

---

# Rethinking Patch Dependence for Masked Autoencoders

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 In this work, we examine the impact of inter-patch dependencies in the decoder of  
2 masked autoencoders (MAE) on representation learning. We decompose the decod-  
3 ing mechanism for masked reconstruction into self-attention between mask tokens  
4 and cross-attention between masked and visible tokens. Our findings reveal that  
5 MAE reconstructs coherent images from visible patches not through interactions  
6 between patches in the decoder but by learning a global representation within the  
7 encoder. This discovery leads us to propose a simple visual pretraining framework:  
8 cross-attention masked autoencoders (CrossMAE). This framework employs only  
9 cross-attention in the decoder to independently read out reconstructions for a small  
10 subset of masked patches from encoder outputs, yet it achieves comparable or  
11 superior performance to traditional MAE across models ranging from ViT-S to  
12 ViT-H. By its design, CrossMAE challenges the necessity of interaction between  
13 mask tokens for effective masked pretraining. Code is available here.

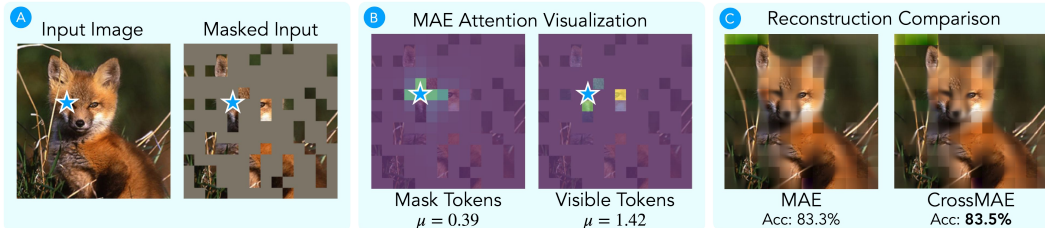
## 14 1 Introduction

15 Masked image modeling [46, 30, 61, 4] has emerged as a pivotal unsupervised learning technique  
16 in computer vision. One such recent work following this paradigm is masked autoencoders (MAE):  
17 given only a small, random subset of visible image patches, the model is tasked to reconstruct the  
18 missing pixels. By operating mostly on this small subset of visible tokens, MAE can efficiently  
19 pre-train high-capacity models on large-scale vision datasets, demonstrating impressive results on a  
20 wide array of downstream tasks [33, 38, 49].

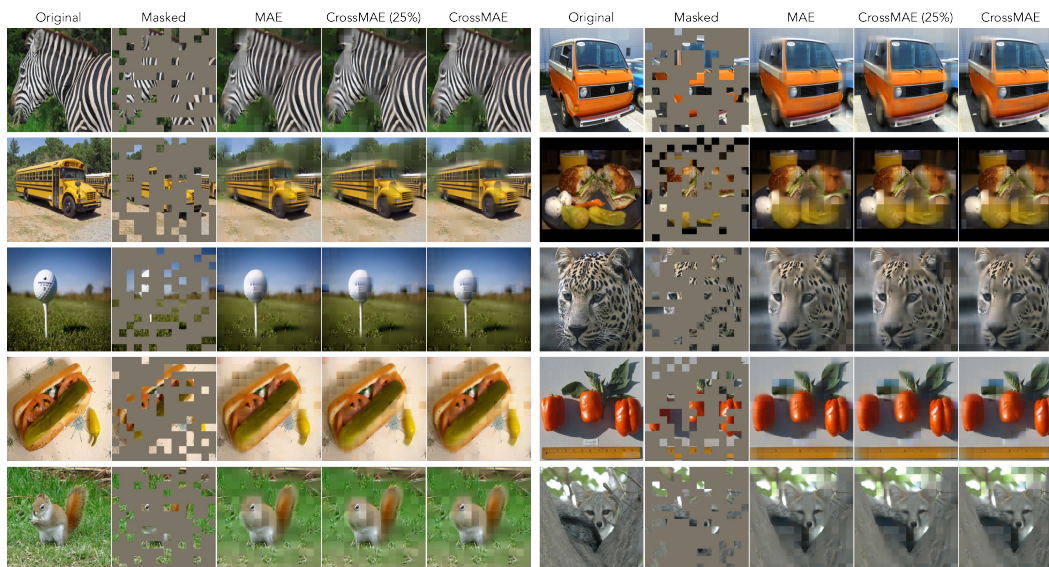
21 The MAE framework employs *self-attention* across the entire model for self-supervised reconstruction  
22 tasks. In this setup, both masked and visible tokens engage in self-attention, not just with each other  
23 but also with themselves, aiming to generate a holistic and context-aware representation. However,  
24 the masked tokens inherently lack information. Intuitively, facilitating information exchange among  
25 adjacent masked tokens should enable the model to synthesize a more coherent image, thereby  
26 accomplishing the task of masked reconstruction and improving representation learning. A question  
27 arises, though: Is this truly the case?

28 We decompose the decoding process of each mask token into two parallel components: self-attention  
29 with other mask tokens, as well as cross-attention to the encoded visible tokens. If MAE relies on  
30 the self-attention with other mask tokens, its average should be on par with the cross-attention. Yet,  
31 the quantitative comparison in Figure 1.(b) shows the magnitude of mask token-to-visible token  
32 cross-attention (1.42) in the MAE decoder evaluated over the entire ImageNet validation set far  
33 exceeds that of mask token-to-mask token self-attention (0.39).

34 This initial observation prompts two questions: **1)** Is the self-attention mechanism among mask  
35 tokens in the decoder necessary for effective representation learning? **2)** If not, can each patch be



**Figure 1: Method Overview.** (A) Masked autoencoder (MAE) starts by masking random patches of the input image. (B) To reconstruct a mask token (marked by the blue star), MAE attends to both the masked tokens (B.Left) and the visible tokens (B.Right). A quantitative comparison over the ImageNet validation set shows that the masked tokens in MAE disproportionately attend to the visible tokens (1.42 vs 0.39), questioning the necessity of attention within mask tokens. (C) We propose CrossMAE, the masked patches are reconstructed from only the cross attention between the masked tokens and the visible tokens. Surprisingly, CrossMAE attains the same or better performance than MAE on ImageNet classification and COCO instance segmentation.



**Figure 2:** Example reconstructions of ImageNet *validation* images. For each set of 5 images, from left to right, are the original image, masked image with a mask ratio of 75%, MAE [30], CrossMAE (trained to reconstruct 25% of image tokens, or 1/3 of the mask tokens), and CrossMAE (trained to reconstruct all masked tokens). Since CrossMAE does not reconstruct them, all model outputs have the visible patches overlaid. Intriguingly, CrossMAE, when trained for partial reconstruction, can decode all mask tokens in one forward pass (shown above), indicating that the encoder rather than the decoder effectively captures global image information in its output tokens. Its comparable reconstruction quality to full-image-trained models suggests that full-image reconstruction might not be essential for effective representation learning.

36 *independently* read out from the encoder output, allowing the reconstruction of only a small subset of  
 37 masked patches, which in turn, accelerates the pretraining without performance degradation?

38 In addressing these questions, we introduce CrossMAE, which diverges from MAE in three ways:

39 **1. Cross-attention for decoding.** Rather than passing a concatenation of mask and visible  
 40 tokens to a *self-attention* decoder, CrossMAE uses mask tokens as queries to read out the masked  
 41 reconstructions from the visible tokens in a *cross-attention decoder*. In this setting, mask tokens  
 42 incorporate information from the visible tokens but do not interact with other mask tokens, thereby  
 43 reducing the sequence length for the decoder and cutting down computational costs.

44 **2. Independent partial reconstruction.** With self-attention removed, the decoding of each mask  
 45 token, based on the encoded features from visible tokens, becomes conditionally independent. This  
 46 enables the decoding of only a fraction of masked tokens rather than the entire image.

47 **3. Inter-block attention.** Due to the separation of visible and mask tokens, we can use features  
 48 from different encoder blocks for each decoder block. Empirically, we find solely relying on the last

49 encoder feature map for reconstruction, the design present in MAE, hurts feature learning. We propose  
50 a lightweight inter-block attention mechanism that allows the CrossMAE decoder to leverage a mix  
51 of low-level and high-level feature maps from the encoder, improving the learned representation.

52 The analysis performed on CrossMAE led to a novel way to understand MAE. Even though the  
53 patches to be reconstructed are independently decoded, our findings demonstrate that *coherent*  
54 reconstruction for each masked patch can be independently read out from the encoder output, without  
55 any interactions among masked tokens in the decoder for consistency (Figure 2). Furthermore, the  
56 downstream performance of the model remains robust even without these interactions (Figure 1.(c),  
57 Tables 1 and 2). Both pieces of evidence confirm that the encoder’s output features already encapsulate  
58 the necessary global context for image reconstruction, while the decoder simply performs a readout  
59 from the encoder output to reconstruct the pixels at the location of each patch.

60 **To sum up, our main contributions are the following:**

61 **1. We present a novel understanding of MAE.** Our findings show that MAE reconstructs coherent  
62 images from visible patches *not through interactions between patches to be reconstructed* in the  
63 decoder but by *learning a global representation within the encoder*. This is evidenced by the model’s  
64 ability to generate coherent images and maintain robust downstream performance without such  
65 interactions, indicating the encoder effectively captures global image information.

66 **2. We advocate replacing self-attention layers with a simple cross-attention readout function.**  
67 Given our discovery that the encoder in MAE already captures a comprehensive global representation,  
68 we propose replacing self-attention layers in the decoder with a more efficient information readout  
69 function. Specifically, we suggest utilizing *cross-attention* to aggregate the output tokens of the  
70 encoder into each input token within the decoder layers *independently*, thereby eliminating the need  
71 for token-to-token communication within the decoder.

72 **3. CrossMAE achieves comparable or superior performance with reduced computational**  
73 **costs** in image classification and instance segmentation compared to MAE on vision transformer  
74 models *ranging from ViT-S to ViT-H*. Code is available here.

## 75 **2 Related Works**

### 76 **2.1 Self-Supervised Learning**

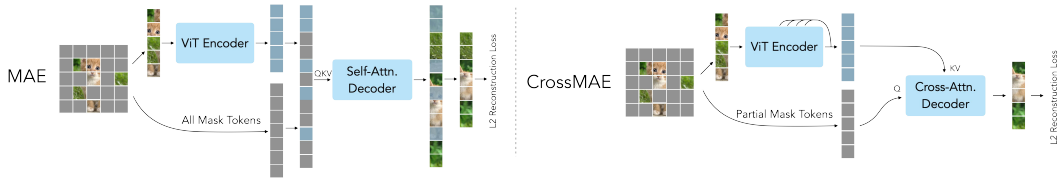
77 In self-supervised representation learning, a model trains on a pretext task where the supervision  
78 comes from the input data itself without labels. Contrastive learning methods learn representations  
79 by contrasting positive and negative samples, such as SimCLR [11], CPC [44], MoCo [29, 12, 13],  
80 CLD [59] and SwAV [7]. Additionally, in BYOL [26], iBOT [65], DINO [8], DINOv2 [45], and  
81 MaskAlign [62] make a student model to imitate a teacher model without negative pairs.

82 Generative modeling, focusing on acquiring a generative model capable of capturing the underlying  
83 data distribution, is an alternative method for self-supervised learning. VAE/GAN [35] merges the  
84 strengths of variational autoencoders and generative adversarial networks to acquire disentangled  
85 representations of data. PixelCNN, PixelVAE, and PixelTransformer [55, 27, 54] generate images  
86 pixel by pixel, taking into account the context of previously generated pixels. Masked modeling, a  
87 large subclass of generative modeling, is discussed in the following subsection. After the pre-training  
88 stage, these generative models can be finetuned for many downstream applications.

### 89 **2.2 Masked Modeling**

90 Masked modeling learns representations by reconstructing a masked portion of the input. Pioneering  
91 works in natural language processing (NLP) present various such pretraining objectives. BERT [19]  
92 and its extensions [41, 34] use a bidirectional transformer and present few-shot learning capabil-  
93 ities from masked language modeling. GPT [47, 48, 5], uses autoregressive, causal masking and  
94 demonstrates multi-task, few-shot, and in-context learning capabilities.

95 Early works in computer vision, such as Stacked Denoising Autoencoders [57] and Context En-  
96 coder [46], investigated masked image modeling as a form of denoising or representation learning.  
97 Recently, with the widespread use of transformer [20] as a backbone vision architecture, where  
98 images are patchified and tokenized as sequences, researchers are interested in how to transfer the  
99 success in language sequence modeling to scale vision transformers. BEiT [3], MAE [30], and Sim-



**Figure 3:** MAE [30] concatenates *all* mask tokens with the visible patch features from a ViT encoder and passes them to a decoder with self-attention blocks to reconstruct the original image. Patches that correspond to visible tokens are then dropped, and an L2 loss is applied to the rest of the reconstruction as the pretraining objective. CrossMAE instead uses cross-attention blocks in the decoder to reconstruct only a subset of the masked tokens.

100 MIM [61] are a few of the early works that explored BERT-style pretraining of vision transformers.  
 101 Compared to works in NLP, both MAE and SimMIM [30, 61] find that a much higher mask ratio  
 102 compared to works in NLP is necessary to learn good visual representation. Many recent works  
 103 further extend masked pretraining to hierarchical architectures [61, 40] and study data the role of data  
 104 augmentation [9, 21]. Many subsequent works present similar successes of masked pretraining for  
 105 video [52, 58, 22, 28], language-vision and multi-modal pretraining [1, 39, 23] and for learning both  
 106 good representations and reconstruction capabilities [60, 37].

107 However, BERT-style pretraining requires heavy use of self-attention, which makes computational  
 108 complexity scale as a polynomial of sequence length. PixelTransformer [54] and DiffMAE [60] both  
 109 use cross-attention for masked image generation and representation learning. Siamese MAE [28]  
 110 uses an asymmetric masking pattern and decodes frames of a video condition on an earlier frame. In  
 111 these settings, *all* masked patches are reconstructed. In this work, we investigate if learning good  
 112 features necessitates high reconstruction quality and if the entire image needs to be reconstructed to  
 113 facilitate representation learning. PCAE [36] progressively discards redundant mask tokens through  
 114 its network, leading to a few tokens for reconstruction. VideoMAEv2 [58] concatenates randomly  
 115 sampled masked tokens with visible tokens and uses self-attention to reconstruct the masked patches.  
 116 In comparison, we minimally modify MAE with a cross-attention-only decoder and masked tokens  
 117 are decoded in a conditional independent way.

### 118 2.3 Applications of Cross-Attention

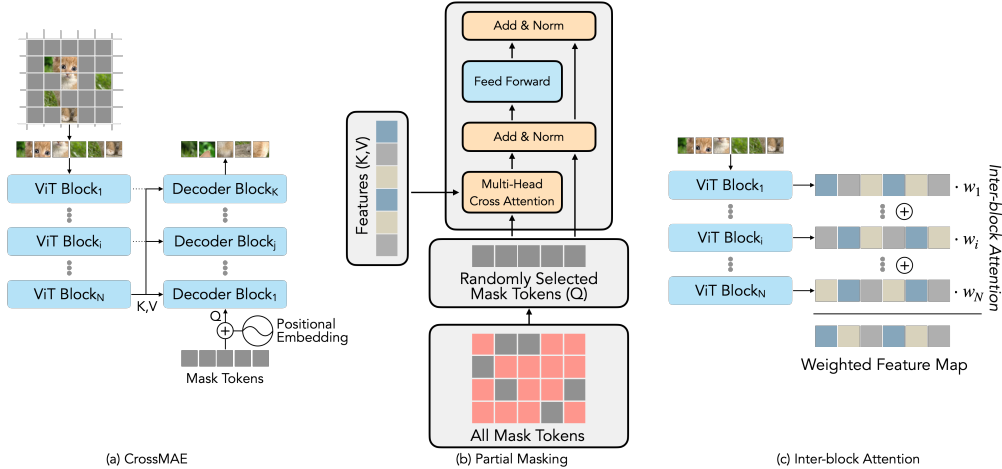
119 In addition to the prevalent use of self-attention in computer vision, cross-attention has shown to be a  
 120 cost-effective way to perform pooling from a large set of visible tokens. Intuitively, cross-attention  
 121 can be seen as a parametric form of pooling, which learnably weighs different features. Touvron  
 122 et al. [53] replace mean pooling with cross-attention pooling and find improvement in ImageNet  
 123 classification performance. Jaegle et al. [32] uses cross-attention to efficiently process large volumes  
 124 of multi-modal data. Cross-attention is also widely used for object detection. Carion et al. [6] utilizes  
 125 query tokens as placeholders for potential objects in the scene. Cheng et al. [16, 15] further extend  
 126 this concept by introducing additional query tokens to specifically tackle object segmentation in  
 127 addition to the query tokens for object detection. Distinct from the prior works, we are interested in  
 128 the role of cross-attention for representation learning in a self-supervised manner.

## 129 3 CrossMAE

130 We start with an overview of vanilla masked autoencoders in Section 3.1. Next, in Section 3.2, we  
 131 introduce the use of cross-attention in place of self-attention in the decoder for testing the necessity  
 132 of interaction between mask tokens for representation learning. In Section 3.3, we discuss how  
 133 eliminating self-attention in the decoding process enables us to reconstruct only a subset of masked  
 134 tokens, leading to faster pretraining. Finally, Section 3.4 presents our inter-block attention mechanism,  
 135 which allows decoder blocks to leverage varied encoder features.

### 136 3.1 Preliminaries: Masked Autoencoders

137 Masked Autoencoders (MAE) [30] pretrain Vision Transformers (ViTs) [20]. Each image input is  
 138 first patchified, and then a random subset of the patches is selected as the visible patches. As depicted  
 139 in Figure 3, the visible patches, concatenated with a learnable class token [CLS], are subsequently



**Figure 4: Overview of CrossMAE.** (a) The vanilla version of CrossMAE uses the output of the last encoder block as the keys and queries for cross-attention. The first decoder block takes the sum of mask tokens and their corresponding positional embeddings as queries, and subsequent layers use the output of the previous decoder block as queries to reconstruct the masked patches. (b) Unlike the decoder block in [56], the cross-attention decoder block does not contain self-attention, decoupling the generation of different masked patches. (c) CrossMAE’s decoder blocks can leverage low-level features for reconstruction via inter-block attention. It weighs the intermediate feature maps, and the weighted sum of feature maps is used as the key and value for each decoder block.

140 fed into the ViT encoder, which outputs a set of feature latents. The latent vectors, concatenated with  
 141 the sum of the positional embeddings of the masked patches and the learnable mask token, are passed  
 142 into the MAE decoder. The decoder blocks share the same architecture as the encoder blocks (i.e.,  
 143 both are transformer blocks with self-attention layers). Note that the number of tokens fed into the  
 144 decoder is the *same* length as the original input, and the decoding process assumes that the decoded  
 145 tokens depend on both visible and masked tokens. Decoder outputs pass through a fully connected  
 146 layer per patch for image reconstruction. After the reconstruction is generated, the loss is applied  
 147 only to the masked positions, while the reconstructions for visible spatial locations are discarded.

148 Recall in Sec. 1 we measure the mean attention value across all attention maps over the ImageNet  
 149 validation set to study the properties of MAE. We grouped the attention values by cross-attention  
 150 and self-attention between visible and masked tokens. We observed that in the decoding process  
 151 of an MAE, mask tokens attend disproportionately to the class token and the visible tokens (see  
 152 Figure 1.(b)). This motivates us to make design decisions and conduct experiments specifically to  
 153 answer the following question: *Can we simplify the decoding process by eliminating self-attention*  
 154 *among masked tokens without compromising the model’s ability to generate coherent images and*  
 155 *perform well on downstream tasks?*

### 156 3.2 Reconstruction with Cross-Attention

157 To address this question, we substitute the self-attention mechanism in the decoder blocks with  
 158 cross-attention, using it as a readout function to decode the latent embedding from the encoder to raw  
 159 pixel values. Specifically, the decoder employs multi-head cross-attention where the queries are the  
 160 output from previous decoder blocks (or the sum of position embedding of the masked patches and  
 161 mask token for the first decoder block). The keys and values are from the encoded features.

162 In the most basic CrossMAE, the output from the final encoder block is used as the key and value  
 163 tokens for all layers of the decoder, as illustrated in Fig. 4(a). Further exploration in Sec.3.4 reveals  
 164 that utilizing a weighted mean of selected encoder feature maps can be beneficial. The residual  
 165 connections in each decoder block enable iterative refinement of decoded tokens as they progress  
 166 through decoder blocks.

167 Diverging from the original transformer architecture [56], our decoder omits the causal self-attention  
 168 layer before the introduction of multi-head cross-attention. This elimination, coupled with the fact  
 169 that layer normalization and residual connections are only applied along the feature axis but not

170 the token axis, enables the independent decoding of tokens. This design choice is evaluated in the  
171 ablation study section to determine its impact on performance.

172 Given the disparity in the dimensions of the encoder and decoder, MAE adapts the visible features to  
173 the decoder’s latent space using an MLP. However, in CrossMAE, as encoder features are integrated  
174 at various decoder blocks, we embed the projection within the multi-head cross-attention module.

175 Cross-attention layers serve as a readout function that decodes the global representation provided  
176 in the encoder’s output tokens to the pixel values within each patch to be reconstructed. However,  
177 CrossMAE does not restrict the architecture to a single cross-attention block. Instead, we stack  
178 multiple cross-attention decoder blocks in a manner more akin to the traditional transformer [56].

### 179 3.3 Partial Reconstruction

180 The fact that CrossMAE uses cross-attention rather than self-attention in the decoder blocks brings  
181 an additional benefit over the original MAE architecture. Recall that mask tokens are decoded inde-  
182 pendently and thus there is no exchange of information between them, to obtain the reconstructions  
183 at a specific spatial location, CrossMAE only needs to pass the corresponding mask tokens to the  
184 cross-attention decoder. This allows partial reconstruction in contrast to the original full-image  
185 reconstruction in the MAE architecture which needs to pass all the masked tokens as the input of the  
186 decoder blocks due to the existence of self-attention in the decoder blocks.

187 To address the second question in Sec. 3.1, rather than decoding the reconstruction for all masked  
188 locations, we only compute the reconstruction on a random subset of the locations and apply the loss  
189 to the decoded locations. Specifically, we name the ratio of predicted tokens to all image tokens as  
190 *prediction ratio* ( $\gamma$ ), and the mask ratio ( $p$ ). Then the prediction ratio is bounded between  $\gamma \in (0, p]$ .  
191 Because we are sampling within the masked tokens uniformly at random and the reconstruction  
192 loss is a mean square error on the reconstructed patches, the expected loss is the same as in MAE,  
193 while the variance is  $(p/\gamma)$  times larger than the variance in MAE. Empirically, we find that scaling  
194 the learning rate of MAE ( $\beta$ ) to match the variance (i.e. setting the learning rate as  $\gamma\beta/p$ ) helps  
195 with model performance. Since cross-attention has linear complexity with respect to the number of  
196 masked tokens, this partial reconstruction paradigm decreases computation complexity. Empirically,  
197 we find that the quality of the learned representations is not compromised by this approach.

### 198 3.4 Inter-block Attention

199 MAE combines the feature of the last encoder block with mask tokens as the input to the self-attention  
200 decoder, which creates an information bottleneck by making early encoder features inaccessible  
201 for the decoder. In contrast, CrossMAE’s cross-attention decoder decouples queries from keys and  
202 values. This decoupling allows different cross-attention decoder blocks to take in feature maps from  
203 different encoder blocks. This added degree of flexibility comes with a design choice for selecting  
204 encoder features for each decoder block. One naive choice is to give the feature of the  $i$ th encoder  
205 block to the last  $i$ th decoder (e.g., feeding the feature of the first encoder to the last decoder), in a  
206 U-Net-like fashion. However, this assumes the decoder’s depth matches the depth of the encoder,  
207 which is not the case for MAE or CrossMAE.

208 Instead of manually matching each decoder block with an encoder feature map, we make the selection  
209 *learnable* and propose inter-block attention for feature fusion for each decoder block (Figure 4(c)).  
210 Analogous to the inter-patch cross-attention that takes a weighted sum of the visible token embeddings  
211 across the patch dimensions to update the embeddings of masked tokens, inter-block attention takes  
212 a weighted sum of the visible token embeddings *across different input blocks* at the same spatial  
213 location to fuse the input features from multiple blocks into one feature map for each decoder block.

214 Concretely, each decoder block takes a weighted linear combination of encoder feature maps  $\{f_i\}$  as  
215 keys and values. Specifically, for each key/value token  $t_k$  in decoder block  $k$  in a model with encoder  
216 depth  $n$ , we initialize a weight  $w^k \in \mathcal{R}^n \sim \mathcal{N}(0, 1/n)$ . Then  $t_k$  is defined as

$$t_k = \sum_{j=1}^n w_j^k f_j. \quad (1)$$

217 In addition to feature maps from different encoder blocks, we also include the inputs to the first  
218 encoder block to allow the decoder to leverage more low-level information to reconstruct the original  
219

Method	ViT-S	ViT-B	ViT-L	ViT-H
Supervised [50]	79.0	82.3	82.6	83.1
DINO [8]	-	82.8	-	-
MoCo v3 [14]	<u>81.4</u>	83.2	84.1	-
BEiT [3]	-	83.2	<u>85.2</u>	-
MultiMAE [2]	-	83.3	-	-
MixedAE [9]	-	<u>83.5</u>	-	-
CIM [21]	<b>81.6</b>	83.3	-	-
MAE [30]	78.9	83.3	<b>85.4</b>	<u>85.8</u>
CrossMAE (25%)	79.2	<b>83.5</b>	<b>85.4</b>	<b>86.3</b>
CrossMAE (75%)	79.3	<b>83.7</b>	<b>85.4</b>	-

**Table 1: ImageNet-1K classification accuracy.** CrossMAE performs on par or better than MAE. All experiments are run with 800 epochs. The best results are in **bold** while the second best results are underlined.

Method	AP <sup>box</sup>		AP <sup>mask</sup>	
	ViT-B	ViT-L	ViT-B	ViT-L
Supervised [38]	47.6	49.6	42.4	43.8
MoCo v3 [14]	47.9	49.3	42.7	44.0
BEiT [3]	49.8	53.3	44.4	47.1
MixedAE [9]	50.3	-	43.5	-
MAE [38]	51.2	54.6	45.5	48.6
CrossMAE	<b>52.1</b>	<b>54.9</b>	<b>46.3</b>	<b>48.8</b>

**Table 2: COCO instance segmentation.** Compared to previous masked visual pretraining works, CrossMAE performs favorably on object detection and instance segmentation tasks.

220 image. We can select a subset of the feature maps from the encoder layers instead of all feature maps.  
 221 This reduces the computation complexity of the system. We ablate this in Table 3d.

222 We show that using the weighted features rather than simply using the features from the last block  
 223 greatly improves the performance of CrossMAE. Intriguingly, in the process of learning to achieve  
 224 better reconstructions, early decoder blocks tend to prioritize information from later encoder blocks,  
 225 while later decoder blocks focus on earlier encoder block information, as demonstrated in Section 4.5.

## 226 4 Experiments

227 We perform self-supervised pretraining on ImageNet-1K, following MAE [30]’s hyperparameter  
 228 settings, only modifying the learning rate and decoder depth. The hyperparameters were initially  
 229 determined on ViT-Base and then directly applied to ViT-Small, ViT-Large, and ViT-Huge. Both  
 230 CrossMAE and MAE are trained for 800 epochs. We provide implementation details and more  
 231 experiments in the appendix.

### 232 4.1 ImageNet Classification

233 **Setup.** The model performance is evaluated with end-to-end fine-tuning, with top-1 accuracy used  
 234 for comparison. Same as in Figure. 2, we compare two versions of CrossMAE: one with a prediction  
 235 ratio of 25% (1/3 of the mask tokens) and another with 75% (all mask tokens). Both models are  
 236 trained with a mask ratio of 75% and a decoder depth of 12.

237 **Results.** As shown in Table 1, CrossMAE outperforms vanilla MAE using the same ViT-B encoder  
 238 in terms of fine-tuning accuracy. This shows that replacing the self-attention with cross-attention  
 239 *does not degrade* the downstream classification performance of the pre-trained model. Moreover,  
 240 CrossMAE outperforms other self-supervised and masked image modeling baselines, *e.g.*, DINO [8],  
 241 MoCo v3 [14], BEiT [3], and MultiMAE [2].

### 242 4.2 Object Detection and Instance Segmentation

243 **Setup.** We additionally evaluate models pretrained with CrossMAE for object detection and instance  
 244 segmentation, which require deeper spatial understanding than ImageNet classification. Specifically,  
 245 we follow ViTDet [38], a method that leverages a Vision Transformer backbone for object detection  
 246 and instance segmentation. We report box AP for object detection and mask AP for instance  
 247 segmentation, following MAE [30]. We compare against supervised pre-training, MoCo-v3 [14],  
 248 BEiT [4], and MAE [30].

249 **Results.** As listed in Table 2, CrossMAE, with the default 75% prediction ratio, performs better  
 250 compared to these baselines, including vanilla MAE. This suggests that similar to MAE, CrossMAE  
 251 performance on ImageNet positively correlates with instance segmentation. Additionally, Cross-  
 252 MAE’s downstream performance scales similarly to MAE as the model capacity increases from ViT-B  
 253 to ViT-L. This observation also supports our hypothesis that partial reconstruction is suprisingly  
 254 sufficient for learning dense visual representation.

Method	Acc. (%)	Mask Ratio	Acc. (%)	Pred. Ratio	Acc. (%)
MAE	83.0	65%	<b>83.5</b>	15%	83.1
CrossMAE	<b>83.3</b>	75%	<b>83.3</b>	25%	83.2
CrossMAE + Self-Attn	83.3	85%	83.3	75%	<b>83.3</b>

(a) **Attention type** in decoder blocks. Adding back self-attention between mask tokens does not improve performance.

(b) **Mask ratio.** CrossMAE has consistent performance across high mask ratios.

(c) **Prediction ratio.** CrossMAE performs well even when only a fraction of mask tokens are reconstructed.

# Feature Maps Fused	Acc. (%)	Decoder Depth	Acc. (%)	Image Resolution	Acc. (%)
1	82.9	1	83.0	224	<b>83.2</b>
3	83.3	4	83.1	448	<b>84.6</b>
6	<b>83.5</b>	8	83.1		
12	<b>83.3</b>	12	<b>83.3</b>		

(d) **Inter-block attention.** A combination of six select encoder feature maps is best.

(e) **Decoder depth.** CrossMAE performance scales with decoder depth.

(f) **Input resolution.** CrossMAE scales to longer input sequences.

**Table 3: Ablations on CrossMAE.** We report fine-tuning performance on ImageNet-1K classification with 400 epochs (*i.e.*, half of the full experiments) with ViT-B/16. MAE performance is reproduced using the official MAE code. Underline indicates the default setting for CrossMAE. **Bold** indicates the best hyperparameter among the tested ones. 1 feature map fused (row 1, Table 3(d)) indicates using only the feature from the last encoder block. We use 25% prediction ratio for both settings in Table 3(f) to accelerate training.

### 255 4.3 Ablations

256 **Cross-Attention vs Self-Attention.** As shown in Table 3a, CrossMAE, with its cross-attention-  
 257 only decoder, outperforms vanilla MAE in downstream tasks as noted in Section 4.1. Additionally,  
 258 combining cross-attention with self-attention does not enhance fine-tuning performance, indicating  
 259 that cross-attention alone is adequate for effective representation learning.

260 **Mask Ratio and Prediction Ratio.** In our experiments with different mask and prediction ratios (*i.e.*,  
 261 the ratio of mask tokens to all tokens and the ratio of reconstructed tokens to all tokens, respectively)  
 262 (see Table 3b and Table 3c), we found that our method’s performance is not significantly affected by  
 263 variations in the number of masked tokens. Notably, CrossMAE effectively learns representations  
 264 by reconstructing as few as 15% of tokens, compared to the 100% required by vanilla MAE, with  
 265 minimal impact on downstream fine-tuning performance, which shows that partial reconstruction is  
 266 sufficient for effective representation learning.

267 **Inter-block Attention.** Our ablation study, detailed in Table 3d, explored the impact of varying the  
 268 number of encoder feature maps in our inter-block attention mechanism. We found that using only  
 269 the last feature map slightly lowers performance compared to using all 12. However, even a partial  
 270 selection of feature maps improves CrossMAE’s performance, with the best results obtained using 6  
 271 feature maps. This indicates that CrossMAE does not require all features for optimal performance.

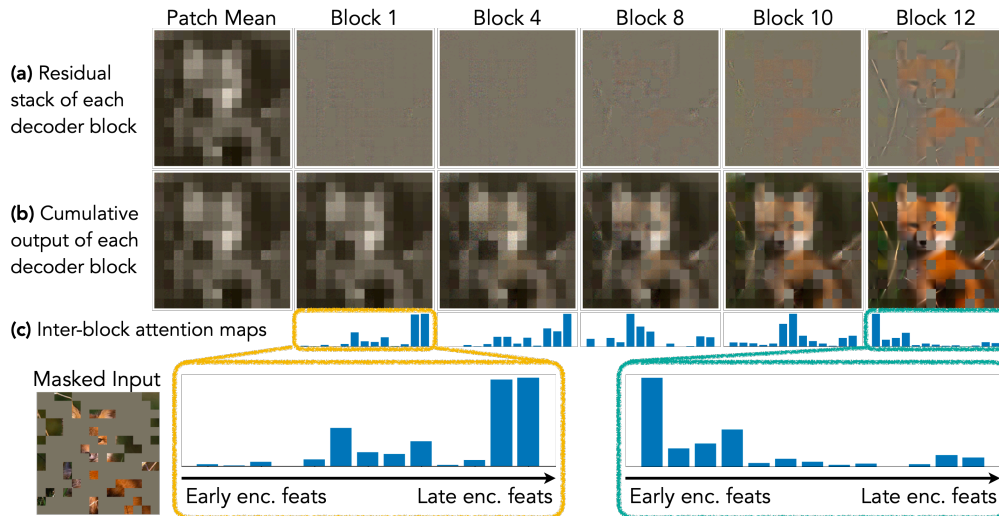
272 **Decoder Depth.** Table 3e shows that a 12-block decoder slightly improves performance compared  
 273 to shallower ones. Remarkably, CrossMAE achieves similar results to MAE with just one decoder  
 274 block, demonstrating its efficiency. Our experiments in Figure 7 that models with lower prediction  
 275 ratios benefit more from deeper decoders.

276 **Input Resolution.** We extend CrossMAE to longer token lengths by increasing the image resolution  
 277 with constant patch size. Escalating the resolution from 224 to 448 increases the token length from  
 278 197 to 785, challenging the scalability of current approaches. Thus, we opt for a CrossMAE variant  
 279 with a 25% prediction ratio. In Table 3f, we observe that the classification accuracy positively  
 280 correlates with the input resolution, indicating that CrossMAE can scale to long input sequences.

### 281 4.4 Training Throughput and Memory Utilization

282 Due to partial reconstruction and confining attention to between mask tokens and visible tokens,  
 283 CrossMAE improves pre-training efficiency over MAE. Results in Table 10 show that the FLOPs





**Figure 5:** We visualize the output of each decoder block. (a-b) **Different decoder blocks play different roles in the reconstruction**, with most details emerging at later decoder blocks, which confirms the motivation for inter-block attention. (c) Visualizations of inter-block attention shows that **different decoder blocks indeed attend to feature from different encoder blocks**, with later blocks focusing on earlier encoder features to achieve reconstruction. The reconstructions are unnormalized w.r.t ground truth mean and std for each patch.

284 reduction does translate to an  $1.54\times$  training throughput and at least 50% reduction in GPU memory  
 285 utilization compared to MAE.

## 286 4.5 Visualizations

287 **Visualizing Per-block Reconstruction.** Rather than only visualizing the final reconstruction, we  
 288 have two key observations that allow us to visualize the work performed by each decoder block:  
 289 1) Transformer blocks have skip connections from their inputs to outputs. 2) The final decoder  
 290 block’s output goes through a linear reconstruction head to produce the reconstruction. As detailed in  
 291 Appendix D, we can factor out each block’s contribution in the final reconstruction with linearity.

292 This decomposition allows expressing the reconstruction as an image stack, where summing up all the  
 293 levels gives us the final reconstruction. As shown in Figure 5 (a,b), we observe that different decoder  
 294 blocks play different roles in reconstruction, with most details emerging at later decoder blocks. This  
 295 justifies the need for low-level features from early encoder blocks, motivating inter-block attention.

296 **Visualizing Inter-block Attention Maps.** As shown in the visualizations of the attention maps of  
 297 inter-block attention in 5(c), CrossMAE naturally leverages the inter-block attention to allow the later  
 298 decoder blocks to focus on earlier encoder features to achieve reconstruction and allow the earlier  
 299 decoder blocks to focus on later encoder features. This underscores the necessity for different decoder  
 300 blocks to attend to different encoder features, correlating with the performance improvements when  
 301 inter-block attention is used.

## 302 5 Discussion and Conclusion

303 In our study, we present a novel understanding of MAE, demonstrating that coherent image recon-  
 304 struction is achieved not through interactions between patches in the decoder but by learning a global  
 305 representation within the encoder. Based on this insight, we propose replacing self-attention layers  
 306 in the decoder with a simple readout function, specifically utilizing cross-attention to aggregate  
 307 encoder outputs into each input token within the decoder layers independently. This approach, tested  
 308 across models ranging from ViT-S to ViT-H, achieves comparable or better performance in image  
 309 classification and instance segmentation with reduced computational requirements, showcasing the  
 310 potential for more efficient and scalable visual pretraining methods. Our findings underscore the  
 311 efficacy of the encoder’s global representation learning, paving the way for streamlined decoder  
 312 architectures in future MAE implementations. CrossMAE’s efficiency and scalability demonstrate  
 313 potential for large-scale visual pretraining, particularly on underutilized in-the-wild video datasets.  
 314 However, our work has not yet explored scaling to models larger than ViT-H, the largest model  
 315 examined in MAE, leaving this for future research.

316 **References**

- 317 [1] Roman Bachmann, David Mizrahi, Andrei Atanov, and Amir Zamir. Multimae: Multi-modal multi-task  
318 masked autoencoders. *arXiv:2204.01678*, 2022.
- 319 [2] Roman Bachmann, David Mizrahi, Andrei Atanov, and Amir Zamir. Multimae: Multi-modal multi-task  
320 masked autoencoders. In *European Conference on Computer Vision*, pages 348–367. Springer, 2022.
- 321 [3] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv*  
322 *preprint arXiv:2106.08254*, 2021.
- 323 [4] Hangbo Bao, Li Dong, and Furu Wei. Beit: Bert pre-training of image transformers. In *ICLR*, 2022.
- 324 [5] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind  
325 Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners.  
326 2020.
- 327 [6] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey  
328 Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*,  
329 pages 213–229. Springer, 2020.
- 330 [7] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsu-  
331 pervised learning of visual features by contrasting cluster assignments. *Advances in neural information*  
332 *processing systems*, 33:9912–9924, 2020.
- 333 [8] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand  
334 Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF*  
335 *international conference on computer vision*, pages 9650–9660, 2021.
- 336 [9] Kai Chen, Zhili Liu, Lanqing Hong, Hang Xu, Zhenguo Li, and Dit-Yan Yeung. Mixed autoencoder for  
337 self-supervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer*  
338 *Vision and Pattern Recognition (CVPR)*, pages 22742–22751, 2023.
- 339 [10] Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever.  
340 Generative pretraining from pixels. 2020.
- 341 [11] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for  
342 contrastive learning of visual representations. In *International conference on machine learning*, pages  
343 1597–1607. PMLR, 2020.
- 344 [12] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive  
345 learning. *arXiv preprint arXiv:2003.04297*, 2020.
- 346 [13] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision  
347 transformers, 2021.
- 348 [14] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision  
349 transformers. *arXiv preprint arXiv:2104.02057*, 2021.
- 350 [15] Bowen Cheng, Alexander G. Schwing, and Alexander Kirillov. Per-pixel classification is not all you need  
351 for semantic segmentation. 2021.
- 352 [16] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-  
353 attention mask transformer for universal image segmentation. In *Proceedings of the IEEE/CVF conference*  
354 *on computer vision and pattern recognition*, pages 1290–1299, 2022.
- 355 [17] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data  
356 augmentation with a reduced search space. arxiv e-prints, page. *arXiv preprint arXiv:1909.13719*, 4, 2019.
- 357 [18] Tri Dao. FlashAttention-2: Faster attention with better parallelism and work partitioning. 2023.
- 358 [19] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirec-  
359 tional transformers for language understanding. 2019.
- 360 [20] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas  
361 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth  
362 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2020.

- 363 [21] Yuxin Fang, Li Dong, Hangbo Bao, Xinggang Wang, and Furu Wei. Corrupted image modeling for  
364 self-supervised visual pre-training. In *The Eleventh International Conference on Learning Representations*,  
365 2023.
- 366 [22] Christoph Feichtenhofer, Haoqi Fan, Yanghao Li, and Kaiming He. Masked autoencoders as spatiotemporal  
367 learners. In *Advances in Neural Information Processing Systems*, 2022.
- 368 [23] Xinyang Geng, Hao Liu, Lisa Lee, Dale Schuurams, Sergey Levine, and Pieter Abbeel. Multimodal  
369 masked autoencoders learn transferable representations. *arXiv preprint arXiv:2205.14204*, 2022.
- 370 [24] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew  
371 Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch sgd: Training imagenet in 1 hour.  
372 *arXiv:1706.02677*, 2017.
- 373 [25] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew  
374 Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv*  
375 *preprint arXiv:1706.02677*, 2017.
- 376 [26] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya,  
377 Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your  
378 own latent—a new approach to self-supervised learning. *Advances in neural information processing systems*,  
379 33:21271–21284, 2020.
- 380 [27] Ishaan Gulrajani, Kundan Kumar, Faruk Ahmed, Adrien Ali Taiga, Francesco Visin, David Vazquez, and  
381 Aaron Courville. Pixelvae: A latent variable model for natural images. *arXiv preprint arXiv:1611.05013*,  
382 2016.
- 383 [28] Agrim Gupta, Jiajun Wu, Jia Deng, and Li Fei-Fei. Siamese masked autoencoders. *arXiv preprint*  
384 *arXiv:2305.14344*, 2023.
- 385 [29] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised  
386 visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*  
387 *recognition*, pages 9729–9738, 2020.
- 388 [30] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders  
389 are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
390 *Recognition (CVPR)*, pages 16000–16009, 2022.
- 391 [31] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q Weinberger. Deep networks with stochastic  
392 depth. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October*  
393 *11–14, 2016, Proceedings, Part IV 14*, pages 646–661. Springer, 2016.
- 394 [32] Andrew Jaegle, Sebastian Borgeaud, Jean-Baptiste Alayrac, Carl Doersch, Catalin Ionescu, David Ding,  
395 Skanda Koppula, Daniel Zoran, Andrew Brock, Evan Shelhamer, et al. Perceiver io: A general architecture  
396 for structured inputs & outputs. *arXiv preprint arXiv:2107.14795*, 2021.
- 397 [33] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao,  
398 Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollar, and Ross Girshick. Segment anything.  
399 In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4015–4026,  
400 2023.
- 401 [34] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut.  
402 Albert: A lite bert for self-supervised learning of language representations. In *International Conference on*  
403 *Learning Representations*, 2020.
- 404 [35] Anders Boesen Lindbo Larsen, Søren Kaae Sønderby, Hugo Larochelle, and Ole Winther. Autoencoding  
405 beyond pixels using a learned similarity metric. In *International conference on machine learning*, pages  
406 1558–1566. PMLR, 2016.
- 407 [36] Jin Li, Yaoming Wang, XIAOPENG ZHANG, Yabo Chen, Dongsheng Jiang, Wenrui Dai, Chenglin Li,  
408 Hongkai Xiong, and Qi Tian. Progressively compressed auto-encoder for self-supervised representation  
409 learning. In *The Eleventh International Conference on Learning Representations*, 2023.
- 410 [37] Tianhong Li, Huiwen Chang, Shlok Kumar Mishra, Han Zhang, Dina Katabi, and Dilip Krishnan.  
411 Mage: Masked generative encoder to unify representation learning and image synthesis. *arXiv preprint*  
412 *arXiv:2211.09117*, 2022.
- 413 [38] Yanghao Li, Hanzi Mao, Ross Girshick, and Kaiming He. Exploring plain vision transformer backbones  
414 for object detection. In *European Conference on Computer Vision*, pages 280–296. Springer, 2022.

- 415 [39] Yanghao Li, Haoqi Fan, Ronghang Hu, Christoph Feichtenhofer, and Kaiming He. Scaling language-image  
416 pre-training via masking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
417 Recognition*, pages 23390–23400, 2023.
- 418 [40] Jihao Liu, Xin Huang, Jinliang Zheng, Yu Liu, and Hongsheng Li. Mixmae: Mixed and masked autoencoder  
419 for efficient pretraining of hierarchical vision transformers. *arXiv:2205.13137*, 2022.
- 420 [41] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis,  
421 Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv  
422 preprint arXiv:1907.11692*, 2019.
- 423 [42] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. 2017.
- 424 [43] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint  
425 arXiv:1711.05101*, 2017.
- 426 [44] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive  
427 coding. *arXiv preprint arXiv:1807.03748*, 2018.
- 428 [45] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre  
429 Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual  
430 features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- 431 [46] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders:  
432 Feature learning by inpainting. In *Proceedings of the IEEE conference on computer vision and pattern  
433 recognition*, pages 2536–2544, 2016.
- 434 [47] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding  
435 by generative pre-training. 2018.
- 436 [48] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language  
437 models are unsupervised multitask learners. 2019.
- 438 [49] Ilija Radosavovic, Tete Xiao, Stephen James, Pieter Abbeel, Jitendra Malik, and Trevor Darrell. Real-world  
439 robot learning with masked visual pre-training. In *Conference on Robot Learning*, pages 416–426. PMLR,  
440 2023.
- 441 [50] Andreas Peter Steiner, Alexander Kolesnikov, Xiaohua Zhai, Ross Wightman, Jakob Uszkoreit, and Lucas  
442 Beyer. How to train your vit? data, augmentation, and regularization in vision transformers. *Transactions  
443 on Machine Learning Research*, 2022.
- 444 [51] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the  
445 inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and  
446 pattern recognition*, pages 2818–2826, 2016.
- 447 [52] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. VideoMAE: Masked autoencoders are data-efficient  
448 learners for self-supervised video pre-training. In *Advances in Neural Information Processing Systems*,  
449 2022.
- 450 [53] Hugo Touvron, Matthieu Cord, Alaaeldin El-Nouby, Piotr Bojanowski, Armand Joulin, Gabriel Synnaeve,  
451 and Hervé Jégou. Augmenting convolutional networks with attention-based aggregation, 2021.
- 452 [54] Shubham Tulsiani and Abhinav Gupta. Pixeltransformer: Sample conditioned signal generation. In  
453 *Proceedings of the 38th International Conference on Machine Learning*, pages 10455–10464. PMLR,  
454 2021.
- 455 [55] Aaron Van den Oord, Nal Kalchbrenner, Lasse Espeholt, Oriol Vinyals, Alex Graves, et al. Conditional  
456 image generation with pixellcn decoders. *Advances in neural information processing systems*, 29, 2016.
- 457 [56] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz  
458 Kaiser, and Illia Polosukhin. Attention is all you need. 2017.
- 459 [57] Pascal Vincent, Hugo Larochelle, Isabelle Lajoie, Yoshua Bengio, Pierre-Antoine Manzagol, and Léon  
460 Bottou. Stacked denoising autoencoders: Learning useful representations in a deep network with a local  
461 denoising criterion. *Journal of machine learning research*, 11(12), 2010.
- 462 [58] Limin Wang, Bingkun Huang, Zhiyu Zhao, Zhan Tong, Yinan He, Yi Wang, Yali Wang, and Yu Qiao.  
463 Videomae v2: Scaling video masked autoencoders with dual masking. In *Proceedings of the IEEE/CVF  
464 Conference on Computer Vision and Pattern Recognition*, pages 14549–14560, 2023.

- 465 [59] Xudong Wang, Ziwei Liu, and Stella X Yu. Unsupervised feature learning by cross-level instance-group  
466 discrimination. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,  
467 pages 12586–12595, 2021.
- 468 [60] Chen Wei, Karttikeya Mangalam, Po-Yao Huang, Yanghao Li, Haoqi Fan, Hu Xu, Huiyu Wang, Cihang  
469 Xie, Alan Yuille, and Christoph Feichtenhofer. Diffusion models as masked autoencoder. In *ICCV*, 2023.
- 470 [61] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han Hu.  
471 Simmim: A simple framework for masked image modeling. In *Proceedings of the IEEE/CVF Conference*  
472 *on Computer Vision and Pattern Recognition (CVPR)*, pages 9653–9663, 2022.
- 473 [62] Hongwei Xue, Peng Gao, Hongyang Li, Yu Qiao, Hao Sun, Houqiang Li, and Jiebo Luo. Stare at what  
474 you see: Masked image modeling without reconstruction. In *Proceedings of the IEEE/CVF Conference on*  
475 *Computer Vision and Pattern Recognition*, pages 22732–22741, 2023.
- 476 [63] Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo.  
477 Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the*  
478 *IEEE/CVF international conference on computer vision*, pages 6023–6032, 2019.
- 479 [64] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk  
480 minimization. In *International Conference on Learning Representations*, 2018.
- 481 [65] Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong. ibot: Image  
482 bert pre-training with online tokenizer. *arXiv preprint arXiv:2111.07832*, 2021.

## 483 A Implementation details

### 484 A.1 Attention Calculation

485 To compare the attention values for mask tokens in vanilla MAE (Figure 1), we trained a ViT-B/16  
486 MAE for 800 epochs using the default hyperparameters provided in [30]. For each image, we  
487 randomly generate a 75% binary mask ( $m$ ) for all tokens, with  $m_i = 1$  representing a token being  
488 masked and  $m_i = 0$  otherwise. During the forward pass of the decoder, for each self-attention  
489 operation, the attention map is stored. This means that for the default MAE, a total of 8 attention  
490 maps, each with 16 attention heads are stored. Based on the mask pattern, we calculate the outer  
491 product ( $m \cdot m^T$ ) for the self-attention among mask tokens, and  $m \cdot (1 - m^T)$  for the cross-attention  
492 from the mask token to the visible tokens. We then calculate the average across all feature maps  
493 and attention heads for self-attention and cross-attention to get the image average values. Lastly, we  
494 averaged across the entire ImageNet validation set to obtain the final values.

### 495 A.2 Inter-Block Attention

496 We tried a few implementations for inter-block attention (IBA) and found the following implementa-  
497 tion to be the fastest and most memory-efficient. In this implementation, we combine inter-block  
498 attention for all encoder layers as a single forward pass of a linear layer. For each decoder block,  
499 we index into the output tensor to extract the corresponding feature map, and a layer norm will be  
500 applied before the feature map is fed into the decoder block. Other alternatives we tried include 1)  
501 performing separate inter-block attentions before each decoder block, and 2) 1x1 convolution on the  
502 stacked encoder feature maps.

503 In MAE, there exists a layer norm after the last encoder feature map before feeding into the decoder.  
504 In our implementation, we only add layer norm after inter-block attention. We find that adding  
505 an additional layer norm before inter-block attention to each encoder feature map does not lead to  
506 improvements in model performance but will significantly increase GPU memory usage.

507 The pseudo-code of inter-block attention is the following:

```
508 1 class InterBlockAttention():  
509 2     def __init__(self, num_feat_maps, decoder_depth):  
510 3         self.linear = Linear(num_feat_maps, decoder_depth, bias=False)  
511 4         std_dev = 1. / sqrt(num_feat_maps)  
512 5         init.normal_(self.linear.weight, mean=0., std=std_dev)  
513 6  
514 7     def forward(self, feature_maps : list):  
515 8         """  
516 9         feature_maps: a list of length num_feat_maps, each with  
517 dimension  
518 10 Batch Size x Num. Tokens x Embedding Dim.  
519 11         """  
520 12         stacked_feature_maps = stack(feature_maps, dim=-1)  
521 13         return self.linear(stacked_feature_maps)
```

522 Additionally, we further investigate the importance of using a cross-attention decoder, where each  
523 decoder block can use different feature maps from the encoder for decoding. In this experiment, we  
524 incorporated IBA into MAE, which uses only a self-attention decoder. Specifically, we concatenate  
525 the interblock attention features with the masked tokens. We then feed the combined features into  
526 MAE’s self-attention decoder. We pre-trained the model and finetuned it for Imagenet classification.  
527 The results are presented in Table. 4, where all models are pre-trained for 400 epochs. We observe that  
528 inter-block attention has negligible performance improvements for MAE, potentially because MAE  
529 only takes in one feature map in its decoder. In contrast, inter-block attention allows cross-attention  
530 layers in CrossMAE to attend to features from different encoder blocks, thanks to its decoupling of  
531 queries with keys and values.

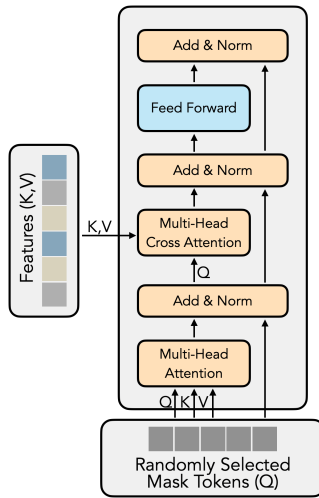
### 532 A.3 Ablation that Adds Self-Attention

533 In Section 4.3 (a), we propose adding self-attention back to CrossMAE as an ablation. In that  
534 particular ablation study, we analyze the effect of self-attention between the masked tokens, which

Method	Acc. (%)
MAE	83.0
MAE + IBA	83.0
CrossMAE (25%)	83.2
CrossMAE (75%)	<b>83.3</b>

**Table 4:** For MAE, inter-block attention has very small differences in terms of finetuning performance, potentially due to the fact that MAE’s decoder only takes in one set of features.

535 can be used to improve the consistency for reconstruction. Specifically, we modify the formulation in  
 536 the original transformer paper [56], where the mask/query tokens are first passed through a multi-  
 537 head self-attention and a residual connection before being used in the multiheaded cross-attention  
 538 with the features from the encoder. The primary difference with the vanilla transformer decoder  
 539 implementation [56] is we do not perform casual masking in the multi-head self-attention. Please  
 540 reference Figure 6 for a more visual presentation of the method.



**Figure 6:** Modification for self-attention ablation

#### 541 A.4 Ablation on Inter-block Attention

542 In Table 3d, the following cases are considered. 1 feature map (row 1) does not use inter-block  
 543 attention. Each decoder block only takes the last feature map from the encoder as the keys and values.  
 544 For scenarios where more than one feature map is used, the output of the patch embedding (input to  
 545 the ViT) is also used.

546 In addition to the simple design of inter-block attention proposed above, we also experimented  
 547 with a variant of inter-block attention by further parameterizing the attention with linear projections.  
 548 Specifically, rather than directly performing weighted sum aggregation to form the features for each  
 549 cross-attention layer in the decoder, we added a linear projection for each encoder feature before the  
 550 feature aggregation. We denote this variant as *CrossMAE+LP*. As shown in the Table. 5 (with ViT-B  
 551 pre-trained for 400 epochs, consistent with the setting in Table. 3), adding a linear projection slightly  
 552 improves the performance. This indicates that it is possible to design variants of readout functions,  
 553 such as through improved inter-block attention, to improve the feature quality of CrossMAE.

Method	Acc. (%)
CrossMAE	83.3
CrossMAE + LP	<b>83.5</b>

**Table 5:** Improving inter-block attention by adding linear projections to the input features. The performance gain indicates that it is possible to design variants of readout functions to improve CrossMAE.

554 **A.5 Hyperparameters**

555 **Pre-training:** The default setting is in Table 6, which is consistent with the official MAE [30]  
 556 implementation. As mentioned in Sec. 3.4, we scale the learning rate by the ratio between mask ratio  
 557 ( $p$ ) and prediction ratio ( $\gamma$ ) to ensure the variance of the loss is consistent with [30]. Additionally, we  
 558 use the linear learning rate scaling rule [25]. This results in  $lr = \gamma * base\_lr * batchsize / (256 * p)$ .  
 559 For Table 1, we use 12 decoder blocks, with mask ratio and prediction ratio both 75%, and interblock  
 560 attention takes in all encoder feature maps. For the 400 epochs experiments in Table 2, we scale the  
 561 warm-up epochs correspondingly. Other hyperparameters, such as decoder block width, are the same  
 562 as MAE.

563 **Finetuning:** We use the same hyperparameters as MAE finetuning. We use global average pooling  
 564 for finetuning. In MAE, the layer norm for the last encoder feature map is removed for finetuning,  
 which is consistent with our pretraining setup. Please refer to Table 7 for more detail.

Config	Value
optimizer	AdamW [43]
base learning rate	1.5e-4
learning rate schedule	cosine decay [42]
batch size	4096
weight decay	0.05
optimizer momentum	$\beta_1, \beta_2 = 0.9, 0.95$ [10]
warm up epoch [24]	20, 40
total epochs	400, 800
augmentation	RandomResizedCrop, RandomHorizontalFlip

**Table 6:** Pretraining Hyperparameters

565

566 **A.6 Compute Infrastructure**

567 Each of the pretraining and finetuning experiments is run on 2 or 4 NVIDIA A100 80GB GPUs. The  
 568 batch size per GPU is scaled accordingly and we use gradient accumulation to avoid out-of-memory  
 569 errors. ViTDet [38] experiments use a single machine equipped with 8 NVIDIA A100 (80GB) GPUs.  
 570 We copy the datasets to the shared memory on the machines to accelerate dataloading. We use  
 571 FlashAttention-2 [18] to accelerate attention calculation.

Config	Value
optimizer	AdamW
base learning rate	1e-3
learning rate schedule	cosine decay
batch size	1024
weight decay	0.05
optimizer momentum	$\beta_1, \beta_2 = 0.9, 0.999$
warm up epoch	5
total epochs	100 (B), 50 (L)
augmentation	RandAug (9, 0.5) [17]
label smoothing [51]	0.1
mixup [64]	0.8
cutmix [63]	1.0
drop path [31]	0.1

**Table 7:** Finetuning Hyperparameters



## 572 B Additional Experiments

### 573 B.1 Linear Probe

574 We provide linear probe comparisons (at 800 epochs) for ViT-Small and ViT-Base in Table. 8. For both  
 575 of these experiments, we run CrossMAE with a prediction ratio of 75% (reconstruction of all masked  
 576 patches). These results show that CrossMAE achieves slightly better linear probe performance than  
 577 vanilla MAE.

Method	ViT-S	ViT-B
MAE	49.7	65.1
CrossMAE	<b>51.5</b>	<b>65.4</b>

**Table 8:** Linear probe experiments of CrossMAE.

### 578 B.2 Masking Strategy

Method	Acc. (%)
Grid Masking	83.2
Random Masking	<b>83.3</b>

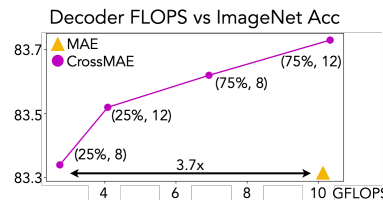
**Table 9:** Ablation of masking strategies.

579 Similar to MAE [30], we here ablate the masking pattern. Instead of random masking, we perform  
 580 grid-wise sampling that “keeps one of every four patches” (see MAE Figure 6). The finetuning  
 581 performance is reported in Table. 9 for ViT-B (at 400 epochs), which shows that grid masking does  
 582 not lead to additional improvements in downstream performance.

## 583 C Runtime and GPU Memory Comparisons with MAE

Method	Memory (MB/GPU)	Runtime (min/epoch)	Acc. (%)
MAE	OOM (>81920)	5.19*	83.3
CrossMAE	<b>41177</b>	<b>3.38</b>	<b>83.5</b>

**Table 10: CrossMAE greatly improves the training throughput and reduces the memory requirements,** lowering the barrier for masked pretraining. Statistics are measured on 2 NVIDIA A100 80GB GPUs. Please refer to Appendix C for comparison details. \*: MAE’s default batch size exceeds the capacity of 4 GPUs, requiring gradient accumulation for runtime measurement.



**Figure 7:** We compare ViT-B which is pre-trained for 800 epochs with different variants of CrossMAE v.s. MAE. For CrossMAE, we vary the prediction ratio  $p$  and number of decoder blocks  $n$ , and we denote each as  $(p, n)$ . While all experiments are run with inter-block attention, CrossMAE has lower decoder FLOPS than MAE [30] and performs on par or better.

584 All experiments in Table 10 are conducted on a server with 4 NVIDIA A100 (80GB) GPUs, with the  
 585 standard hyperparameters provided above for pretraining. NVLink is equipped across the GPUs. We  
 586 use the default setting for MAE and set the global batch size to 4096. For CrossMAE, we also use  
 587 the default setting with a prediction ratio 0.25, and this takes around 41GB memory per GPU without  
 588 gradient accumulation (i.e., local batch size is set to 1024 samples per GPU). However, the same  
 589 local batch size results in out-of-memory (OOM), which indicates that the total memory requirement  
 590 is larger than the available memory for each GPU (80GB). To run MAE on same hardware, we  
 591 thus employ gradient accumulation with a local batch size of 512 to maintain the global batch size.  
 592 The benchmark runs each method and measures the average per epoch runtime as well as the max  
 593 memory allocation for 10 training epochs. Our experiments in Figure 7 show that models with lower  
 594 prediction ratios benefit more from deeper decoders. Our model performs on par or better when  
 595 compared to MAE, with up to  $3.7\times$  lower decoder FLOPS.

596 **D Visualizing the Contributions per Decoder Block**

597 We propose a more fine-grained visualization approach that allows us to precisely understand the  
 598 effect and contribution of each decoder block.

599 Two key observations enable per-block visualization: **1)** Transformer blocks have residual connections  
 600 from their inputs to outputs. Let  $f_i$  be the output and  $g_i(\cdot)$  the residual function of decoder  $i$ , so  
 601  $f_i = f_{i-1} + g_i(f_{i-1})$ . **2)** The final decoder block’s output goes through a reconstruction head  $h$ ,  
 602 which is linear, consisting of a layer-norm and a linear layer, to produce the reconstruction. With  
 603  $D$  as the decoder depth,  $f_0$  the initial input, and  $y$  the final output,  $y$  is recursively defined as  
 604  $y = h(f_{D-1} + g_D(f_{D-1}))$ , which simplifies due to the linearity of  $h$ :

$$\begin{aligned}
 \mathbf{y} &= h(f_0 + g_1(f_0) + \cdots + g_D(f_{D-1})) \\
 &= \underbrace{h(f_0)}_{\text{Pos Embed. + Mask Token}} + \underbrace{h(g_1(f_0))}_{\text{Block 1}} + \cdots + \underbrace{h(g_D(f_{D-1}))}_{\text{Block D}}
 \end{aligned}$$

605 This decomposition allows us to express the reconstruction as an image stack, where the sum of all  
 606 the levels gives us the final reconstruction. We present the visualization in Figure 5.

## 607 **NeurIPS Paper Checklist**

### 608 **1. Claims**

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

611 Answer: [\[Yes\]](#)

612 Justification: The claims in the abstract are justified in the method and the experiments  
613 section.

614 Guidelines:

- 615 • The answer NA means that the abstract and introduction do not include the claims  
616 made in the paper.
- 617 • The abstract and/or introduction should clearly state the claims made, including the  
618 contributions made in the paper and important assumptions and limitations. A No or  
619 NA answer to this question will not be perceived well by the reviewers.
- 620 • The claims made should match theoretical and experimental results, and reflect how  
621 much the results can be expected to generalize to other settings.
- 622 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
623 are not attained by the paper.

### 624 **2. Limitations**

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

626 Answer: [\[Yes\]](#)

627 Justification: The limitations of the work have been discussed in the Discussion and Conclu-  
628 sion section.

629 Guidelines:

- 630 • The answer NA means that the paper has no limitation while the answer No means that  
631 the paper has limitations, but those are not discussed in the paper.
- 632 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 633 • The paper should point out any strong assumptions and how robust the results are to  
634 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
635 model well-specification, asymptotic approximations only holding locally). The authors  
636 should reflect on how these assumptions might be violated in practice and what the  
637 implications would be.
- 638 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
639 only tested on a few datasets or with a few runs. In general, empirical results often  
640 depend on implicit assumptions, which should be articulated.
- 641 • The authors should reflect on the factors that influence the performance of the approach.  
642 For example, a facial recognition algorithm may perform poorly when image resolution  
643 is low or images are taken in low lighting. Or a speech-to-text system might not be  
644 used reliably to provide closed captions for online lectures because it fails to handle  
645 technical jargon.
- 646 • The authors should discuss the computational efficiency of the proposed algorithms  
647 and how they scale with dataset size.
- 648 • If applicable, the authors should discuss possible limitations of their approach to  
649 address problems of privacy and fairness.
- 650 • While the authors might fear that complete honesty about limitations might be used by  
651 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
652 limitations that aren't acknowledged in the paper. The authors should use their best  
653 judgment and recognize that individual actions in favor of transparency play an impor-  
654 tant role in developing norms that preserve the integrity of the community. Reviewers  
655 will be specifically instructed to not penalize honesty concerning limitations.

### 656 **3. Theory Assumptions and Proofs**

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

659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711

Answer: [NA]

Justification: This work offers observations and hypotheses justified with empirical results.

Guidelines:

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

#### 4. Experimental Result Reproducibility

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

Answer: [Yes]

Justification: Our code, which reproduces our results, is provided through an anonymous link in the abstract.

Guidelines:

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

#### 5. Open access to data and code

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

715 Answer: [Yes]

716 Justification: Our method is evaluated on open datasets that are publicly available.

717 Guidelines:

- 718 • The answer NA means that paper does not include experiments requiring code.
- 719 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
720 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 721 • While we encourage the release of code and data, we understand that this might not be  
722 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
723 including code, unless this is central to the contribution (e.g., for a new open-source  
724 benchmark).
- 725 • The instructions should contain the exact command and environment needed to run to  
726 reproduce the results. See the NeurIPS code and data submission guidelines ([https:  
727 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 728 • The authors should provide instructions on data access and preparation, including how  
729 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 730 • The authors should provide scripts to reproduce all experimental results for the new  
731 proposed method and baselines. If only a subset of experiments are reproducible, they  
732 should state which ones are omitted from the script and why.
- 733 • At submission time, to preserve anonymity, the authors should release anonymized  
734 versions (if applicable).
- 735 • Providing as much information as possible in supplemental material (appended to the  
736 paper) is recommended, but including URLs to data and code is permitted.

## 737 6. Experimental Setting/Details

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

741 Answer: [Yes]

742 Justification: We follow the hyperparam selection from MAE. The hyperparams introduced  
743 by our work, such as the mask ratio and the number of feature maps used, are ablated.

744 Guidelines:

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

## 750 7. Experiment Statistical Significance

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

753 Answer: [No]

754 Justification: Error bars are not reported because they would be too computationally expen-  
755 sive.

756 Guidelines:

- 757 • The answer NA means that the paper does not include experiments.
- 758 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
759 dence intervals, or statistical significance tests, at least for the experiments that support  
760 the main claims of the paper.
- 761 • The factors of variability that the error bars are capturing should be clearly stated (for  
762 example, train/test split, initialization, random drawing of some parameter, or overall  
763 run with given experimental conditions).

- 764 • The method for calculating the error bars should be explained (closed form formula,  
765 call to a library function, bootstrap, etc.)
- 766 • The assumptions made should be given (e.g., Normally distributed errors).
- 767 • It should be clear whether the error bar is the standard deviation or the standard error  
768 of the mean.
- 769 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
770 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
771 of Normality of errors is not verified.
- 772 • For asymmetric distributions, the authors should be careful not to show in tables or  
773 figures symmetric error bars that would yield results that are out of range (e.g. negative  
774 error rates).
- 775 • If error bars are reported in tables or plots, The authors should explain in the text how  
776 they were calculated and reference the corresponding figures or tables in the text.

## 777 8. Experiments Compute Resources

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

781 Answer: [Yes]

782 Justification: We described the compute requirements in Appendix A.6. We do not use  
783 GPUs from a cloud provider.

784 Guidelines:

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

## 793 9. Code Of Ethics

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

796 Answer: [Yes]

797 Justification: The research conforms to the NeurIPS Code of Ethics.

798 Guidelines:

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

## 804 10. Broader Impacts

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

807 Answer: [Yes]

808 Justification: This paper aims to advance the field of self-supervised learning. Like other self-  
809 supervised learning methods, our work may have various societal implications. However,  
810 we do not believe any specific consequences need to be highlighted in this context.

811 Guidelines:

- 812 • The answer NA means that there is no societal impact of the work performed.
- 813 • If the authors answer NA or No, they should explain why their work has no societal  
814 impact or why the paper does not address societal impact.

- 815
- 816
- 817
- 818
- 819
- 820
- 821
- 822
- 823
- 824
- 825
- 826
- 827
- 828
- 829
- 830
- 831
- 832
- 833
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
  - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
  - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
  - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 834 11. Safeguards

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

838 Answer: [NA]

839 Justification: The paper does not pose such risks.

840 Guidelines:

- 841
- 842
- 843
- 844
- 845
- 846
- 847
- 848
- 849
- 850
- The answer NA means that the paper poses no such risks.
  - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
  - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
  - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 851 12. Licenses for existing assets

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

855 Answer: [Yes]

856 Justification: The code and datasets used in this work follow the original MAE work.

857 Guidelines:

- 858
- 859
- 860
- 861
- 862
- 863
- 864
- 865
- 866
- 867
- 868
- The answer NA means that the paper does not use existing assets.
  - The authors should cite the original paper that produced the code package or dataset.
  - The authors should state which version of the asset is used and, if possible, include a URL.
  - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
  - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
  - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.

- 869
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- 870
- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.
- 871
- 872

873 **13. New Assets**

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

876 Answer: [NA]

877 Justification: The paper does not release new assets.

878 Guidelines:

- The answer NA means that the paper does not release new assets.
  - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
  - The paper should discuss whether and how consent was obtained from people whose asset is used.
  - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
- 879
- 880
- 881
- 882
- 883
- 884
- 885
- 886

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

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

891 Answer: [NA]

892 Justification: The paper does not involve crowdsourcing or research with human subjects.

893 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
  - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
  - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.
- 894
- 895
- 896
- 897
- 898
- 899
- 900
- 901

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

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

908 Answer: [NA]

909 Justification: The paper does not involve crowdsourcing or research with human subjects.

910 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
  - Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
  - We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
  - For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.
- 911
- 912
- 913
- 914
- 915
- 916
- 917
- 918
- 919
- 920