unified vision-language understanding and generation. In *ICML*, 2022.
- 366 [22] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr  
367 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Proceedings*  
368 *of the European Conference on Computer Vision (ECCV)*, pages 1209–1218, 2014.
- 369 [23] Xingbin Liu, Jinghao Zhou, Tao Kong, Xianming Lin, and Rongrong Ji. Exploring target  
370 representations for masked autoencoders. *arXiv preprint arXiv:2209.03917*, 2022.
- 371 [24] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining  
372 Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings*  
373 *of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.
- 374 [25] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv*  
375 *preprint arXiv:1608.03983*, 2016.
- 376 [26] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*  
377 *arXiv:1711.05101*, 2017.
- 378 [27] Namuk Park, Wonjae Kim, Byeongho Heo, Taekyung Kim, and Sangdoon Yun. What do  
379 self-supervised vision transformers learn? *arXiv preprint arXiv:2305.00729*, 2023.
- 380 [28] Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. Relational knowledge distillation. In  
381 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages  
382 3967–3976, 2019.
- 383 [29] Zhiliang Peng, Li Dong, Hangbo Bao, Qixiang Ye, and Furu Wei. BEiT v2: Masked image  
384 modeling with vector-quantized visual tokenizers. 2022.
- 385 [30] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
386 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
387 models from natural language supervision. In *International Conference on Machine Learning*  
388 *(ICML)*, 2021.
- 389 [31] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and  
390 Yoshua Bengio. Fitnets: Hints for thin deep nets. *arXiv preprint arXiv:1412.6550*, 2014.
- 391 [32] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng  
392 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual  
393 recognition challenge. *International journal of computer vision (IJCV)*, 115(3):211–252, 2015.
- 394 [33] Alexander Strehl and Joydeep Ghosh. Cluster ensembles—a knowledge reuse framework for  
395 combining multiple partitions. *Journal of machine learning research*, 3(Dec):583–617, 2002.
- 396 [34] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable  
397 effectiveness of data in deep learning era. In *Proceedings of the IEEE International Conference*  
398 *on Computer Vision (ICCV)*, 2017.
- 399 [35] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive representation distillation. *arXiv*  
400 *preprint arXiv:1910.10699*, 2019.

- 401 [36] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and  
402 Herve Jegou. Training data-efficient image transformers & distillation through attention. In  
403 *International Conference on Machine Learning (ICML)*, volume 139, pages 10347–10357, July  
404 2021.
- 405 [37] Grant Van Horn, Oisin Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig  
406 Adam, Pietro Perona, and Serge Belongie. The inaturalist species classification and detection  
407 dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
408 pages 8769–8778, 2018.
- 409 [38] Shaoru Wang, Jin Gao, Zeming Li, Xiaoqin Zhang, and Weiming Hu. A closer look at self-  
410 supervised lightweight vision transformers, 2023.
- 411 [39] Wenguan Wang, Tianfei Zhou, Fisher Yu, Jifeng Dai, Ender Konukoglu, and Luc Van Gool. Ex-  
412 ploring cross-image pixel contrast for semantic segmentation. *arXiv preprint arXiv:2101.11939*,  
413 2021.
- 414 [40] Longhui Wei, Lingxi Xie, Wengang Zhou, Houqiang Li, and Qi Tian. Mvp: Multimodality-  
415 guided visual pre-training. In *Computer Vision–ECCV 2022: 17th European Conference, Tel  
416 Aviv, Israel, October 23–27, 2022, Proceedings, Part XXX*, pages 337–353. Springer, 2022.
- 417 [41] Yixuan Wei, Han Hu, Zhenda Xie, Zheng Zhang, Yue Cao, Jianmin Bao, Dong Chen, and  
418 Baining Guo. Contrastive learning rivals masked image modeling in fine-tuning via feature  
419 distillation. *Tech Report*, 2022.
- 420 [42] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing  
421 for scene understanding. In *Proceedings of the European Conference on Computer Vision  
422 (ECCV)*, pages 418–434, 2018.
- 423 [43] Zhenda Xie, Zigang Geng, Jingcheng Hu, Zheng Zhang, Han Hu, and Yue Cao. Revealing the  
424 dark secrets of masked image modeling. *arXiv preprint arXiv:2205.13543*, 2022.
- 425 [44] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and  
426 Han Hu. Simmim: A simple framework for masked image modeling. In *Proceedings of the  
427 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9653–9663, 2022.
- 428 [45] Hongwei Xue, Peng Gao, Hongyang Li, Yu Qiao, Hao Sun, Houqiang Li, and Jiebo Luo. Stare at  
429 what you see: Masked image modeling without reconstruction. *arXiv preprint arXiv:2211.08887*,  
430 2022.
- 431 [46] Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. Ernie:  
432 Enhanced language representation with informative entities. In *ACL*, pages 1441–1451, 2019.
- 433 [47] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio  
434 Torralba. Semantic understanding of scenes through the ade20k dataset. *International Journal  
435 on Computer Vision (IJCV)*, 127:302–321, 2019.