

---

# Appendix for Hybrid Distillation: Connecting Masked Autoencoders with Contrastive Learners

---

Anonymous Author(s)

Affiliation

Address

email

Table A1: Compared with more baselines using ViT-B as the backbone.  $\star$ : using MAE+DeiT teachers.  $\dagger$ : using MAE+CLIP teachers.

Method	COCO		ADE20K
	AP <sup>box</sup>	AP <sup>Mask</sup>	
Distill-DeiT	47.7	42.1	47.3
Distill-MAE	49.1	43.1	47.8
Distill-CLIP	49.5	43.5	50.3
FD-DeiT [11]	47.0	41.6	47.9
FD-MAE [11]	48.1	42.6	47.0
FD-CLIP [11]	49.2	43.3	50.5
dBOT-DeiT [9]	47.5	41.9	47.9
dBOT-MAE [9]	49.3	43.5	48.2
Hybrid Distill $\star$	50.3	44.2	49.1
Hybrid Distill $\dagger$	<b>50.6</b>	<b>44.4</b>	<b>51.5</b>

## 1 A More Experimental Results

### 2 A.1 Compared with More Baselines

3 Tab. A1 compares Hybrid Distill with two other methods, *i.e.*, dBOT [9] and FD [11], which employ  
4 asymmetric designs in distillation. We conduct distilling for 300 epochs based on their corresponding  
5 official codes<sup>1</sup>. We omit the dBOT-CLIP result since dBOT specifically removes the asymmetric  
6 designs for CLIP, thus its distillation process is similar to our Distill-CLIP baseline. As shown in  
7 Tab. A1, their benefits towards symmetrical distillation are not always significant, and the performance  
8 is inferior to our Hybrid Distill, which validates the effectiveness of our framework.

### 9 A.2 Results with Cascade Mask-RCNN

10 Tab. A2 further presents the object detection and instance segmentation results of Hybrid Distill with  
11 Cascade Mask-RCNN, which allows for a direct comparison with dBOT [9], as they also provide  
12 1600-epoch distillation results under this setting. As shown, 300-epoch Hybrid Distill with MAE and  
13 DeiT teachers can achieve 53.0 AP<sup>box</sup>, outperforming 1600-epoch dBOT-DeiT (52.5 AP<sup>box</sup>) and  
14 dBOT-MAE (52.7 AP<sup>box</sup>). Additionally, 300-epoch Hybrid Distill with MAE and CLIP teachers  
15 achieves 53.4 AP<sup>box</sup>, which is also very close to the 1600-epoch dBOT-CLIP result (53.6 AP<sup>box</sup>).  
16 The above results reflect that due to the better properties obtained, Hybrid Distill can obtain promising  
17 results with fewer training epochs.

<sup>1</sup>dBOT [9]: <https://github.com/liuxingbin/dbot/>. FD [11]: <https://github.com/SwinTransformer/Feature-Distillation/>. Since FD does not provide codes for downstream verification, we uniformly perform verification under our downstream frameworks.

Table A2: Object detection and instance segmentation results with *Cascade Mask-RCNN*. \*: using MAE+DeiT teachers. †: using MAE+CLIP teachers.

Method	Epoch	AP <sup>box</sup>	AP <sup>mask</sup>
Distill-DeiT	300	50.4	43.4
Distill-MAE	300	51.9	44.7
Distill-CLIP	300	52.4	45.0
dBOT-DeiT [9]	2 × 800	52.5	-
dBOT-MAE [9]	2 × 800	52.7	-
dBOT-CLIP [9]	1 × 1600	53.6	-
Hybrid Distill*	300	53.0	45.6
Hybrid Distill†	300	<b>53.4</b>	<b>45.9</b>

Table A3: Hybrid Distill uses MAE and DINO as teachers. Object Detection and instance segmentation results are reported with *Mask-RCNN*, following the setting in Tab. 1 of our main paper.

Method	AP <sup>box</sup>	AP <sup>mask</sup>
MAE	48.4	42.6
DINO	46.8	41.5
Distill-DINO	47.5	41.9
Distill-MAE	49.1	43.1
Hybrid Distill	<b>49.6</b>	<b>43.5</b>

### 18 A.3 Hybrid Distillation with DINO

19 Tab. A3 test the results of our Hybrid Distill using the MAE and DINO teachers. Under this setting,  
 20 Hybrid Distill achieves 49.6 AP<sup>box</sup> and 43.5 AP<sup>mask</sup>. Although still superior to the baselines,  
 21 results with DINO are not as good as those with CLIP and DeiT. We analyze that this is because the  
 22 discrimination of DINO is weaker than DeiT and CLIP, which makes its complementarity with MAE  
 23 also weaker than the latter two. The visualization in Fig.A1 provides evidence for this. On the one  
 24 hand, we notice that the average attention distance of DINO itself is lower than that of DeiT and  
 25 CLIP in the final layer. On the other, the attention maintenance of the final layer after distillation is  
 26 weaker compared with that obtained by DeiT and CLIP.

### 27 A.4 More Ablation Studies

28 **The choice of hyperparameter  $\alpha$ .** Tab. A4 ablates different setting of  $\alpha$ . It can be concluded  
 29 that adding additional MIM supervision can lead to performance improvement towards not using  
 30 MIM supervision ( $\alpha = 0$ ), regardless of the value of  $\alpha$ . While setting  $\alpha$  to 1.0 can bring the best  
 31 performance for both MAE+DeiT and MAE+CLIP teachers. Using the CLIP teacher achieves more  
 32 stable performance since CLIP itself has higher quality compared with DeiT, while DeiT relies more  
 33 on the help of MAE.

34 **Token masking strategy and local optima.** Tab. A5 further reveals that the proposed progressive  
 35 redundant token masking strategy in Hybrid Distill can prevent the student from falling into local  
 36 optima. As shown, when the token mask is removed and the distillation epoch is prolonged from  
 37 100 to 300, no performance gains are observed. This phenomenon has also been observed in [4]. We  
 38 analyze that over-fitting is the root cause of this problem and introducing token masks can alleviate it  
 39 since they can play a regulatory role. The performance gains achieved by the token masks provide  
 40 clear support for their effectiveness.

## 41 B Further Discussion about Diversity and Discrimination

### 42 B.1 Asymmetric Encoder Designs

43 Fig. A2 studies the asymmetric encoder designs used in FD, *i.e.*, adding additional learnable param-  
 44 eters and relative position bias to the attention layers of the student. As shown, the asymmetric encoder  
 45 (Fig. A2(c)) *de facto* improves diversity compared to using only the symmetric encoder (Fig. A2(b)).

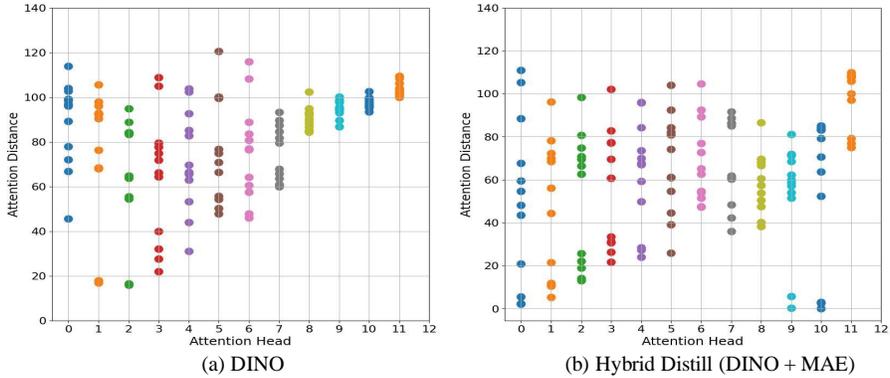


Figure A1: Average head distance of different (a) DINO baseline and (b) Hybrid Distill with MAE and DINO as teachers.

Table A4: Ablation on the hyperparameter  $\alpha$  which controls the contribution of two teacher models. (a)  $T_c(x)$ : DeiT,  $T_m(x)$ : MAE.

$\alpha$	0	0.1	0.3	0.5	0.7	1.0
AP <sup>box</sup>	47.5	48.2	49.3	49.3	49.5	<b>50.0</b>
AP <sup>mask</sup>	41.8	42.6	43.4	43.4	43.5	<b>43.9</b>

(b)  $T_c(x)$ : CLIP,  $T_m(x)$ : MAE.

$\alpha$	0	0.1	0.3	0.5	0.7	1.0
AP <sup>box</sup>	49.1	49.9	49.8	50.1	50.2	<b>50.4</b>
AP <sup>mask</sup>	43.1	43.8	43.8	43.9	44.1	<b>44.1</b>

Table A5: The token masking strategy for alleviating over-fitting.  $\star$ : using MAE+DeiT teachers.  $\dagger$ : using MAE+CLIP teachers.

Method	Epoch	Masking	AP <sup>box</sup>	AP <sup>mask</sup>
Hybrid Distill $\star$	100/300		<b>50.0</b> /50.0	<b>43.9</b> /44.0
Hybrid Distill $\star$	100/300	✓	49.9/ <b>50.3</b>	43.8/ <b>44.2</b>
Hybrid Distill $\dagger$	100/300		<b>50.4</b> /50.2	<b>44.1</b> /44.1
Hybrid Distill $\dagger$	100/300	✓	50.2/ <b>50.6</b>	43.9/ <b>44.4</b>

46 However, compared to the DeiT teacher (Fig. A2(a)), it does not bring noticeable diversity gains.  
 47 Therefore, we conclude that the diversity brought by the asymmetric encoder is not always significant.

## 48 B.2 Mask Feature Reconstruction in dBOT

49 Fig. A3 compares two variants of dBOT, *i.e.*, with the same asymmetric decoder design but conducting  
 50 direct feature distillation and mask feature reconstruction, respectively. It can be seen that the two  
 51 tasks bring no significant differences, *i.e.*, the diversity is increased and the discrimination is lost  
 52 regardless of the task. These visualizations further support our claim in Sec. 2.3 and Sec 2.4 of our  
 53 main paper.

## 54 B.3 Reducing the Number of the Asymmetric Decoder Layers

55 Fig. A4 investigates the effect of reducing the number of asymmetric decoder layers. We find that  
 56 even with a reduced number of decoder layers, the discrimination in the last layer of the encoder still  
 57 cannot be maintained. Therefore, we abandon this asymmetric decoder design in our Hybrid Distill  
 58 to avoid losing discrimination.

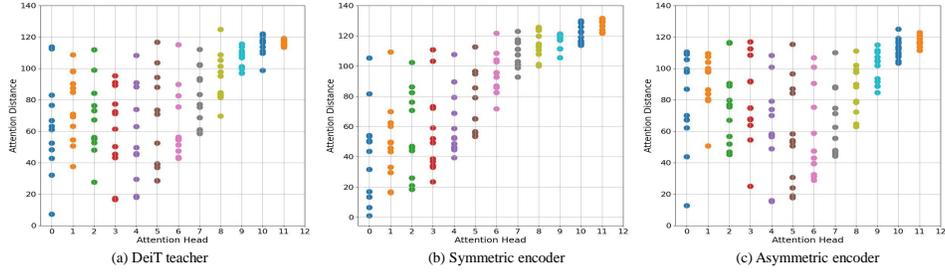


Figure A2: Average head distance of (a) DeiT teacher and student models with (b) symmetric encoder and (c) asymmetric encoder.

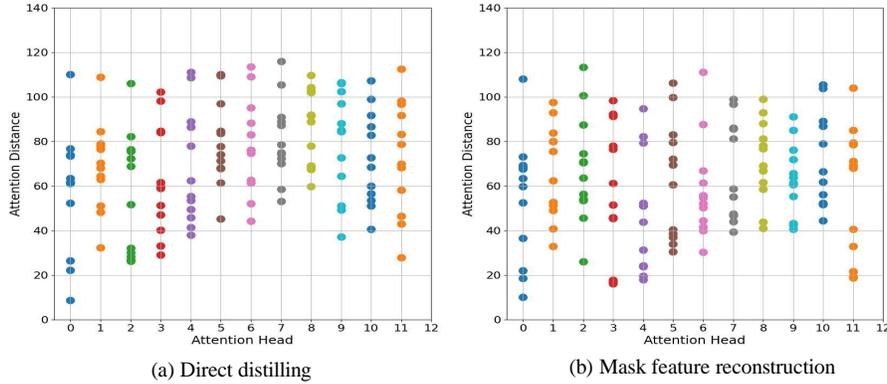


Figure A3: Average head distance of different dBOT variants that conduct (a) direct feature distillation and (b) mask feature reconstruction, respectively.

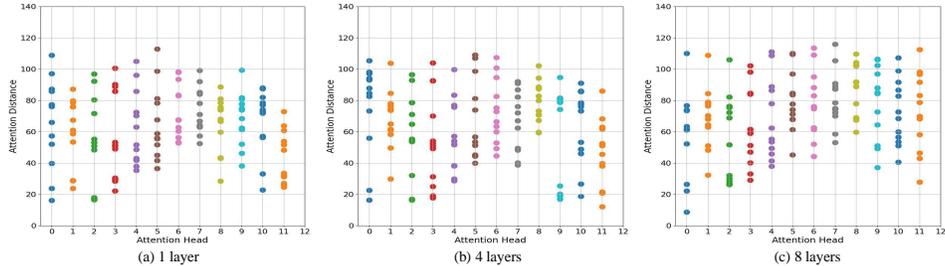


Figure A4: Average head distance of using (a) 1, (b) 4, and (c) 8 asymmetric decoder layers, respectively.

## 59 C Implementation Details for Different Downstream Tasks

60 **Classification.** We report the fine-tuning results on ImageNet-1K. Following dBOT [9], the learning  
 61 rate is set to  $3e-4$  and the batch size is set to 256. We also report results on CIFAR100 [7], Cars [6],  
 62 and iNaturalist19 [10]. For these datasets, the batch size is 768 and the learning rate is  $7.5e-6$ .

63 **Object detection and instance segmentation.** Following [1], we fine-tune the student model on  
 64 COCO [8] using the Mask-RCNN [5] framework. We train the network with the 1x schedule and the  
 65 learning rate is set to  $3e-4$  for ViT-B and  $2e-4$  for ViT-L. We also provide the 1x results using the  
 66 Cascade Mask-RCNN framework in the appendix, and the learning rate is set to  $3e-4$ .

67 **Semantic segmentation.** The semantic segmentation evaluation is conducted on ADE20K [13].  
 68 Following [1, 2], we use ViT [3] with UperNet [12] framework and fine-tune the model for 160k  
 69 iterations. The batch size, learning rate, and weight decay are set to 16,  $4e-4$ , and 0.05, respectively.

## 70 D Limitation

71 Hybrid Distill jointly utilizes two teacher models to guide the representation learning of the student.  
72 Although exhibiting promising properties and results, the additional overhead of introducing two  
73 teachers may be a limitation. Fortunately, since the teacher model does not require gradient updates,  
74 the training cost of Hybrid Distill does not increase significantly, *i.e.*, the training time of Hybrid  
75 Distill with ViT-B backbone is around 1.2 times longer than that of using a single teacher. Besides,  
76 Hybrid Distill can achieve better performance with much fewer training epochs, as shown in Tab. A2.  
77 From this perspective, Hybrid Distill in turn reduces the training cost. Another possible limitation is  
78 that Hybrid Distill does not improve CLIP as much as DeiT after introducing the MAE teacher, and  
79 we analyze that it may be caused by the gap between the pre-training capacities of CLIP and MAE  
80 teachers. We look forward to better MIM models that can further facilitate our work.

## 81 E Reproducibility

82 We will release our source code once this paper is accepted.

## 83 References

- 84 [1] Xiaokang Chen, Mingyu Ding, Xiaodi Wang, Ying Xin, Shentong Mo, Yunhao Wang, Shumin  
85 Han, Ping Luo, Gang Zeng, and Jingdong Wang. Context autoencoder for self-supervised  
86 representation learning. *arXiv preprint arXiv:2202.03026*, 2022.
- 87 [2] Yabo Chen, Yuchen Liu, Dongsheng Jiang, Xiaopeng Zhang, Wenrui Dai, Hongkai Xiong, and  
88 Qi Tian. Sdae: Self-distillated masked autoencoder. In *ECCV*, 2022.
- 89 [3] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,  
90 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al.  
91 An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint*  
92 *arXiv:2010.11929*, 2020.
- 93 [4] Yuxin Fang, Wen Wang, Binhui Xie, Quan Sun, Ledell Wu, Xinggang Wang, Tiejun Huang,  
94 Xinlong Wang, and Yue Cao. Eva: Exploring the limits of masked visual representation learning  
95 at scale. *arXiv preprint arXiv:2211.07636*, 2022.
- 96 [5] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *ICCV*, pages  
97 2961–2969, 2017.
- 98 [6] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-  
99 grained categorization. In *Proceedings of the IEEE international conference on computer vision*  
100 *workshops*, pages 554–561, 2013.
- 101 [7] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.  
102 2009.
- 103 [8] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr  
104 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Proceedings*  
105 *of the European Conference on Computer Vision (ECCV)*, pages 1209–1218, 2014.
- 106 [9] Xingbin Liu, Jinghao Zhou, Tao Kong, Xianming Lin, and Rongrong Ji. Exploring target  
107 representations for masked autoencoders. *arXiv preprint arXiv:2209.03917*, 2022.
- 108 [10] Grant Van Horn, Oisín Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig  
109 Adam, Pietro Perona, and Serge Belongie. The inaturalist species classification and detection  
110 dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
111 pages 8769–8778, 2018.
- 112 [11] Yixuan Wei, Han Hu, Zhenda Xie, Zheng Zhang, Yue Cao, Jianmin Bao, Dong Chen, and  
113 Baining Guo. Contrastive learning rivals masked image modeling in fine-tuning via feature  
114 distillation. *Tech Report*, 2022.

- 115 [12] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing  
116 for scene understanding. In *Proceedings of the European Conference on Computer Vision*  
117 (*ECCV*), pages 418–434, 2018.
- 118 [13] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio  
119 Torralba. Semantic understanding of scenes through the ade20k dataset. *International Journal*  
120 *on Computer Vision (IJCV)*, 127:302–321, 2019.