

The Role of Masking for Efficient Supervised Knowledge Distillation of Vision Transformers

Seungwoo Son  , Jegwang Ryu  , Namhoon Lee  , and Jaeho Lee 

Pohang University of Science and Technology (POSTECH)
`{swson, jgryu, namhoon, jaeho.lee}@postech.ac.kr`
<https://maskedkd.github.io/>

Abstract. Knowledge distillation is an effective method for training lightweight vision models. However, acquiring teacher supervision for training samples is often costly, especially from large-scale models like vision transformers (ViTs). In this paper, we develop a simple framework to reduce the supervision cost of ViT distillation: masking out a fraction of input tokens given to the teacher. By masking input tokens, one can skip the computations associated with the masked tokens without requiring any change to teacher parameters or architecture. We find that masking patches with the lowest student attention scores is highly effective, saving up to 50% of teacher FLOPs without any drop in student accuracy, while other masking criterion leads to suboptimal efficiency gains. Through in-depth analyses, we reveal that the student-guided masking provides a good curriculum to the student, making teacher supervision easier to follow during the early stage and challenging in the later stage.

Keywords: Knowledge Distillation · Vision Transformer · Token Pruning

1 Introduction

Large-scale vision transformers (ViTs; [8]) are becoming increasingly popular as a backbone for a wide range of visual tasks [20, 32, 35, 51], leading to an increased need for effective ways to compress these models. The classic idea of *knowledge distillation* has been found to be very effective for this purpose [15]. Many recent works have found that distilling the knowledge of large, pre-trained ViTs provides a substantial boost in the prediction quality of lightweight ViTs [13, 46, 48, 50, 56]. Furthermore, it has been shown that distillation can be combined with other model compression techniques, such as pruning or quantization, as a post-processing step that helps in recovering the pre-compression level of ViT accuracy [24, 53].

However, the computational cost of acquiring teacher supervision (i.e., *supervision cost*) for distillation can be prohibitively expensive. One needs to make multiple predictions for training samples on the teacher ViT, which typically requires more computation than what is needed to process the student model (see left of Fig. 1). For large-scale distillation tasks, such as distilling from ViT-G on ImageNet dataset, the supervision cost can be as large as 10^4 TPUv3-days [55].

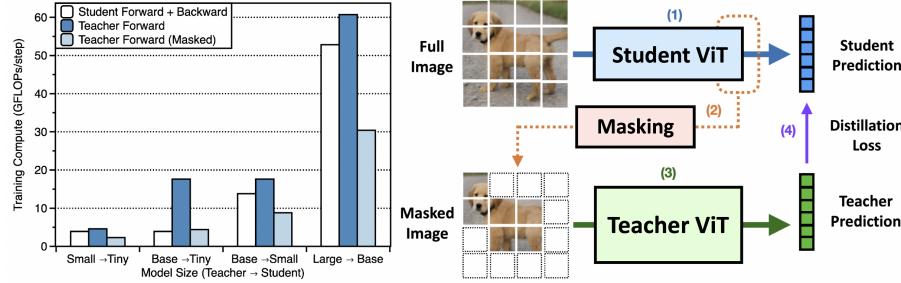


Fig. 1: (\Leftarrow) Supervision cost is expensive. We compare the per-step supervision cost of the teacher ViT that sees full/masked images, with the training FLOPs of the student. We mask the teacher input to a point where there is no student accuracy drop. Supervision cost is larger than student FLOPs, and masking can save a great amount. (\Rightarrow) **MaskedKD illustrated.** MaskedKD works in four steps: (1) Student predicts on the full image. (2) Mask the image with the student’s attention score. (3) Teacher predicts on the masked image. (4) Match the teacher-student logits.

Existing works attempt to alleviate this issue by pre-computing teacher’s predictions for all training data and reusing them during distillation [39, 54]. However, this approach not only requires ample storage to save the teacher’s predictions, but also tends to degrade student accuracy. The degradation is due to its limited ability to account for various data augmentations, such as *mixup* or RandAugment; distilling with the teacher’s predictions on unaugmented sample to the student seeing augmented data is detrimental to performance [3, 39].

To address this problem, we rethink the value of *masking* for reducing the supervision cost of knowledge distillation. Masking a fraction of input image tokens given to a vision transformer lets us skip any computations associated with the tokens, without having to modify the parameters or structures of the model. Thus, if we can mask input tokens to the pre-trained teacher in a way that the *supervision quality* (*i.e.*, the performance increase of student) remains similar, we can save much supervision cost. In short, we ask:

“How does masking affect the supervision quality of a pre-trained teacher?”

This question critically differs from two existing lines of work that also leverage the computational benefit of masking: token pruning [4, 37] and mask-based self-supervised learning (SSL) [5, 14]. Unlike token pruning, this paper seeks for the best masking strategy for retaining the supervision quality instead of the prediction quality of the teacher model. Unlike mask-based SSL, we apply masking on the inputs to the pre-trained teacher models for the supervision; in SSL, one typically masks the input tokens of the student (instead of the teacher) to use masking as a pretext task, and does not utilize pre-trained teachers. It turns out that such differences in objectives and settings necessitate a significant change in the algorithm design, as we will show in Sec. 5.

Contribution. In this paper, we develop a very simple yet effective approach to dramatically reduce the supervision cost of supervised ViT distillation, called

MaskedKD (Masked Knowledge Distillation). MaskedKD masks out teacher input tokens based on the patch saliency scores computed with the student model attention (see right of Fig. 1). This student-guided masking strategy is highly effective and versatile; MaskedKD can be combined with various distillation algorithms to cut down the ViT supervision cost by 25-50% without any degradation in the student accuracy, over various choices of model architectures.

What makes the proposed student-guided masking strategy effective? To answer this question, we conduct an in-depth comparative analysis over a number of different masking criteria. Our observations can be summarized as follows.

- ***Student-guided masking provides a distillation curriculum.*** The student-guided masking enhances the student training by playing two distinct roles. During the early phase, masking forces the teacher to give noisier supervision that is less challenging for the student to mimic (similar to teacher assistants [29]), accelerating the student training. During the later phase, masking works as a means of data augmentation, letting teachers provide diverse supervisions based on multiple views of the image (Sec. 5.1).
- ***Mask the tokens at input, not in intermediate layers.*** Unlike in token pruning literature, we find that reducing the number of input tokens is more effective than gradually removing tokens in the intermediate layers [4, 11]. While gradual removal retains the prediction quality of teachers better, the supervision quality is preserved better by masking at input (Sec. 5.2).
- ***Mask the teacher, not the student.*** For supervised knowledge distillation, we observe that masking the student model substantially degrades the final student accuracy, even at a very low masking ratio. This critically differs from the standard mask-based SSL literature, where the student always sees masked input and is trained to imitate the outputs of weight-tied teachers that see full (or masked) inputs [1, 5, 7, 23, 58] (Sec. 5.3).

2 Related Work

Distilling ViTs. Existing work on ViT distillation mostly focuses on “what to transfer” from the teacher to the student. [43] considers distilling the architectural bias of convolutional networks to enhance the student’s data-efficiency. More recent works consider distilling patch-level information of large ViTs to enhance the student performance, *e.g.*, attention-based information [48, 50, 57] or manifold structure of patches [13]. Unlike these works, our work focuses on the computational efficiency of ViT distillation.

Token removal. Removing tokens for easing the computational burden of transformer-based language models has been first proposed by [11], and much effort has followed to design similar methods for ViTs. [28, 37, 52] give algorithms to train a model that gradually removes intermediate tokens as the layer goes deeper, dramatically reducing the inference cost of the model. Instead of completely discarding tokens, [21, 26, 27] propose to combine intermediate tokens into another, instead of removing them. [4, 19] introduce a drop-in token merging module that can enhance the inference efficiency without any additional training. These works

Table 1: Comparison with mask-based SSL algorithms. We compare the masking strategy of MaskedKD with SSL algorithms that use masking as a pretext task. MaskedKD differs dramatically from these works in many senses, including the teacher type, masked model, and masking criterion. (\star : student sees multiple random crops).

	Distillation setup			Masked		Masking Criterion	Teacher type
	Super. dist.	Self-super. dist.	Teacher	Student			
MAE [14]	✗	✗	-	✓		Random	Pixel
DINO [5]	✗	✓	✗	★		Random	EMA
MSN [1]	✗	✓	✗	✓		Random	EMA
MaskFeat [49]	✗	✓	✗	✓		Random	HOG/DINO
SdAE [7]	✗	✓	✓	✓	✓	Random	EMA
PCAE [23]	✗	✓	✓	✓	✓	Activation	EMA
ccMIM [58]	✗	✓	✓	✓	✓	Shared Attn.	EMA
I-JEPA [2]	✗	✓	✓	✓	✓	Random	EMA
MaskedKD (Ours)	✓	✗	✓	✗	Student's Attn.	Supervised	

focus on preserving the prediction quality of a model. Our paper, in contrast, focuses on preserving the *supervision quality* of a teacher model.

Saving the supervision cost. A recent line of work attempts to reduce the supervision cost by re-using teacher supervisions. In particular, [39] draws inspirations from ImageNet re-labeling technique [54] to pre-compute teacher’s predictions on multiple random crops of training samples; the crop information and the supervision are stored as additional attributes of the sample. Then, the student is trained by drawing randomly cropped data and corresponding supervisions, and using distillation loss on these samples to update the model. However, this approach has limited applicability for the cases where more diverse data augmentations are used, leading to a degraded student accuracy [3, 39]. On the other hand, our work saves the supervision cost by reducing the teacher inference cost directly, applicable to the standard on-the-fly distillation scenario.

Masking as a self-supervision. Masking has been actively studied in the self-supervised learning (SSL) literature as a pretext task: the model takes the masked image as an input and is trained to predict missing parts [14] or to make similar predictions with a weight-tied model that sees the full image [1]; the latter can be viewed as a form of *self-distillation* [12], by treating the weight-tied model as a teacher and the original model as a student. Various mask-based SSL algorithms, based on diverse masking mechanisms, have been proposed; see Tab. 1 for a partial summary, and [33] for a general overview. In such works, the primary purpose is on maximizing the SSL performance rather than reducing the computation. To our knowledge, none of the works exclusively studies the computational benefit of masking for distillation, *decoupled* from how well the masking serves as the pretext task. In this work, we find that many common masking practices in SSL are suboptimal for reducing supervision cost in supervised distillation.

3 Framework: Masking for Efficient Supervised Distillation

We now describe the proposed MaskedKD framework, a very simple yet effective approach for reducing supervision cost by masking the teacher input tokens.

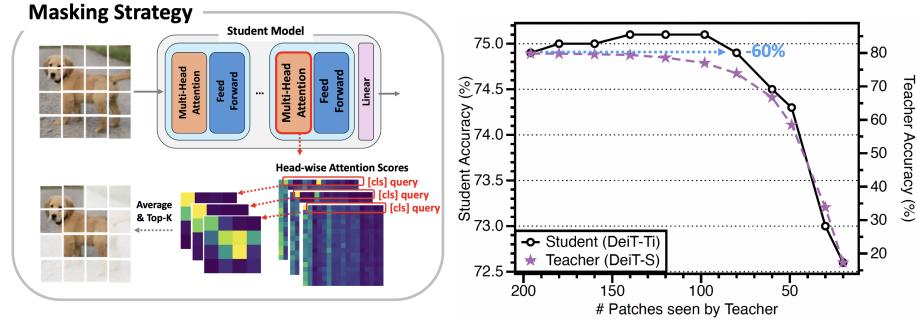


Fig. 2: (\Leftarrow) Student-guided masking illustrated. We mask the teacher’s input tokens based on the student’s attention scores of the class token query in its final layer. (\Rightarrow) **Accuracy vs. # patches seen.** Masking the teacher substantially degrades the teacher accuracy. The student accuracy, however, slightly increases first, and then starts decreasing after masking over 50% of the patches.

Formally, let $f_S(\cdot)$ and $f_T(\cdot)$ be the student and teacher ViTs, respectively. Given a training sample that consists of image-label pair (x, y) , we generate a masked version x_{mask} of the input image by removing some patches of the image x based on some masking criterion. Then, we train the student by distilling the knowledge of the teacher that predicts on the masked image, while the student keeps making predictions on the full image. In other words, the training loss is

$$\ell(x, y) = \ell_{\text{CE}}(f_S(x), y) + \lambda \cdot \ell_{\text{KD}}(f_S(x), f_T(x_{\text{mask}})), \quad (1)$$

where ℓ_{CE} denotes the cross-entropy loss, ℓ_{KD} denotes the distillation loss, and $\lambda \geq 0$ is a balancing hyperparameter. The distillation loss ℓ_{KD} can be chosen flexibly depending on the base distillation algorithm we want to use. For example, one can use the KL-divergence for the classic logit distillation [15], or the ℓ_2 distance between activations for feature distillation [50].

Remark. We note that it is essential for the final student accuracy that (1) we mask tokens at the input, and (2) we keep the student input unmasked. This choice is quite different from the masking criteria that are known to be effective for token pruning or mask-based SSL literature. We defer the detailed discussion to Secs. 5.2 and 5.3, respectively.

3.1 Masking Criterion: Student-guided Saliency Score

To select the patches to be masked, we propose to use the patch saliency score based on the student attention. More specifically, we use the last layer attention scores of the student ViT as the patch saliency score (see left of Fig. 2). This choice is inspired by the observations of [5, 6] that the last layer attention tends to be more well-aligned with human semantics.

More concretely, suppose that the input image is split into N non-overlapping patches for the student and teacher ViTs.¹ The student’s last multi-head attention layer processes $N + 1$ tokens, including the class token. For each attention head, the attention score from the class token to patch tokens is computed as

$$\mathbf{a}^{(h)} = \text{Softmax} \left((q_{\text{cls}}^\top k_1, q_{\text{cls}}^\top k_2, \dots, q_{\text{cls}}^\top k_N) / \sqrt{d} \right), \quad h \in \{1, 2, \dots, H\} \quad (2)$$

where q_{cls} is the query vector of the class token, k_i is the key vector of the i -th image patch token, d is the length of query and key vectors, and H is the number of attention heads. Given these head-wise attention scores, we compute the final patch saliency score by taking the average $\bar{\mathbf{a}} = (\sum_{h=1}^H \mathbf{a}^{(h)}) / H$.

Empirically, this masking criterion allows us to reduce the supervision cost by 25-50%, without degrading the student model accuracy (see right of Fig. 2). Experimental results over a wider range of setups will be given in Sec. 4.

Why use student attention? In principle, the score prevents masking away the core features that the student pays the most attention to, and thus allows the teacher to give well-tailored supervision to the student. Our in-depth analysis reveals that this tailored supervision provides a good training curriculum for the student; we defer the discussions to Sec. 5.1.

3.2 Computational Considerations

Computing saliency scores. Computationally, the proposed student-guided saliency score is very efficient. Computing the score requires almost no computational overhead, as the student attention scores $\mathbf{a}^{(h)}$ are already computed during the student inference. In fact, the only additional computation is taking the mean of these scores, which requires $N \cdot H$ FLOPs. This is very small even for very large students, adding less than 4.1 kFLOPs per sample for ViT-G/14.

Parallelism. Student’s attention score-based masking requires the student forward to precede the teacher forward. At first glance, this looks like an impediment toward parallel training, as it may introduce unnecessary idle time in teacher devices. However, the problem can be easily resolved by adopting micro-batching [17], readily available in most deep learning frameworks such as PyTorch or TensorFlow; see Appendix F for details.

4 Experiments

In this section, we validate that the proposed MaskedKD can save 25-50% of the supervision cost (measured in FLOPs) over a wide range of settings, without any degradation in accuracy. This section is organized as follows.

- **Sec. 4.1** describes the experimental setup for the main result.

¹ If the number of patches is different, we compute the scores for student patches, and conduct bilinear interpolation to get the scores for teacher patches.

- **Sec. 4.2** reports the performance of MaskedKD when applied for supervised ViT distillation on ImageNet-1k, over various choices of teacher and student models and base distillation algorithms. (Tab. 2).
- **Sec. 4.3** discusses the performance of MaskedKD over the boundary of supervised ViT distillation; in particular, we consider distilling the audio spectrogram transformer [10] and the distillation for the self-supervised training [5].

We also provide several additional experimental results in the appendix, including distillation from teachers trained with higher-resolution images (Appendix C), distilling only the linear classifiers (Appendix D), and measuring the effect of data augmentation in distillation (Appendix E).

4.1 Experimental Setup

As our main experimental setup, we consider the task of supervised ViT distillation with the ImageNet-1k classification dataset [38]. We experiment over various choices of models and base distillation algorithms.

Students. We use ViTs of various sizes, trained from scratch using the DeiT training recipes [43]; DeiT tend to perform better than vanilla ViTs, without requiring extra training dataset. For large-scale student models (ViT-L or ViT-H), we start from the MAE checkpoint, and fine-tune for 50 epochs [14].

Teachers. We consider several different choices of teacher:

1. DeiT (default) [43]: We use a model that has a larger size than the student.
2. CaiT [45]: A model that has a slightly different architecture than ViT.
3. CLIP [35]: A contrastively trained visual-language model that has been trained on a proprietary dataset not available to the student.
4. MAE [14]: A teacher that has been fine-tuned from the MAE; we use this for comparing with distillation algorithms that require self-supervised teachers.

Base distillation algorithms. We apply MaskedKD to these algorithms.

1. Logit (default) [15]: A popular, versatile algorithm that uses the model output for distillation (and thus usable to distill from proprietary teachers).
2. Manifold [13]: A feature distillation algorithm that regularizes the student features to have the same patch-level manifold structure as the teacher features.
3. Attention (adapted from [48]): A feature distillation algorithm that distills attention scores. Unlike [48], we apply it for the supervised distillation setup.
4. G2SD [16]: A two-stage distillation algorithm that distills during both pre-training and fine-tuning. We apply MaskedKD during the fine-tuning.

Patch size. By default, each 224×224 image is divided into total 196 patches of size 16×16 pixels. The DeiT3-H and MAE-ViT-H models use 256 patches of size 14×14 . Whenever there is a mismatch between the number of patches between the teacher and the student, we bilinearly interpolate the student’s attention score on student patches to compute the score for teacher patches.

Seeds. We report an average over 3 independent trials, except for the large-scale experiments like CLIP and those with large/huge student models.

Other details. Other experimental details are given in Appendix A.

Table 2: MaskedKD on supervised ViT distillation. MaskedKD dramatically reduces the supervision cost without degrading the student accuracy. “MaskedKD $_{\kappa\%}$ ” means that we keep only $\kappa\%$ of tokens from the teacher input. “Acc.” denotes the ImageNet top-1 accuracy of the student model. “img/s” and “PFLOPs” denote the throughput and the total supervision cost of the teacher throughout the training. † denotes that the model has an additional distillation token as an input; for this model, we do not report the performance of a model trained without distillation.

Student	Teacher	Method	Acc.	img/s	PFLOPs	Student	Teacher	Method	Acc.	img/s	PFLOPs		
DeiT-S	DeiT-Ti	-	No distillation	72.0	-	-	DeiT-S	DeiT-B	-	No distillation	79.9	-	
		Logit	75.0	1790	1770	Logit	80.8	750	6757	Logit	81.3	-	
		+MaskedKD $_{50\%}$	75.2	3702(x2.1)	866(-51%)	+MaskedKD $_{75\%}$	80.9	1038(x1.4)	4932(-27%)	+MaskedKD $_{50\%}$	81.0	1536(x2.0)	
		+MaskedKD $_{40\%}$	74.9	4642(x2.6)	707(-60%)	+MaskedKD $_{50\%}$	81.0	1536(x2.0)	3349(-50%)	+MaskedKD $_{50\%}$	81.3	-	
	DeiT3-S	Manifold	75.0	-	-	DeiT3-B	+MaskedKD $_{75\%}$	81.4	1038(x1.4)	4932(-27%)	+MaskedKD $_{50\%}$	81.3	1536(x2.0)
		+MaskedKD $_{75\%}$	75.2	2514(x1.4)	1282(-28%)	+MaskedKD $_{50\%}$	81.3	1536(x2.0)	3349(-50%)	+MaskedKD $_{50\%}$	81.8	-	
		Attention	75.3	-	-	+MaskedKD $_{50\%}$	81.8	-	-	+MaskedKD $_{50\%}$	83.5	248	
	DeiT-Ti	+MaskedKD $_{50\%}$	75.3	3702(x2.1)	866(-51%)	+MaskedKD $_{75\%}$	83.5	337(x1.4)	17301(-27%)	+MaskedKD $_{50\%}$	83.6	512(x2.1)	
		Logit	75.1	-	-	+MaskedKD $_{50\%}$	83.6	512(x2.1)	11744(-50%)	+MaskedKD $_{75\%}$	83.6	1038(x1.4)	
		+MaskedKD $_{50\%}$	75.2	3612(x2.0)	866(-51%)	+MaskedKD $_{50\%}$	83.6	1038(x1.4)	3287(-27%)	+MaskedKD $_{50\%}$	85.9	-	
DeiT-B	DeiT3-S	Logit	75.1	-	-	+MaskedKD $_{75\%}$	86.2	133	10723	+MaskedKD $_{75\%}$	86.2	179(x1.3)	
		+MaskedKD $_{50\%}$	75.2	3612(x2.0)	866(-51%)	+MaskedKD $_{50\%}$	86.2	179(x1.3)	7991(-25%)	+MaskedKD $_{50\%}$	86.1	271(x2.0)	
		+MaskedKD $_{25\%}$	74.8	6894(x3.9)	439(-75%)	+MaskedKD $_{50\%}$	86.1	271(x2.0)	5301(-51%)	+MaskedKD $_{75\%}$	86.5	-	
	DeiT3-B	Logit	74.4	750	6757	+MaskedKD $_{75\%}$	86.5	-	-	+MaskedKD $_{50\%}$	86.5	179(x1.3)	
		+MaskedKD $_{50\%}$	74.7	1536(x2.0)	3349(-50%)	+MaskedKD $_{50\%}$	86.5	179(x1.3)	7991(-25%)	+MaskedKD $_{50\%}$	86.3	271(x2.0)	
		+MaskedKD $_{25\%}$	74.7	2882(x3.8)	1696(-75%)	+MaskedKD $_{50\%}$	86.3	271(x2.0)	5301(-51%)	+MaskedKD $_{75\%}$	86.9	-	
	CaiT-S24	Logit	75.1	528	3591	+MaskedKD $_{75\%}$	86.9	-	-	+MaskedKD $_{75\%}$	87.2	-	
		+MaskedKD $_{75\%}$	75.2	807(x1.5)	2583(-28%)	+MaskedKD $_{50\%}$	86.9	-	-	+MaskedKD $_{50\%}$	87.2	-	
		+MaskedKD $_{50\%}$	75.2	1018(x1.9)	1728(-52%)	+MaskedKD $_{50\%}$	86.9	-	-	+MaskedKD $_{75\%}$	87.2	-	
MAE-ViT-L 1	CLIP-B/16	Manifold	75.7	-	-	+MaskedKD $_{75\%}$	86.9	-	-	+MaskedKD $_{75\%}$	87.2	-	
		+MaskedKD $_{75\%}$	75.9	807(x1.5)	2583(-28%)	+MaskedKD $_{50\%}$	86.9	-	-	+MaskedKD $_{50\%}$	87.2	-	
		Logit	73.9	-	-	+MaskedKD $_{50\%}$	86.9	-	-	+MaskedKD $_{50\%}$	87.2	-	
	ConvViT-B	+MaskedKD $_{75\%}$	75.2	1017(x1.4)	4932(-27%)	+MaskedKD $_{50\%}$	86.9	-	-	+MaskedKD $_{75\%}$	87.2	-	
		+MaskedKD $_{50\%}$	75.2	1536(x2.0)	3349(-50%)	+MaskedKD $_{50\%}$	86.9	-	-	+MaskedKD $_{50\%}$	87.2	-	
		Logit	73.8	432	8543	+MaskedKD $_{50\%}$	86.9	-	-	+MaskedKD $_{75\%}$	87.2	-	
	Swin-Ti	+MaskedKD $_{50\%}$	74.2	590(x1.4)	5540(-33%)	+MaskedKD $_{50\%}$	86.9	-	-	+MaskedKD $_{50\%}$	87.2	-	
		Logit	81.7	750	6757	+MaskedKD $_{75\%}$	87.2	179(x1.3)	7991(-25%)	+MaskedKD $_{50\%}$	87.1	512(x2.0)	
		+MaskedKD $_{75\%}$	81.8	1038(x1.4)	4932(-27%)	+MaskedKD $_{50\%}$	87.1	512(x2.0)	5301(-51%)	+MaskedKD $_{75\%}$	87.1	5301(-51%)	

4.2 Main Results

In Tab. 2, we provide the performances of MaskedKD when applied to various supervised distillation scenarios.² From the table, we observe that we can safely remove 25–50% of the patches from the teacher input, without sacrificing the student accuracy; in some cases, we can even remove 60–75% of the patches. We also observe that, in most cases, masking a small fraction of patches is beneficial for the performance of the trained student (although with a very small boost). Such performance gain from masking is most pronounced when the size gap between the teacher and the student is large (DeiT3-B → DeiT-Ti). From this observation, we hypothesize that the masking may have an effect similar to making a (low-capacity) *teaching assistant* model [29], which can teach the student better than an overly large teacher. We validate this intuition in Sec. 5.

4.3 Extending the boundaries of MaskedKD

Here, we test the performance of the MaskedKD for the tasks other than the supervised ViT distillation. In particular, we apply MaskedKD to (1) the distil-

² For G2SD, the student accuracy we obtained for the baseline (81.5%) is slightly lower than what is reported in the original paper. We could not reproduce the reported accuracy, although we used the official code.

Table 3: Audio classification. We apply MaskedKD on distilling the audio spectrogram transformer [10] for the audio classification on ESC-50 [34].

Student	Method	Acc. (%)	TFLOPs
	No distillation	85.1	-
AST-S	Logit	86.3	92.8
	+MaskedKD 83%	86.4	75.9
	+MaskedKD 50%	85.7	44.0

Table 4: Masking DINO. We apply MaskedKD on the self-supervised training procedure of DINO [5]. “Linear” denotes the linear probing accuracy and “PFLOPs” denotes the total teacher FLOPs.

Model	Method	Linear (%)	PFLOPs
	DINO	73.7	589
DeiT-S	+MaskedKD _{87%}	73.9	507
	+MaskedKD _{77%}	74.0	446
	+MaskedKD _{66%}	74.2	384
	+MaskedKD _{61%}	73.6	354

lation of audio spectrogram transformers, and (2) the self-supervised learning algorithms that utilize some form of self-distillation.

Distilling audio transformers. Here, we check whether the MaskedKD can also be applied for an efficient distillation of the transformer that processes data from other domains. To this end, we consider distilling an audio spectrogram transformer (AST) [10] on the ESC-50 audio classification dataset [34]. AST has an identical architecture to DeiT, but uses a modified training recipe and hyperparameters tailored for processing the spectrograms of the speech data. We give a more detailed description of the experimental setups in Appendix A.

Tab. 3 compares the performance of the vanilla logit distillation algorithm against the distillation with MaskedKD. We observe that the MaskedKD can indeed reduce the number of patches used by 17% without sacrificing the performance (masking 100 patches out of 600). The supervision cost reduction, however, is not as much as in the vision domain. The performance drops by 0.6% if we use only 50% of the patches, unlike most vision transformers.

Self-distillation for self-supervised learning. Here, we apply MaskedKD to reduce the computations of a distillation-based self-supervised learning algorithm, DINO [5]. A line of self-supervised learning literature aims to train useful image representations by distilling the knowledge from a teacher model—generated as a moving average of the student—that sees a differently augmented version from the student [12]. Since its advent, such *self-training* algorithms for self-supervised learning have been one of the major use cases of knowledge distillation. In the context of ViT, DINO [5] is one of the most prominent algorithms in the direction. DINO employs a teacher that sees the full image, a student that sees the full image as well, and additional students having smaller field of views (FOVs); then, DINO performs the self-distillation to regularize the teacher and students to give similar outputs.

We apply MaskedKD to DINO as follows: We use only the attention scores of a student that sees the full image to compute the patch saliency, and mask the teacher input. As the student models with smaller FOVs have a much smaller computational cost for inference, masking the teacher can save a substantial portion of the whole training cost.

Tab. 4 compares the performance of the vanilla DINO against MaskedKD, where we use ImageNet-1k dataset for both the pre-training phase (100 epochs) and fine-tuning phase. We observe that we can mask away more than 30% of

the patches from the teacher input without a degradation in the quality of the learned representation; the quality of representation is measured by the accuracy achievable by linear probing, i.e., only the linear classifier is fine-tuned for the task. We can save the teacher computation accordingly, cutting down the total teacher FLOPs from 589PFLOPs to 384PFLOPs when we use 66% of the patches.

5 A Closer Look at the Masked Knowledge Distillation

In this section, we take a closer look at the proposed MaskedKD. We conduct an in-depth analysis on why the proposed student-guided masking works well, and demystify several design choices for the MaskedKD. In particular, we answer:

- How should we mask? Student-guided masking is an effective strategy that enhances the student optimization via implicit curriculum; the mask makes training easier during the early phase, and difficult in the late phase (Sec. 5.1).
- Where should we mask? Masking tokens at input is effective for preserving supervision quality, while removing tokens at intermediate layers is better for keeping high prediction quality (Sec. 5.2).
- Who should we mask? Even at the low masking ratio, masking student degrades the final accuracy after distillation, while masking teacher does not. This contrasts with mask-based SSL, where masking student is essential (Sec. 5.3).

We provide additional ablations in Sec. 5.4.

5.1 How Should We Mask?

Answer: Student-guided masking. We compare various masking mechanisms, and find that the student-guided masking (used in MaskedKD) is the most effective masking criterion. In particular, we consider four different mechanisms:

- (a) Student (ours): The mask used in MaskedKD, based on student attention.
- (b) Teacher: Same as “Student,” but uses the teacher attention instead of student.
Note that, practically, this is computationally inefficient since computing this requires an additional teacher forward.
- (c) DINO: Same as “Student,” but uses the attention of DINO [5]. This procedure is known to retain highly semantically meaningful patches.
- (d) Random: Randomly mask the patches with uniform probability.

We compare how the student and (masked) teacher perform, in the following setup: We distill the knowledge of DeiT-Base teacher to a DeiT-Small student [43]. We train on ImageNet dataset, using the logit distillation [15]. By default, we set the masking ratio to 50%. The results (Fig. 3, left) show that student-guided masking is the only masking criterion that achieves similar to or better accuracy than the vanilla logit distillation.

Follow-up Question. What makes the student-guided masking work well?

Answer: It provides a good curriculum for distillation. In the early stage, student-guided masking makes teacher supervision less challenging-to-fit. In the late stage, student-guided masking provides diverse views to the teacher, helping the student learn more comprehensively from the teacher.

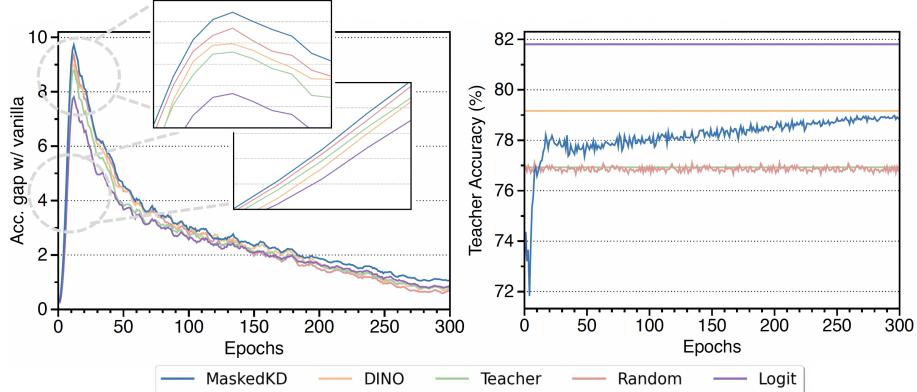


Fig. 3: (↔) **Student accuracy gain from distillation, with various masks.** The student-guided masking achieves the most rapid student accuracy increase in the early phase (averaged over three seeds). (⇒) **Accuracy of masked teachers.** The student-guided masking leads to the lowest teacher accuracy during the early phase.

To see what happens in the early stage, we report in Fig. 3 the training curves for the student accuracy and teacher accuracy. We first observe that, regardless of how we mask, masking the teacher is beneficial in the early epochs. Interestingly, the scale of the early gain seems to be negatively correlated with the teacher accuracy; MaskedKD and Random are the worst in teacher accuracy, but best in terms of student accuracy. These observations suggest that the early benefits of masked distillation may be attributed to the implicit curriculum it provides; masking makes the teachers less accurate, making it easier for the students to mimic their predictions [18, 25, 29]. In this sense, MaskedKD can be viewed as a computation-efficient way to implement curriculum learning.

In the late stage, we find that the student-guiding lets the teacher provide supervision on diverse views of the image, which is beneficial for the student performance. As Tab. 5 demonstrates, having the teacher supervise on more diverse views in the late phase can improve the student performance. In this sense, the student-guided masking is very effective, as it allows the teacher to predict on diverse views, while not losing focus on the core features (see Fig. 4).

Remark. An intriguing observation is that the student-guided masking degrades the teacher accuracy more severely than random masking in the early stage (Fig. 3, right). To explain this phenomenon, we analyze the mask structures given by randomly initialized students (Fig. 5). We make two observations: (1) Even at random initialization, the student tends to mask similar patches at the same time, often masking away all patches of the foreground object at once (upper row). In contrast, random masking rarely masks out the entire foreground object. (2) During the early stage, the student-guided masking tends to preserve more peripheral patches than central patches (lower row). This, in the early stage, makes it difficult for the teacher to supervise on foreground objects, which are typically located at the center of the given image.

Table 5: Random patch selection helps DINO masking in the late training. We compare DINO with a “DINO+Random,” a version which, for the later half of the training, we draw 40% of the patches with high DINO attention and the remaining 10% randomly. We observe that randomization in late training helps.

DINO	DINO + Random
80.7%	80.9% (+0.2%)

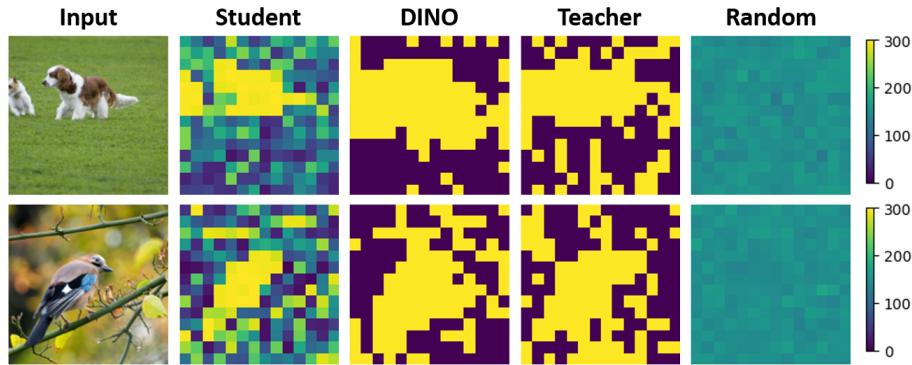


Fig. 4: Student-guided masking lets teacher supervise on diverse views. We visualize the utilization frequency of each patch throughout the training. The student-guided masking lets teacher predict on diverse patches (unlike “Teacher” or “DINO”), while conveying the core semantic information of the image (unlike “Random”).

5.2 Where Should We Mask?

Answer: At the input, not in the intermediate layers. We find that the proposed MaskedKD, which masks the teacher at the input, provides better supervision than teachers whose tokens have been removed in the intermediate layers. In particular, we compare with ToMe [4], one of the most popular token removal algorithms; this algorithm does not require any further training of the teacher model, and thus can provide a fair comparison. We compare two algorithms in the setup where we distill the knowledge of DeiT-Small to DeiT-Tiny on ImageNet.

The experimental results are given in the left two panes of Fig. 6. Here, we observe that MaskedKD outperforms ToMe-based distillation. This, however, does not mean that student-masked teachers also predict better; ToMe teacher achieves better accuracy in the high-computation regime. Interestingly, MaskedKD also predicts better in the low-computation regime. Essentially, this is because one needs to remove a larger number of tokens in the intermediate layers to have similar reductions in computations.

5.3 Who Should We Mask?

Answer: Mask the teacher, and not the student. In the right of Fig. 6, we compare the student accuracy of MaskedKD with its variant where the student is

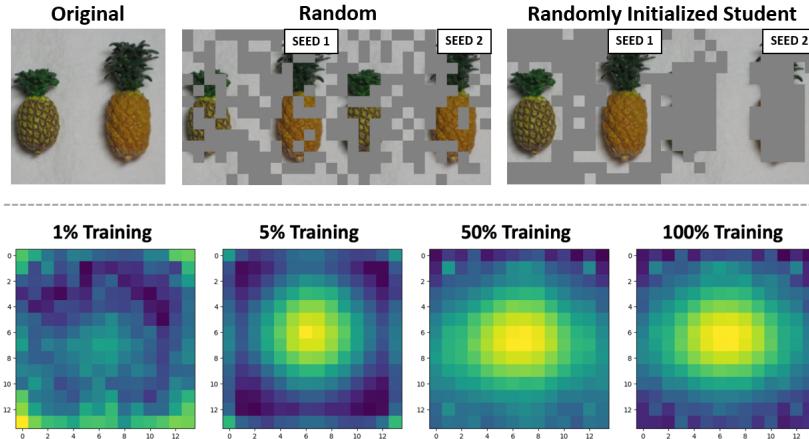


Fig. 5: (↑) **Randomly initialized students can mask similar patches at once.** Randomly initialized students tend to mask all similar patches at the same time, often removing all foreground or background objects at once. We provide additional examples in Appendix B. (↓) **Student shifts attention from periphery to center.** We visualize the patch selection frequency of MaskedKD at different stages of training. Early in training, the student attend more on peripheral patches. As the training proceeds, the student shifts the attention to the central region.

masked instead of the teacher. We observe that masking the student during distillation immediately degrades the student performance, even at a very low masking rate. In contrast, masking the teacher (as in MaskedKD) does not degrade the student accuracy until one masks over 50% of all patches. This illustrates the crucial difference in the role of masking in supervised distillation and mask-based self-supervised learning algorithms, e.g., [1, 7].

5.4 Additional ablations

In Tab. 6, we ablate various components of the MaskedKD. From the table, we draw three conclusions. (a) Using the class-patch attention leads to better masking than using (an average of) the patch-patch attention score. Using class-patch attention also introduces less computational overhead. (b) Using the attention score from the last layer works better than using the attention score from the preceding layers. This observation is well-aligned with the observation by [5] that the attention becomes more focused and semantic in the last transformer block. (c) The student’s attention score is well-correlated with the accuracy of the student. Making the teacher predict on bottom- k patches introduces a large degradation in the student performance.

6 Discussion

In this work, we develop a simple yet effective framework to reduce the supervision cost of the ViT distillation. The proposed MaskedKD masks the teacher input

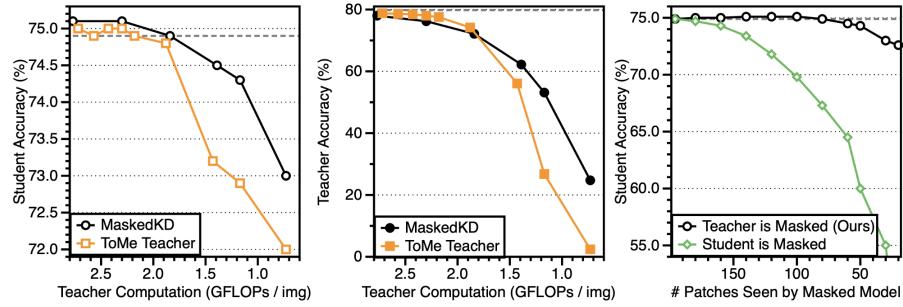


Fig. 6: (\leftarrow, \uparrow) MaskedKD vs. ToMe. Applying ToMe to the teacher also reduces the supervision cost, but MaskedKD achieves a better accuracy-computation tradeoff than ToMe (\leftarrow). This happens both in high-compute and low-compute regime, regardless of which teacher predicts better (\uparrow). **(\Rightarrow) Masking teacher vs. student.** Masking the student substantially degrades the student accuracy, but masking the teacher does not.

Table 6: Additional ablations. We validate various components of the MaskedKD: Extracting ^(a)class-patch tokens from the ^(b)last layer of the student, and removing all patches except for ^(c)top-k score. We use DeiT-S as a teacher and DeiT-Ti as a student, and mask away 50% of all tokens. The default MaskedKD is marked in purple .

method	top-1	top-5	method	top-1	function	top-1	top-5
[cls]-patch	75.1	92.1	first (1)	72.3	top-k	75.1	92.1
patch-patch	74.6	91.9	middle (6)	75.0	random	74.6	91.3
(a) Class or patch token?						bottom-k	71.7
(b) Which layer?						(c) Random & Bottom-k	

based on the student model attention, and can reduce the supervision cost by 25-50% over a wide range of setups. A key limitation of the present work is its scope: MaskedKD is specialized for the supervised distillation of transformer-based models. Generalizing the proposed framework to cover a broader range of models and training setups is an important future research direction.

Potential negative societal impact. The proposed framework selects a fraction of tokens for distillation using an attention-based mechanism. As the mechanism implicitly determines which feature is important and which is not for distillation, it may potentially capture and strengthen the spurious correlations in the dataset.

Acknowledgment. This work was partly supported by the National Research Foundation of Korea (NRF) grant funded by the Korean government (MSIT) (RS2023-00213710, RS2023-00210466), and the Institute of Information & communications Technology Planning & Evaluation (IITP) grant funded by the Korean government (MSIT) (RS-2019-II191906, Artificial Intelligence Graduate School Program (POSTECH), RS-2022-II220959, Few-Shot learning of Causal Inference in Vision and Language for Decision Making), and POSCO Creative Ideas grant (2023Q024, 2023Q032).

References

1. Assran, M., Caron, M., Misra, I., Bojanowski, P., Bordes, F., Vincent, P., Joulin, A., Rabbat, M., Ballas, N.: Masked Siamese networks for label-efficient learning. In: Proceedings of the European Conference on Computer Vision (2022)
2. Assran, M., Duval, Q., Misra, I., Bojanowski, P., Vincent, P., Rabbat, M., LeCun, Y., Ballas, N.: Self-supervised learning from images with a joint-embedding predictive architecture. In: Proceedings of the International Conference on Computer Vision. pp. 15619–15629 (2023)
3. Beyer, L., Zhai, X., Royer, A., Markeeva, L., Anil, R., Kolesnikov, A.: Knowledge distillation: A good teacher is patient and consistent. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2022)
4. Bolya, D., Fu, C.Y., Dai, X., Zhang, P., Feichtenhofer, C., Hoffman, J.: Token merging: Your ViT but faster. In: International Conference on Learning Representations (2023)
5. Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., Joulin, A.: Emerging properties in self-supervised vision transformers. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 9650–9660 (2021)
6. Chefer, H., Schwartz, I., Wolf, L.: Optimizing relevance maps of vision transformers improves robustness. In: Advances in Neural Information Processing Systems (2022)
7. Chen, Y., Liu, Y., Jiang, D., Zhang, X., Dai, W., Xiong, H., Tian, Q.: SdAE: Self-distillated masked autoencoder. In: Proceedings of the European Conference on Computer Vision (2022)
8. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. In: International Conference on Learning Representations (2020)
9. Fang et al.: You only look at one sequence: Rethinking transformer in vision through object detection. In: NIPS (2021)
10. Gong, Y., Chung, Y.A., Glass, J.: AST: Audio spectrogram transformer. In: Interspeech (2021)
11. Goyal, S., Choudhury, A.R., Raje, S.M., Chakaravarthy, V.T., Sabharwal, Y., Verma, A.: POWER-BERT: Accelerating BERT inference via progressive word-vector elimination. In: Proceedings of the International Conference on Machine Learning (2020)
12. Grill, J.B., Strub, F., Altché, F., Tallec, C., Richemond, P.H., Buchatskaya, E., Doersch, C., Pires, B.A., Guo, Z.D., Azar, M.G., Piot, B., Kavukcuoglu, K., Munos, R., Valko, M.: Bootstrap your own latent: A new approach to self-supervised learning. In: Advances in Neural Information Processing Systems (2020)
13. Hao et al.: Learning efficient vision transformers via fine-grained manifold distillation. In: NeurIPS (2022)
14. He, K., Chen, X., Xie, S., Li, Y., Dollár, P., Girshick, R.: Masked autoencoders are scalable vision learners. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2022)
15. Hinton, G., Vinyals, O., Dean, J.: Distilling the knowledge in a neural network. arXiv preprint arXiv:1503.02531 (2015)
16. Huang, W., Peng, Z., Dong, L., Wei, F., Jiao, J., Ye, Q.: Generic-to-specific distillation of masked autoencoders. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2023)

17. Huang, Y., Cheng, Y., Bapna, A., Firat, O., Chen, M.X., Chen, D., Lee, H., Ngiam, J., Le, Q.V., Wu, Y., Chen, Z.: GPipe: Efficient training of giant neural networks using pipeline parallelism. In: Advances in Neural Information Processing Systems (2019)
18. Jin, X., Peng, B., Wu, Y., Liu, Y., Liu, J., Liang, D., Yan, J., Hu, X.: Knowledge distillation via route constrained optimization. In: Proceedings of the International Conference on Computer Vision (2019)
19. Kim, M., Gao, S., Hsu, Y.C., Shen, Y., Jin, H.: Token fusion: Bridging the gap between token pruning and token merging. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 1383–1392 (2024)
20. Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A.C., Lo, W.Y., Dollár, P., Girshick, R.: Segment anything. In: Proceedings of the International Conference on Computer Vision (2023)
21. Kong, Z., Dong, P., Ma, X., Meng, X., Niu, W., Sun, M., Ren, B., Qin, M., Tang, H., Wang, Y.: SPViT: Enabling faster vision transformers via soft token pruning. In: European Conference on Computer Vision (2021)
22. Lee, Y., Chen, A.S., Tajwar, F., Kumar, A., Yao, H., Liang, P., Finn, C.: Surgical fine-tuning improves adaptation to distribution shifts. In: International Conference on Learning Representations (2023)
23. Li, J., Wang, Y., Zhang, X., Chen, Y., Jiang, D., Dai, W., Li, C., Xiong, H., Tian, Q.: Progressively compressed auto-encoder for self-supervised representation learning. In: International Conference on Learning Representations (2022)
24. Li, Y., Xu, S., Zhang, B., Cao, X., Gao, P., Guo, G.: Q-ViT: Accurate and fully quantized low-bit vision transformer. In: Advances in Neural Information Processing Systems (2022)
25. Li, Z., Li, X., Yang, L., Zhao, B., Song, R., Luo, L., Li, J., Yang, J.: Curriculum temperature for knowledge distillation. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 37, pp. 1504–1512 (2023)
26. Liang, Y., Ge, C., Tong, Z., Song, Y., Wang, J., Xie, P.: Not all patches are what you need: Expediting vision transformers via token reorganizations. In: International Conference on Learning Representations (2022)
27. Marin, D., Chang, J.H.R., Ranjan, A., Prabhu, A., Rastegari, M., Tuzel, O.: Token pooling in vision transformers for image classification. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (2023)
28. Meng, L., Li, H., Chen, B.C., Lan, S., Wu, Z., Jiang, Y.G., Lim, S.N.: AdaViT: Adaptive vision transformers for efficient image recognition. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2022)
29. Mirzadeh, S.I., Farajtabar, M., Li, A., Levine, N., Matsukawa, A., Ghasemzadeh, H.: Improved knowledge distillation via teacher assistant. In: Proceedings of the AAAI Conference on Artificial Intelligence (2020)
30. Park, N., Kim, W., Heo, B., Kim, T., Yun, S.: What do self-supervised vision transformers learn? In: International Conference on Learning Representations (2023)
31. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al.: PyTorch: An imperative style, high-performance deep learning library. In: Advances in Neural Information Processing Systems (2019)
32. Peebles, W., Xie, S.: Scalable diffusion models with transformers. In: Proceedings of the International Conference on Computer Vision (2023)
33. Peng, Z., Dong, L., Bao, H., Wei, F., Ye, Q.: A unified view of masked image modeling. Transactions on Machine Learning Research (2023)

34. Piczak, K.J.: ESC: Dataset for environmental sound classification. In: Proceedings of the ACM International Conference on Multimedia (2015)
35. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., Sutskever, I.: Learning transferable visual models from natural language supervision. In: Proceedings of the International Conference on Machine Learning (2021)
36. Ramon van Handel: On the spectral norm of Gaussian random matrices. Transactions of AMS (2017)
37. Rao, Y., Zhao, W., Liu, B., Lu, J., Zhou, J., Hsieh, C.J.: DynamicViT: Efficient vision transformers with dynamic token sparsification. In: Advances in Neural Information Processing Systems (2021)
38. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Li, F.F.: ImageNet large scale visual recognition challenge. International Journal of Computer Vision (2015)
39. Shen, Z., Xing, E.P.: A fast knowledge distillation framework for visual recognition. In: Proceedings of the European Conference on Computer Vision (2022)
40. Song et al.: Vidit: An efficient and effective fully transformer-based object detector. In: ICLR (2022)
41. Steiner, A., Kolesnikov, A., Zhai, X., Wightman, R., Uszkoreit, J., Beyer, L.: How to train your ViT? Data, augmentation, and regularization in vision transformers. Transactions on Machine Learning Research (2022)
42. Strudel, R., Garcia, R., Laptev, I., Schmid, C.: Segmenter: Transformer for semantic segmentation. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 7262–7272 (2021)
43. Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., Jégou, H.: Training data-efficient image transformers & distillation through attention. In: Proceedings of the International Conference on Machine Learning (2021)
44. Touvron, H., Cord, M., Jégou, H.: DeiT III: Revenge of the ViT. In: Proceedings of the European Conference on Computer Vision (2022)
45. Touvron, H., Cord, M., Sablayrolles, A., Synnaeve, G., Jégou, H.: Going deeper with image transformers. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2021)
46. Vasu, P.K.A., Pouransari, H., Faghri, F., Vemulapalli, R., Tuzel, O.: MobileCLIP: Fast image-text models through multi-modal reinforced training. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2024)
47. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. In: Advances in Neural Information Processing Systems (2017)
48. Wang, K., Yang, F., van de Weijer, J.: Attention distillation: Self-supervised vision transformer students need more guidance. In: British Machine Vision Conference (2022)
49. Wei, C., Fan, H., Xie, S., Wu, C.Y., Yuille, A., Feichtenhofer, C.: Masked feature prediction for self-supervised visual pre-training. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2022)
50. Wu, H., Gao, Y., Zhang, Y., Lin, S., Xie, Y., Sun, X., Li, K.: Self-supervised models are good teaching assistants for vision transformers. In: Proceedings of the International Conference on Machine Learning (2022)
51. Yang, L., Kang, B., Huang, Z., Xu, X., Feng, J., Zhao, H.: Depth anything: Unleashing the power of large-scale unlabeled data. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2024)

52. Yin, H., Vahdat, A., Alvarez, J.M., Mallya, A., Kautz, J., Molchanov, P.: A-ViT: Adaptive tokens for efficient vision transformer. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2022)
53. Yu, S., Chen, T., Shen, J., Yuan, H., Tan, J., Yang, S., Liu, J., Wang, Z.: Unified visual transformer compression. In: International Conference on Learning Representations (2022)
54. Yun, S., Oh, S.J., Heo, B., Han, D., Choe, J., Chun, S.: Re-labeling ImageNet: From single to multi-labels, from global to localized labels. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2021)
55. Zhai, X., Kolesnikov, A., Houlsby, N., Beyer, L.: Scaling vision transformers. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2022)
56. Zhang, C., Han, D., Qiao, Y., Kim, J.U., Bae, S.H., Lee, S., Hong, C.S.: Faster segment anything: Towards lightweight SAM for mobile applications. arXiv preprint arXiv:2306.14289 (2023)
57. Zhang, J., Peng, H., Wu, K., Liu, M., Xiao, B., Fu, J., Yuan, L.: MiniViT: Compressing vision transformers with weight multiplexing. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12145–12154 (2022)
58. Zhang, S., Zhu, F., Zhao, R., Yan, J.: Contextual image masking modeling via synergized contrasting without view augmentation for faster and better visual pretraining. In: International Conference on Learning Representations (2023)

A Details on the experimental setup

Training recipe and data augmentations. For training student ViTs, we follow the settings of DeiT [43], except for the tiny model; we find that Tiny ViT tends to achieve better accuracy with less data augmentations, similar to [41]. In particular, we only use random resized crops and horizontal flips for tiny model, without RandAugment, *mixup* and *cutmix*.

Hardware. Throughputs reported in this paper are measured on a single NVIDIA RTX A6000 graphic card using FP32 weights/activations and the batch size 128.

Distillation hyperparameters. For the balancing hyperparameter λ and the temperature scale τ (for logits), we use $(1.0, 1)$ unless other noted otherwise; we have tuned the hyperparameters over the search space $\{0.1, 1.0, 9.0\} \times \{1, 2, 3, 4\}$, and observe that $(1.0, 1)$ works well throughout all setups. Similar observations about the hyperparameters have been made in [13, 50, 57].

Manifold distillation. While the original paper only uses CaiT as the teacher model, we also experiment with DeiT teachers; we use the same hyperparameters for training with DeiT teachers.

Attention distillation. We adapt the attention distillation to the supervised learning setup by distilling the attention from all layers (with the same scaling factors), rather than distilling only the last layer.

Audio experiments. We use the official code of AST³ and follow the same training procedure; the model’s performance is evaluated by averaging 5 seed validation results.

DINO. We follow the smaller-scale experimental setup for DINO, available at the official code repository⁴, where we train for 100 epochs. We also halved the number of GPUs (from 8 to 4) and the per-GPU batch size (from 256 to 128), due to the limitations in the computing resource available.

³ <https://github.com/YuanGongND/ast>

⁴ The “vanilla DINO training” in <https://github.com/facebookresearch/dino/>

B Additional visualizations of masking by randomly initialized students

In Figure 7, we present further examples of images where 50% of the patches have been masked by DeiT-Small students initialized randomly.

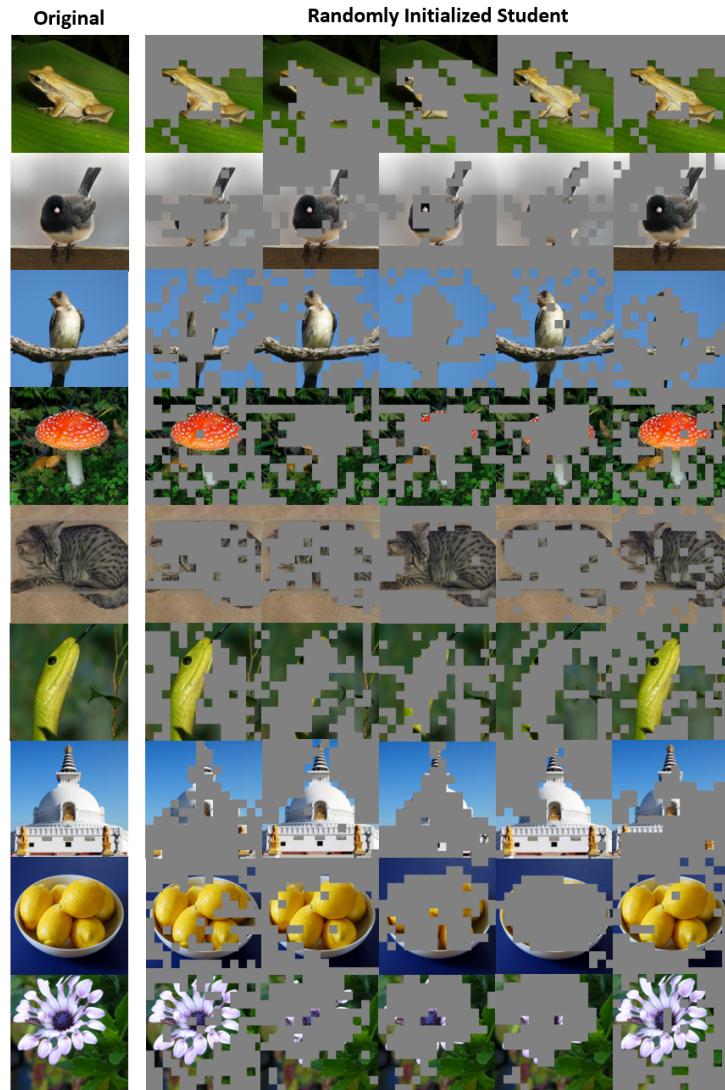


Fig. 7: Images masked by randomly initialized student

Table 7: MaskedKD with teachers that take higher-resolution images. We validate the effectiveness of MaskedKD with the teacher model that takes in higher-resolution images as inputs. Teachers are trained with 384×384 resolution images.

Student	Teacher	Method	Acc.	PFLOPs
DeiT-Ti	DeiT3-S @ 384	Logit	74.9	5694
		+MaskedKD _{75%}	75.3	4270 (-28%)
		+MaskedKD _{50%}	74.9	2725 (-54%)
		+MaskedKD _{34%}	74.6	1814 (-70%)
DeiT-Ti	DeiT3-S @ 224	Logit	74.9	1771
		+MaskedKD _{50%}	75.1	866 (-51%)
		+MaskedKD _{25%}	74.8	439 (-75%)

C Teachers trained with higher-resolution images

One of the key factors that govern the inference compute of a model is the *input resolution*. While the models that are trained on higher-resolution images tend to work better [44], the models tend to require much higher computational cost. We ask whether the proposed MaskedKD can be used to distill the knowledge from the teachers that are trained on higher-resolution images, without inducing an excessive overhead in the training cost.

To this end, we consider the following MaskedKD pipeline for distilling from the ViT teachers that takes a higher-dimensional input.⁵ More specifically, we consider distilling ViTs trained with 384×384 images, which are processed by dividing into total 576 patches of size 16×16 .

- (1) Given a low-resolution image, forward the image through the student model to get the predictions and the mean attention scores for the patches; here, if we use 224×224 input image, then we get 196 scores.
- (2) Interpolate the mean attention scores to generate the scores for larger number of patches (e.g., 576 patches). We use bilinear interpolation.
- (3) Interpolate the input image to generate a high-resolution image and mask it with the interpolated mean attention scores computed in the previous step.
- (4) Forward the masked image through the teacher model and use the prediction to regularize the student training.

We give the experimental results in Tab. 7. Here, we observe that one can successfully distill the knowledge from a teacher trained with higher-resolution images. Also, we observe that MaskedKD successfully reduces the computational burden for distillation, without any performance drop. We note, however, that we did not count the FLOPs that are required for the interpolation procedures in step (2,3). Also, the forward FLOPs for the teacher trained with higher resolution images are too big, so that the performance of the MaskedKD-distilled model does not match the performance of the model distilled using vanilla KD + low-resolution teacher that has a similar training FLOPs.

⁵ We note that, up to our knowledge, this is the first attempt distilling a high-resolution teacher to a low-resolution student.

Table 8: MaskedKD for distilling only the linear classifiers. We apply MaskedKD to a scenario where we only fine-tune the linear classifier of the student, whose (frozen) feature map has been pre-trained with self-supervision.

Student	Teacher	Method	Acc.	PFLOPs
DINO-ViT-S	DeiT-B	-	76.9	-
		No distillation	76.9	-
		Logit	77.5	2252
	DeiT3-L	+MaskedKD 75%	77.5	1644(-27%)
		+MaskedKD 50%	77.5	1116(-49%)
	DeiT3-L	Logit	77.5	7892
		+MaskedKD 75%	77.6	5767(-27%)
		+MaskedKD 50%	77.5	3914(-50%)
DINO-ViT-B	DeiT-B	-	77.9	-
		No distillation	77.9	-
		Logit	78.1	2252
	DeiT3-L	+MaskedKD 75%	78.1	1644(-27%)
		+MaskedKD 50%	78.1	1116(-49%)
	DeiT3-L	Logit	78.2	7892
		+MaskedKD 75%	78.3	5765(-27%)
		+MaskedKD 50%	78.2	3914(-50%)

D MaskedKD for distilling linear classifiers of self-supervised models

We also consider the scenario where we distill the knowledge of a supervisedly trained teacher to a student whose feature map has been pre-trained with a self-supervised learning scheme and frozen; during the distillation, we only fine-tune the linear classifier of the student. In such scenario, reducing the teacher computation gains even greater importance, as the computational cost for the student backward is greatly diminished. Previous studies on knowledge distillation primarily focuses on cases where the entire model is fine-tuned. However, several recent studies show that fine-tuning the entire model may be suboptimal in some cases [22, 30].

Tab. 8 provides the experimental results. We observe that KD still provides performance boost under this setup, and the efficiency gain of MaskedKD takes place again. One interesting observation is that using a large teacher (DeiT3-L) for a relatively much smaller student (ViT-S) does not degrade the performance in this case, unlike in the typical case where we fine-tune all layers of the student model.

Table 9: Flexibility in data augmentation. Our method overcomes FastKD’s limitation, i.e., being restricted to simple data augmentations, resulting in an improved student model performance. “Simple” refers to applying basic augmentations: random resized crop and horizontal flip. “Hard” means additionally performing RandAugment, *mixup* and *cutmix*.

Student	Augmentation	Teacher	Method	Acc.	PFLOPs
DeiT-S	Simple	DeiT-B	-	No distillation	71.3 -
			Logit +MaskedKD 50%	79.7 80.2	6757 3349(-50%)
	Hard	DeiT-B	-	No distillation	79.9 -
			Logit +MaskedKD 50%	80.8 81.0	6757 3349(-50%)

E MaskedKD and data augmentations

One of the key advantages of the MaskedKD comparing with the FastKD [39] is that MaskedKD can be applied to distillation scenarios where we use heavy data augmentation schemes. As FastKD requires pre-computing and storing the teacher predictions for all augmented samples, computational benefits of the FastKD may be greatly undermined by considering heavier and more diverse augmentations. In this section, we perform a basic sanity check that (1) such heavy data augmentations are indeed useful in KD scenarios,⁶ and (2) MaskedKD works well with heavy augmentations.

Tab. 9 gives the experimental results. We find that, even for relatively small-scale student models such as DeiT-S, the data augmentation greatly boost the model performance. For undistilled students, the gain can be as large as 8.6%. For the models trained with basic logit distillation, the gain is 1.1%. We also observe that MaskedKD preserves the student model accuracy with both light and heavy data augmentations.

⁶ Under non-KD contexts, [41] makes a similar observation.

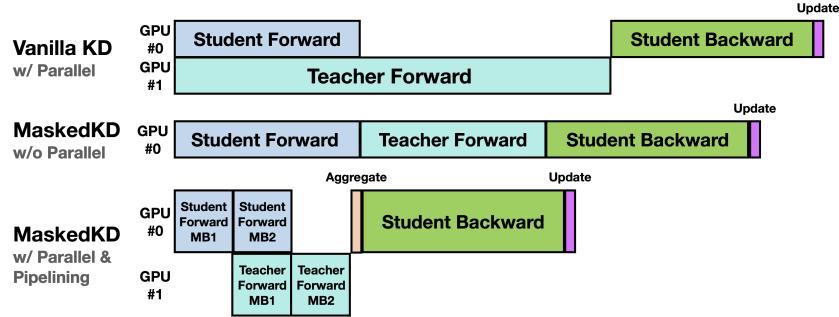


Fig. 8: Right : Training time breakdown of MaskedKD. (Top): In vanilla KD, the student and teacher forwards can be processed in parallel. (Middle): In MaskedKD, the teacher forward may wait until the student forward ends, and yet, the teacher forward can be reduced a lot. (Bottom): Pipeline parallelism allows MaskedKD to efficiently utilize multiple GPUs.

F Pipelining MaskedKD

To apply the MaskedKD, we need to compute the patch saliency before the teacher inference stage. As the saliency score is computed at the last layer of the student, the teacher forward cannot begin until the student forward is completed. Thus, a naïve parallelization strategy of running the teacher and the student model on two separate GPUs may be somewhat less effective than in the vanilla knowledge distillation (Fig. 8, top).

However, this *does not* imply that (1) there is no speedup, or (2) there is no effective parallelization strategy. First, we note that the teacher forward is usually the key bottleneck in knowledge distillation, often taking much longer than the student forward. MaskedKD dramatically reduces the time for teacher forward, so that the MaskedKD running the student and teacher forward in series can be faster than the teacher forward of the vanilla KD in some cases (Fig. 8, middle). Second, when using multiple GPUs, we can use pipelining to fill up the bubbles, as in GPipe [17]. More specifically, one can divide each training data batch into smaller mini-batches and make forward inferences on them sequentially. This division allows the teacher GPU to access the data before the student forward completes on all mini-batches (Fig. 8, bottom).

G Large Teacher with More Masking vs. Small Teacher with Less Masking

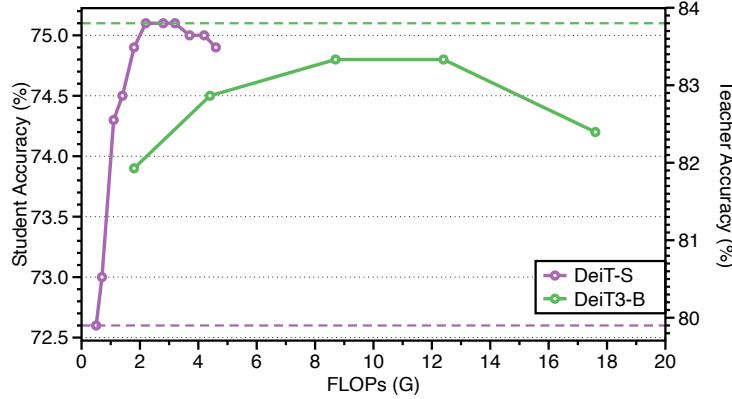


Fig. 9: Differences of the MaskedKD results by model size. We compare two teachers: DeiT-S and DeiT3-B. The colored dashed lines denote the accuracy of the teacher model. The performance of the teacher model is depicted by the dashed line, indicating its level of performance (Right). The line with dots represents the performance of the distilled student model performance (Left).

Masking the teacher input gives a new way to trade the student performance for the training efficiency, in addition to changing the model size. In this section, we compare the performance-efficiency tradeoff curve of two different-sized teacher models. In particular, we examine whether a larger teacher that uses less patch can give a more cost-efficient guidance than a smaller teacher that uses more patches. For this purpose, we compare the performance of the DeiT-Tiny students trained with DeiT-S and DeiT3-B teachers that use varying fraction of patches (Fig. 9). We observe that masking does not give a dramatic change to the answer to the question “which teacher is most cost-efficient?” DeiT-S teacher dominates DeiT3-B teacher at all per-iteration FLOPs level.

H Implementation

The following is the pseudo code of our “MaskedKD Engine” in PyTorch [31]:

```

def maskedkd_engine(image, labels, student, teacher, num_keep):
    """
    num_keep: the int number of patches to keep.
    """

    output, attn = student(image)

    len_keep = torch.topk(attn.mean(dim=1)[:, 0, 1:], num_keep).indices

    def teacher_inference(image, teacher, len_keep):
        x = teacher.patch_embed(image)
        B, _, D = x.shape  # batch, length, dim

        cls_tokens = teacher.cls_token.expand(B, -1, -1)
        x = torch.cat((cls_tokens, x), dim=1)
        x = x + teacher.pos_embed

        cls_save = x[:, 0, :].unsqueeze(dim=1)
        x = x[:, 1:, :]
        index = len_keep.unsqueeze(-1).repeat(1, 1, D)
        x = torch.gather(x, dim=1, index=index)
        x = torch.cat((cls_save, x), dim=1)

        for blk in teacher.blocks:
            x = blk(x)

        x = teacher.norm(x)
        output = teacher.head(x)
        return output

    t_output = teacher_inference(image, teacher, len_keep)
    loss = CE(output, labels) + KL_DIV(output, t_output)

    return loss

```

This returns the loss. By simply adding the patch selection stage, we can greatly improve efficiency and achieve substantial gains.

I Saving Costs

Table 10: Analyzing Computational Cost. We examine how the computational cost varies when applying MaskedKD. The arrow symbol represents the difference in FLOPs when utilizing MaskedKD.

Layer	Complexity	Computation (GFLOPs)		
		DeiT-S @ 384	DeiT-B	DeiT-S
Softmax-Attention	$\mathcal{O}(LN^2M)$	3.1 → 0.8	0.7 → 0.2	0.4 → 0.1
Projections	$\mathcal{O}(LNM^2)$	4.1 → 2.0	5.6 → 2.8	1.4 → 0.7
MLP	$\mathcal{O}(LNM^2)$	8.2 → 4.1	11.2 → 5.6	2.8 → 1.4
Total	$\mathcal{O}(LNM(M + N))$	15.3 → 6.9	17.5 → 8.6	4.5 → 2.2

We analyze how much cost, in terms of FLOPs and memory, could be saved by applying MaskedKD. ViT is encoder-only transformer [47], which is mainly consisted of a multi-head self-attention layers and a multi-layer perceptron layers. There are also many details which takes very tiny portion in terms of calculation, such as the embedding layer, residual connection, bias, GeLU, or layer normalization and we will ignore it in this section. We denote $\phi(n, d)$ as a function of FLOPs with respect to the number of tokens n and the embedding dimension d . For example, in case of DeiT-B, n is 197 and d is 768. For self-attention layer, the FLOPs mainly comes from two parts: (1) The projection of Q, K, V matrices and the self-attention outputs $\phi_{\text{proj}}(n, d) = 4nd^2$, (2) The calculation of the softmax-attention $\phi_{\text{SA}}(n, d) = 2n^2d$.

The FLOPs of MLP layers comes from two fully-connected (FC) layers. Two FC layers have a difference of four times in dimension. Therefore, the FLOPs for MLP layer is $\phi_{\text{FC}}(n, d) = 8nd^2$.

By combining the self-attention layer and the MLP layer, we can get the total FLOPs of one ViT block.

$$\phi_{\text{BLK}}(n, d) = \phi_{\text{proj}}(n, d) + \phi_{\text{FC}}(n, d) + \phi_{\text{SA}}(n, d) = 4nd^2 + 2n^2d + 8nd^2 = 12nd^2. \quad (3)$$

Since there is 12 layers in case of DeiT-Base, the total FLOPs is

$$12 * \phi_{\text{BLK}}(n, d) = 12 * (12nd^2 + 2n^2d) = 144nd^2 + 24n^2d. \quad (4)$$

In Table 11, we compare the RAM used by Logit and MaskedKD, when distilling DeiT-B to DeiT-S with batch size 128. Masking reduces the amount of intermediate activations computed during teacher forward; the reduction ratio can change by using different batch size or memory offloading.

Table 11: Memory Usage Comparison.

Method	Student	Teacher	Total
Logit	6426MB	6499MB	12925MB
+MaksedKD _{50%}	6426MB	3412MB	9838MB (-24%)

J Applicable to detection or segmentation

In Table 12, we present the application of the MaskedKD to object detection and segmentation tasks. For detection, we distill YOLOS-B [9] to YOLOS-Ti. The model is initially pretrained on the COCO dataset and then finetuned on the PASCAL VOC dataset, utilizing algorithms from ViDT [40]. For segmentation, we use the manifold distillation [13] to distill from Segmente-S [42] to Segmente-Ti on ADE20k. By applying MaskedKD, we save training costs through masking during distillation without sacrificing performance in detection and segmentation tasks as well.

Table 12: MaskedKD to Detection and Segmentation Tasks.

Method	Task: Detection	Task: Segmentation
	AP @ VOC val	mIoU @ ADE20k
Baseline	45.4	38.7
+MaskedKD	45.9 (50% masked)	38.8 (25% masked)

K Theoretical analyses

We provide theoretical analyses to explain the results in Figure 5, offering the following proposition to explain why randomly initialized attentions can generate meaningful masks:

Proposition 1. *Let $\mathbf{c} \in \mathbb{R}^d$, $W_q, W_k \in \mathbb{R}^{d \times d}$ be the class token and the query/key weight matrices whose entries are i.i.d. initialized as $\mathcal{N}(0, 1/d)$. Let $f(\cdot)$ be the pre-softmax attention of the class token to another token, i.e., $f(\mathbf{x}) := (W_q \mathbf{c}) \cdot (W_k \mathbf{x}) / \sqrt{d}$. Then, for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, we have*

$$\mathbb{E} \|f(\mathbf{x}) - f(\mathbf{y})\|_2 \leq C_0 \cdot (\log \log d / \sqrt{d}) \cdot \|\mathbf{x} - \mathbf{y}\|_2,$$

where C_0 is a constant independent of d .

Proof. We proceed as follows.

$$\begin{aligned} \|f(\mathbf{x}) - f(\mathbf{y})\|_2 &= \frac{1}{\sqrt{d}} \|(\mathbf{x} - \mathbf{y})^\top W_k^\top W_q \mathbf{c}\|_2 \\ &\leq \frac{1}{\sqrt{d}} \|\mathbf{x} - \mathbf{y}\|_2 \cdot \|W_k\| \cdot \|W_q\| \cdot \|\mathbf{c}\|_2, \end{aligned}$$

where $\|\cdot\|$ denotes the spectral norm. Taking expectations to both sides, we get

$$\mathbb{E}\|f(\mathbf{x}) - f(\mathbf{y})\|_2 \leq \frac{1}{\sqrt{d}}\|\mathbf{x} - \mathbf{y}\|_2 \cdot \underbrace{(\mathbb{E}\|W_k\|)^2}_{:=T_1} \cdot \underbrace{\mathbb{E}\|\mathbf{c}\|_2}_{:=T_2}$$

where we have used the fact that W_q, W_k, \mathbf{c} are independent, and W_q, W_k have identical distributions. From the standard results on the spectral norm of random Gaussian matrices [36, Theorem 2], we know that $T_1 \asymp \log \log d$. Invoking Jensen's inequality, we know that $T_2 \leq \sqrt{\mathbb{E}\|\mathbf{c}\|_2^2} = 1$. Thus, we get what we want.