

# Local Masked Reconstruction for Efficient Self-Supervised Learning on High-resolution Images

Jun Chen<sup>1</sup>\*, Faizan Farooq Khan<sup>1</sup>\*, Ming Hu<sup>2</sup>, Ammar Sherif<sup>3</sup>,  
Zongyuan Ge<sup>2</sup>, Boyang Li<sup>4</sup>, Mohamed Elhoseiny<sup>1</sup>  
{jun.chen, faizan.khan, mohamed.elhoseiny}@kaust.edu.sa  
{ming.hu, zongyuan.ge}@monash.edu  
{libo0001@gmail.com, asherif@nu.edu.eg}  
<sup>1</sup>King Abdullah University of Science and Technology  
<sup>2</sup>Monash University <sup>3</sup>Nile University <sup>4</sup>Nanyang Technological University

## Abstract

*Self-supervised learning for computer vision has progressed tremendously and improved many downstream vision tasks, such as image classification, semantic segmentation, and object detection. Among these, generative self-supervised vision learning approaches, such as MAE and BEiT, show promising performance. However, their global reconstruction mechanism is computationally demanding, especially for high-resolution images. The computational cost increases extensively when scaled to a large-scale dataset. To address this issue, we propose local masked reconstruction (LoMaR), a simple yet effective approach that reconstructs image patches from small neighboring regions. The strategy can be easily integrated into any generative self-supervised learning techniques and improves the trade-off between efficiency and accuracy compared to reconstruction over the entire image. LoMaR is 2.5× faster than MAE and 5.0× faster than BEiT on 384×384 ImageNet pretraining and surpasses them by 0.2% and 0.8% in accuracy, respectively. It is 2.1× faster than MAE on iNaturalist pretraining and gains 0.2% in accuracy. On MS COCO, LoMaR outperforms MAE by 0.5 AP<sup>box</sup> on object detection and 0.5 AP<sup>mask</sup> on instance segmentation. It also outperforms MAE by 0.2% on semantic segmentation. Our code and pretrained models are available at: <https://github.com/junchen14/LoMaR>.*

## 1. Introduction

Recently, self-supervised learning [2, 6, 13, 14, 20, 30, 34, 45, 61] has achieved enormous success in learning representations conducive to downstream applications, such as

image classification and object detection. Among these, several generative methods, such as Masked Autoencoder (MAE) [30] and Bidirectional Encoder Representation from Image Transformers (BEiT) [6], which reconstruct the input image from a small portion of image patches, have demonstrated excellent performance.

However, a significant bottleneck of MAE and BEiT is their high demand for compute, as they reconstruct masked image patches from global information and operate on a large number of image patches. For example, pretraining an MAE-Huge network on ImageNet under  $224 \times 224$  resolution takes 34.5 hours on 128 TPU-v3 GPUs. BEiT [6] training is even slower due to the cost of the discrete variational autoencoder.

High-resolution images further exacerbate this issue due to the  $\mathcal{O}(n^2)$  time complexity of the Transformer model on  $n$  image patches. For example, pretraining MAE on  $384 \times 384$  images consumes 4.7 times the compute time of  $224 \times 224$  counterpart. However, high-resolution images are essential in many tasks, such as object detection. Thus, improving the efficiency of pretraining holds the promise to unleash additional performance gains under pretraining with a much larger dataset or higher-resolution images.

We observe that most reconstruction in MAE relies on local information only. In Fig. 1, we visualize the attention weights (white indicating high attention) when reconstructing a target image patch. From a pretrained MAE<sub>Large</sub> model, we extract the attention weights from the decoder layers 2, 4, 6, and 8. The model mainly attends to patches close to the target patch, which motivates us to limit the range of attention used in the reconstruction.

Hence, we propose a new model, dubbed **Local Masked Reconstruction** or LoMaR. The model restricts the attention region to a small window, such as  $7 \times 7$  image patches, which is sufficient for reconstruction. Similar approaches

\*Equal Contribution

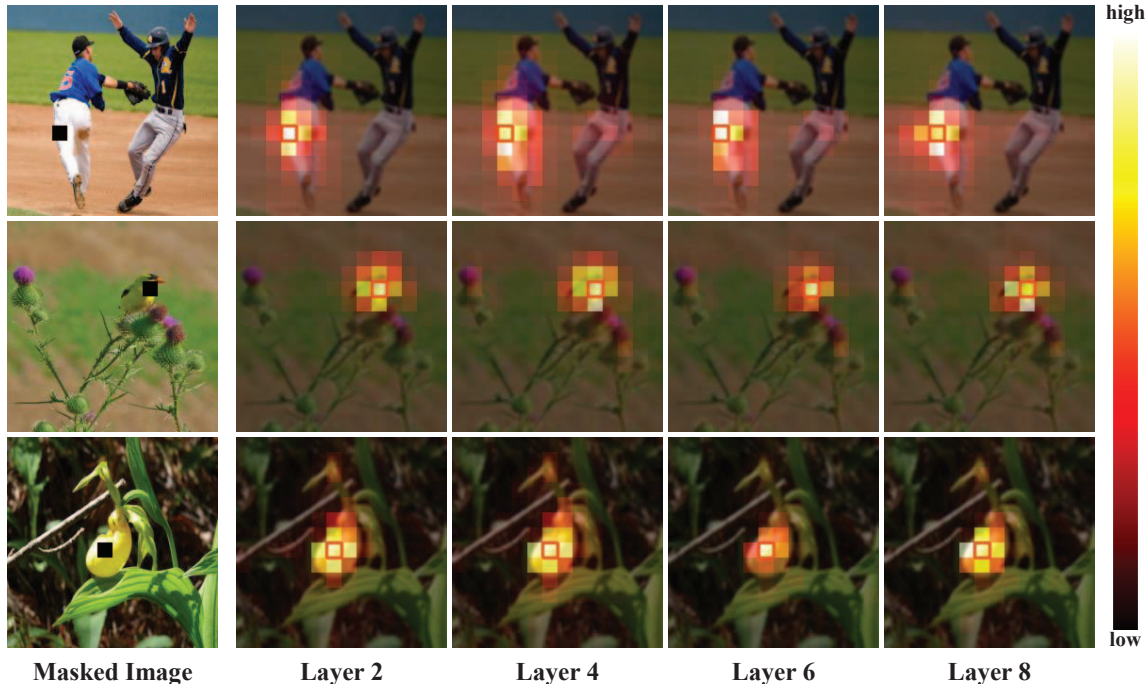


Figure 1. We visualize the attention patterns employed by  $\text{MAE}_{\text{Large}}$  [30] in the reconstruction of a random target patch, indicated by orange. Patches that are important for prediction are usually close to the target patch. We selected the images randomly from the ImageNet-1K [22] Val set.

[21, 55, 64] have been seen in NLP problems that need to process long sequences. The small windows have also been explored in vision domains for higher training and inference speed [42, 63]. However, unlike prior work in vision transformers, which create shifting windows with fixed coordinates for each image, we sample several windows with random locations, which can better capture the objects in different spatial areas.

In Figure 2, we compare LoMaR and MAE and note two significant differences: a) We sample a region with  $k \times k$  patches to perform masked reconstruction instead of from the total number of patches. Instead of reconstructing the masked patches from the 25% visible patches globally located in the image, we find that it is sufficient to recover the missing information with only some local visual clues. b) We replace the heavy-weight decoder in MAE with a lightweight MLP head. We feed all image patches directly into the encoder, including masked and visible patches. In comparison, only the visible patches are fed to the encoder in MAE. Experiments show that these architectural changes bring more performance gain to the local masked reconstruction in small regions.

After conducting extensive experiments, we found that

- LoMaR is more efficient than other baselines in pre-training on high-resolution images since its computation cost is invariant to the different image resolutions.

However, other approaches have a quadratic computational cost to the image resolution increase, which leads to much more expensive pretraining. For example, for pretraining on  $448 \times 448$  images, LoMaR is  $3.1 \times$  faster than MAE and  $5.3 \times$  faster than BEiT while achieving higher classification performance.

- LoMaR also has a strong generalization ability on object detection and semantic segmentation tasks. It outperforms MAE by  $0.5 \text{ AP}^{\text{box}}$  under ViTDet [40] framework for object detection. Also, it outperforms MAE by 0.2 points under UperNet [62] for semantic segmentation.
- LoMaR is efficient and can be easily integrated into any other generative self-supervised learning approach. Equipping our local masked reconstruction learning mechanism into BEiT can improve its ImageNet-1K classification performance from 83.2 to 83.4 Top-1 accuracy, costing only 35.8% of its original pretraining time.

## 2. Related Work

**Self-supervised Learning For Images.** The past few years have witnessed a boom in self-supervised learning. Existing techniques can be broadly categorized into discriminative [8] and generative [27, 46]. The prominent discrimina-

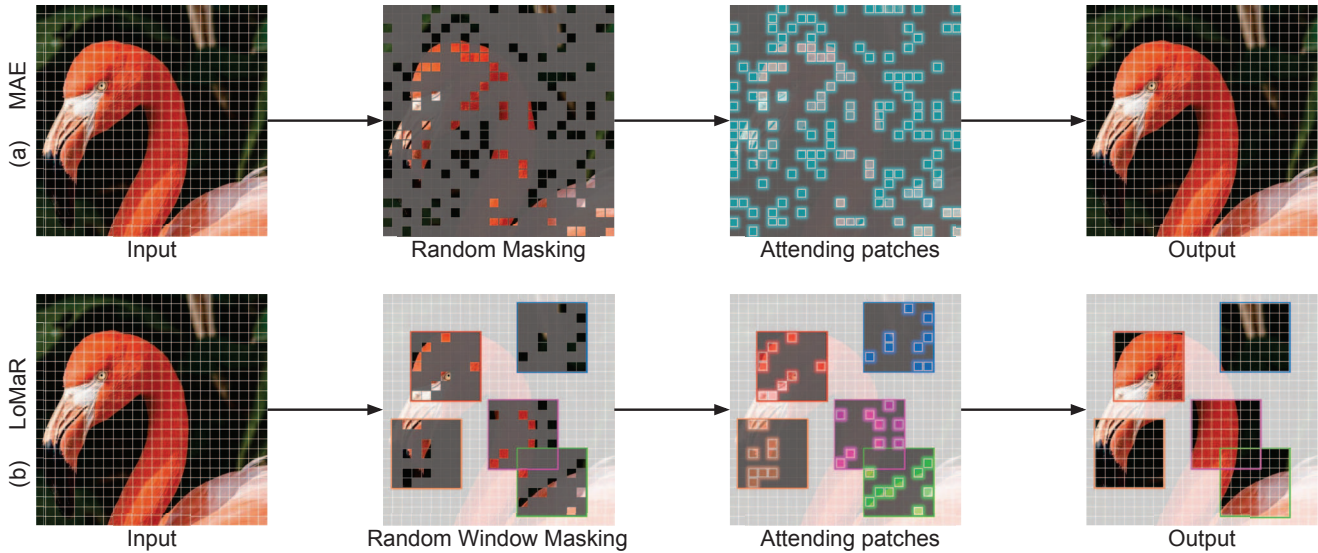


Figure 2. **Contrasting the masking and reconstruction strategy between MAE and LoMaR.** During the pretraining, MAE randomly masks 75% patches as masking and reconstructs them by attending to the remaining visible patches. For LoMaR, it randomly samples several small regions and masks a random subset of patches from each region, e.g. 80%. The masked patches will only attend to the visible patches inside each region for reconstruction. In contrast to MAE, LoMaR usually samples less visible patches per image.

tive approach, instance discrimination, distinguishes different views of the same data instance from other instances [15, 19, 28, 31, 45, 61]. The most representative works include BYOL [28], MOCO [18, 20, 31] and SimCLR [15, 16]. VICRegL [7] performs the contrasting learning in local and global features. Other approaches such as SwAV [12] and DINO [13] heavily rely on multi-crop strategy, for instance, discrimination. The generative approach includes autoregressive prediction [14, 29] and autoencoders, which we discuss next.

**Autoencoders for Representation Learning.** An autoencoder, which aims to learn a representation from which the original input can be reconstructed, has been a popular choice for representation learning since the dawn of deep learning [9, 10, 33]. The autoencoding problem is inherently ill-posed due to the existence of a trivial solution: a network entirely composed of identity mappings. Hence, it is necessary to apply some form of regularization, such as sparsity [52], input corruption [59], probability priors [37, 53], or adversarial discriminators [43].

In particular, denoising autoencoder [59], which attempts to recover the original input from a corrupted version, has received significant research attention. Variations include solving jigsaw puzzles [44], color restoration [38, 66], spatial relation recovery [23], inpainting [47], and so on. Recently, BEiT [6] proposes to encode image patches as a dictionary using dVAE [51] and predict the encoding of missing patches. PeCo [24] further improves BEiT by enforcing perceptual similarity from dVAE.

MAE [30] reconstructs directly the missing pixels. CiM [26] replaces image patches with plausible alternatives and learns to recover the original and predict which patches are replaced. Data2vec [5] performs self-supervised learning across multi-modalities. MultiMAE [3] shows that multi-modality pretraining can be more training-efficient than single-modality. Unlike prior methods, LoMaR achieves linear complexity by restricting self-attention to localized patches, rather than relying on global reconstructions with quadratic complexity.

### 3. Approach

LoMaR relies on a stack of Transformer [58] blocks to pre-train on a large number of unlabeled images by recovering the missing patches from corrupted images, similar to MAE [30]. Still, LoMaR differentiates from MAE in several key places. In this section, we first revisit the MAE model and then describe the differences between LoMaR and MAE.

#### 3.1. Background: Masked Autoencoder

The Masked Autoencoder (MAE) model [30], employs an asymmetric encoder-decoder architecture. The encoder takes in a subset of patches from an image and outputs latent representations for the patches. From those, the decoder reconstructs the missing patches. For an input image with resolution  $h \times w$ , MAE first divides it into a sequence of non-overlapping patches. Then, MAE randomly masks a large proportion (e.g., 75%) of image patches; see the upper side

of Fig. 2. The positional encodings are added to each patch to indicate their spatial location. MAE first encodes the remaining patches into the latent representation space. Then, it feeds the latent representations with placeholders for the masked patches into the decoder, which carries out the reconstruction. For each reconstructed image, MAE uses the mean squared error (MSE) with the original image in the pixel space as the loss function.

### 3.2. Local Masked Reconstruction (LoMaR)

We describe LoMaR by contrasting it with MAE from the following perspectives.

**Local vs. Global Masked Reconstruction.** MAE reconstructs each missing patch with patches sampled from the entire image. However, as indicated by Fig 1, usually only the patches in the proximity of the target patch contribute significantly to the reconstruction, suggesting that local information is sufficient for reconstruction. Therefore, we perform the random window masking and reconstruction on patches within a small region, shown in the bottom side of Fig. 2. Specifically, we perform the random window masking by sampling several small regions from the image and restricting the masked patches to only attend to their local surrounding visible patches, as we highlighted. Our experiments find that a region size of  $7 \times 7$  patches leads to the best trade-off between accuracy and efficiency. On the other hand, similar to convolutional networks [32, 54], LoMaR has the translation invariance property due to the usage of small windows sampled in random spatial locations each iteration.

From the complexity perspective, local masking and reconstruction are more computationally efficient than global masking and reconstruction of MAE because there are fewer tokens for operation. Suppose each image can be divided into  $h \times w$  patches. The time complexity for computing the self-attention is  $\mathcal{O}(h^2w^2)$ . The complexity is quadratic to the number of patches and hard to scale up with large  $hw$ . However, for our local masked reconstruction, we sample  $n$  windows where each contains  $m \times m$  patches; Its computational complexity is  $\mathcal{O}(nm^4)$ , which has linear time complexity if we fix  $m \times m$  as a constant window size. It can reduce the computational cost significantly if  $nm^4 \ll h^2w^2$ . For example, for a  $448 \times 448$  image, the cost of self-attention calculation is reduced from  $448^2 \times 448^2$  in the case of MAE to  $4 \times 7^2 \times 7^2$  in the case of LoMaR when we sample 4 views of  $7 \times 7$  patches.

**Architecture.** Instead of the asymmetric encoder-decoder of MAE, LoMaR only applies a simple Transformer encoder architecture. We input all the visible and masked patches under a sampled region into the encoder and reconstruct the masked patches through a simple MLP layer. Although feeding the masked patches into the encoder can be deemed a less efficient operation than MAE, which only in-

puts masked patches into the decoder, we find that inputting the masks in the early stages can enhance the visual representation and make it more robust to mask reconstruction from the smaller regions. It might be because the encoder can convert the masked patches back to their original RGB representation after multiple encoder layers interact with the other visible patches. Those recovered masks in the hidden layers can implicitly contribute to the image representation. Therefore, LoMaR preserves the mask patches as the encoder input.

**Implementation.** Given an image, we first divide it into several non-overlapping patches. Each patch is linearly projected into an embedding. We randomly sample several square-shaped regions of  $K \times K$  patches at different spatial locations. We then zeroed out a fixed percentage of patches within each region. After that, we feed all the patches, including visible and masked ones, from each region to the encoder in raster order. We also apply the relative positional encoding [60] into our model, enabling the translation-invariant property for the local masked reconstruction. We convert the latent representations from the encoder output to their original feature dimension with a simple MLP head and then compute the mean squared error with the normalized ground-truth image.

## 4. Experiments

We examine the performance of LoMaR by pretraining and finetuning on ImageNet-1K [22] dataset with the following procedure. First, we perform the self-supervised pretraining on the ImageNet-1K training dataset without label information. Then, we finetune the pre-trained model on ImageNet-1K with supervision from the labels. During finetuning, we feed all the image patches to the model and take the average of their features as the final representation for classification. We follow the same experimental settings as MAE [30]; detailed hyperparameters can be found in the supplementary material.

### 4.1. Experiments on High-resolution Images

We evaluate our model, MAE [30] and BEiT [6] on ImageNet [22] and Inaturalist [57] datasets. We pre-train and finetune on high-resolution images such as  $384 \times 384$  and  $448 \times 448$  images. We follow MAE’s default settings during pretraining; Sample 75% patches as masks. For LoMaR, we set the number of views to 6 and 9 for resolutions of 384 and 448 on the ImageNet dataset and sample 8 and 12 views for resolutions of 384 and 448 on the iNaturalist dataset. We pretrain all the models with 300 epochs and finetune them under the same image resolution.

We summarize the results in Table 1. The results demonstrate that LoMaR outperforms other models with substantially less pretraining time, which scales linearly with the window numbers. In contrast, the pretraining time of MAE

Method	Resolution	ImageNet 1K			Inaturalist		
		Time (h) ↓	Top-1 Acc	Speed-up	Time (h) ↓	Top-1 Acc	Speed-up
BEiT [6]	384	~408	83.7	1.0×	~93	81.1	1.0×
MAE	384	~203	84.3	2×	~46	81.2	2.0×
LoMaR (ours)	384	~ <b>81</b>	<b>84.5</b>	<b>5.0×</b>	~ <b>22</b>	<b>81.4</b>	<b>4.2×</b>
BEiT [6]	448	~595	84.1	1.0×	~121	82.3	1.0×
MAE	448	~345	84.5	1.7×	~70	<b>82.4</b>	1.7×
LoMaR (ours)	448	~ <b>113</b>	<b>84.7</b>	<b>5.3×</b>	~ <b>32</b>	82.3	<b>3.8×</b>

Table 1. High-resolution image pretraining and classification results on ImageNet-1K dataset [22] and Inaturalist [57]. The pretraining times are all computed on 4 NVIDIA 80GB A100 GPUs. We take BEiT [6] as the comparison baseline when computing the speed-up for MAE [30] and LoMaR. Our LoMaR can always achieve comparable or higher performance with at least a 3.8× speed-up than BEiT and a 2.2× speed-up than MAE on both datasets.

Methods	Epochs	Res	Time (h) ↓	Top-1 Acc	Speed-up
No Pretraining	-	224	-	82.3	-
DINO [13]	300	224	-	82.8	-
MoCov3 [20]	600	224	-	83.2	-
MSN [1]	600	224	-	83.4	-
CAE [17]	300	224	-	83.6	-
CAE [17]	800	224	-	83.8	-
CAE [17]	1600	224	-	83.9	-
MMAE [4]	1600	224	-	83.3	-
SemMAE [39]	800	224	-	83.3	-
BEiT [6] [51]	300	224	~107	82.9	1.0×
MAE* [30]	400	224	~58	83.1	1.8×
LoMaR	300	224	~49	83.3	2.2×
LoMaR <sub>8×8</sub>	400	224	~66	84.3	1.6×

Table 2. Image classification results on the ImageNet-1K (IN1K) dataset [22]. All baselines excluding LoMaR<sub>8×8</sub> adopt the ViT B/16 model [25] and are pretrained on 224×224 images. LoMaR<sub>8×8</sub> applies ViT B/8 as the backbone. \* denotes our reproduced results based on the officially released code and pretrained models. The pretraining times are all computed on 4 NVIDIA 80GB A100 GPUs.

and BEiT scales quadratically as the resolution increases. As a result, LoMaR is 2.5× faster than MAE (accuracy +0.2%) and 5.0× faster than BEiT (accuracy +0.8%) on 384×384 images, and for the resolution of 448×448, it is 3.1× faster than MAE (accuracy +0.2%) and 5.3× faster than BEiT (accuracy +0.6%). On Inaturalist, LoMaR is 2.1× faster than MAE model (accuracy +0.2%) and 4.2× faster than BEiT (accuracy +0.3%), and it can produce comparable performance with the other baselines but is 3.8× faster than BEiT and 2.2× faster than MAE.

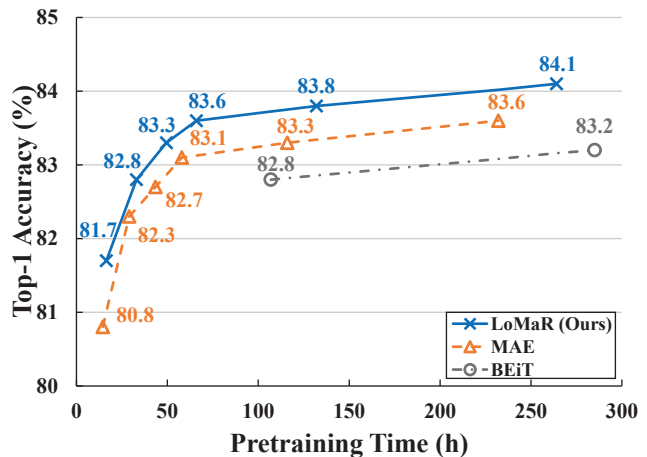


Figure 3. **Computational efficiency evaluation:** We compute their ImageNet-1K top-1 accuracy per pretraining time for low-resolution images 224×224.

## 4.2. Experiments on Low-resolution images

Table 2 summarizes the results of different self-supervised learning approaches. All models are pretrained in self-supervised fashion on ImageNet-1K [22] under the 224×224 resolution and finetuned on labeled ImageNet-1K. LoMaR reaches the best result of MAE, 83.6%, after only 400 epochs of pretraining. When pretrained for 1,600 epochs, its performance further improves to 84.1%; this can also be seen in Fig. 3 after ≈250 hours of pretraining. When finetuned under the 384×384 resolution, LoMaR reaches an accuracy of 85.4%, 0.6% higher than the best baseline. Overall, LoMaR outperforms strong baselines with less pretraining time.

**Efficiency analysis.** We train LoMaR, MAE [30] and BEiT [6] baselines with different pretraining epochs [100, 200, 300, 400, 800, 1,600] on 224×224 images. We compare their pretraining time v.s top-1 accuracy in Fig. 3. We

Backbone	ViTDet* [40]		ViTAE [65]	
	AP <sup>box</sup>	AP <sup>mask</sup>	AP <sup>box</sup>	AP <sup>mask</sup>
MAE	51.1*	45.4*	51.6	45.8
LoMaR	51.4	45.7	51.8	46.0
LoMaR <sub>384</sub>	<b>51.6</b>	<b>45.9</b>	<b>52.0</b>	<b>46.2</b>

Table 3. **Object detection and instance segmentation results** on MS COCO under two ViT frameworks. \* denotes reproduced results with the code from [65]. LoMaR<sub>384</sub> denotes the model pretrained on 384×384 images. Other models are pretrained on 224×224 images.

carefully tuned all models to achieve the best load balancing between the GPU and the CPU and the maximal image throughput during training. We do this by adjusting ghost batch size [35] while keeping the total batch size constant for all models. Compared to baselines, we observe that LoMaR consistently achieves the same or higher accuracy in less pretraining time. Specifically, pretraining MAE achieves 83.6% accuracy but takes about 232 hours. LoMaR reaches the same accuracy within ~66 hours of pretraining, 3.5× faster. BEiT requires 285 pretraining hours to get 83.2% accuracy. In contrast, LoMaR obtains a similar result within ~49 hours, translating to 5.8× time savings.

CAE [17], MMAE [4], and SEMMAE [39] build on top of MAE [30] architecture adding additional components. At best, the time complexity of these approaches is similar to that of MAE, so for high-resolution pre-training, LoMaR is much more efficient. MSN [1] differs from MAE as it matches the representation of an image with randomly masked patches to the representation of the original unmasked image. We calculate the time per epoch for MSN at a high resolution of 384 and 448 and compare it with LoMaR; for 384, MSN takes 0.64 hours/epoch, while for LoMaR, it takes 0.27 hours/epoch. For 448, MSN takes 0.88 hours/epoch, while LoMaR takes 0.38 hours/epoch. This comparison also demonstrates the computational efficiency of our LoMaR model.

**Pretraining on small patches.** We also evaluate our model on smaller patches with 8×8 pixels instead of the usual 16×16 pixels. We employ ViT B/8 [25] as a backbone in Table 2. We pre-train LoMaR with 7×7 windows (4 views per image) on 224×224 images for 400 epochs. It is worth noting that this incurs the same computation time (about 66 hours) as 16×16 patches. The model accuracy after finetuning reaches to 84.3% top-1 accuracy. However, similar experiments are costly for MAE and BEiT, as smaller patches substantially increase the number of patches for operation and lead to the high cost of self-attention. In our experiments, the pretraining of MAE with their official code under smaller patches crashes due to nu-

Models	Pre-train Data	ADE20K
supervised	IN1K w/ labels	47.4
BEiT [6]	IN1K+DALLE	47.1
MAE [30]	IN1K	48.1
LoMaR	IN1K	47.8
LoMaR <sub>384</sub>	IN1K	<b>48.3</b>

Table 4. **Semantic segmentation on ADE20K [67] (mIoU).** All the baselines are computed under UperNet [62] framework. LoMaR<sub>384</sub> denotes the results under pretraining images with the resolution of 384×384. Other baselines are pretrained on 224×224 images.

Method	Time(h)↓	Top-1 Acc	Speed-up
BEiT	~285	83.2	1×
BEiT+ window masking	~102	<b>83.4</b>	2.8×

Table 5. The results of applying our method on the BEiT approach.

merical issues.

### 4.3. Object Detection and Instance Segmentation.

We finetune our model end-to-end on MS COCO [41] for the object detection and instance segmentation tasks. We replace the ViT backbone with our pretrained LoMaR model in the ViTDet [40] and ViTAE [65] frameworks. We report object detection results in AP<sup>box</sup> and instance segmentation results in AP<sup>mask</sup>.

We provide the results in Table 3. It shows the consistent improvement of LoMaR on the COCO object detection benchmark. Under ViTDet, LoMaR surpasses MAE by 0.3 AP<sup>box</sup> and 0.3 AP<sup>mask</sup>. When applying the LoMaR pretrained for 1,600 epochs under the 384×384 resolution, it further improves to 51.6 AP<sup>box</sup> and 45.9 AP<sup>mask</sup>. In the ViTAE framework, LoMaR improves over MAE by 0.4 AP<sup>box</sup> and 0.4 AP<sup>mask</sup>, respectively.

### 4.4. Semantic Segmentation

We evaluate our model on the semantic segmentation benchmark, ADE20K [67], and compare with the baselines in Table 4. We train the UperNet [62] model with our pretrained LoMaR as initialization. When applying the LoMaR pretrained on images with 384×384 resolution, it consistently outperforms MAE by 0.2 points. This shows the consistent improvement of our LoMaR over the MAE and BEiT baselines. Additionally, it demonstrates the usefulness of high-resolution image pretraining.

### 4.5. Integration to BEiT

Our core idea, local masked reconstruction, can be easily integrated into other generative self-supervised learn-

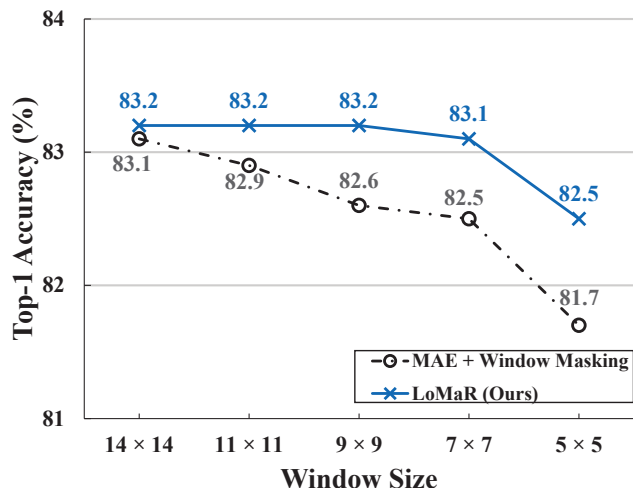


Figure 4. Comparison between LoMaR simple encoder and MAE asymmetric encoder-decoder architectures on our random window masking strategy. The window sizes vary from  $14 \times 14$  to  $5 \times 5$ .

ing methods. To examine its effectiveness in a different paradigm, we integrate it to BEiT [6]. Specifically, we randomly sample four  $7 \times 7$  windows, feed them into the BEiT model, and pretrain for 300 epochs while retaining all other experimental settings as the original BEiT. Results in Table 5 show that this strategy improves the accuracy from 82.8% to 83.4%, which is higher than the original BEiT, and speeds up the training by  $2.8 \times$ .

#### 4.6. Ablative Experiments

We conducted many ablation experiments to explore properties such as the window size, masking ratio, and architecture design and share our findings in this section. We performed all the ablation experiments under 4 NVIDIA 80GB A100 GPUs with the same setting for fair comparisons, and all the experiments were obtained by pretraining on  $224 \times 224$  images.

**Architecture.** Fig. 4 compares different architectures, including a simple encoder (with both visible and masked patches as input) and MAE, an asymmetric encoder-decoder architecture with a local window. Initially, we sample 75% patches as the masks following the guidance of MAE. By default, we use absolute positional encoding (APE) for both architectures. We ablate these two architectures with different masked reconstruction windows, showing that a simple encoder can consistently outperform the asymmetric encoder-decoder. Moreover, the performance gap is further magnified when we decrease the window size from 14 to 7. This suggests a simple encoder is more robust to smaller window sizes than MAE-like architecture.

**Efficiency vs. window size.** We test LoMaR with multiple different window sizes such as  $5 \times 5$ ,  $7 \times 7$ ,  $9 \times 9$ ,  $11 \times 11$

Window Size	$5 \times 5$	$7 \times 7$	$9 \times 9$	$11 \times 11$	$14 \times 14$
Views	8	4	3	2	1

Table 6. **Window size utilized** : The number of views per image, as utilized by LoMaR for different window sizes.

and  $14 \times 14$ . One caveat is that the smaller window covers much fewer visible patches than the larger ones, which creates unfair comparisons. To encourage fairness, we assign different numbers of views for each window size, as we demonstrated in Table 6. Thereby, all conditions have a similar number of visible patches in training.

From the results in Fig. 4, we can observe that while there is a performance drop when decreasing the window size from  $14 \rightarrow 5$ , the performance does not change much for other sizes, only  $83.2 \rightarrow 83.1$ , when decreasing the region size from  $14 \times$  to  $7$ . However, the restricted attention region decreases the total pretraining time from 120 to 66 hours. This means that pretraining on  $7 \times 7$  window size can roughly  $2 \times$  speed up the pretraining process with minimal performance change. Therefore, window size  $7 \times 7$  can be deemed an optimal trade-off for local masked reconstruction.

**RPE vs. APE.** Relative positional encoding (RPE) has been widely used in previous works, including BEiT [6]. We also employ the RPE [60] in LoMaR. We observe that it can bring 0.4 top-1 accuracy gain from 83.1 to 83.5. Therefore, we set RPE as the default setting for LoMaR in future experiments.

**Mask ratio.** We also explore the best mask ratio under the local masked reconstruction scenario (see Fig. 6). We train the previous best setting of our LoMaR on different mask ratios, ranging from 30% to 90%. The results show that too low (30%) or too high (90%) mask ratios are not optimal since they over-simplify or complicate the training task. We found that the 80% mask ratio can result in the best performance, differentiating from the 60% mask ratio observed in MAE for best finetuning performance. With this motivation, we employ the 80% mask ratio in the rest of our experiments.

#### 4.7. Visualization of Reconstructed Images

We qualitatively show the reconstruction performance of our pre-trained model in Fig. 5. We randomly sample several images from ImageNet-1K [22] and MS COCO [41]. After that, we sample a region containing  $7 \times 7$  patches in every image and zero out 80% patches in the window for reconstruction. It can be seen that LoMaR is capable of generating plausible images, which also confirms our initial conjecture that the missing patches can be recovered from the local surrounding patches alone.

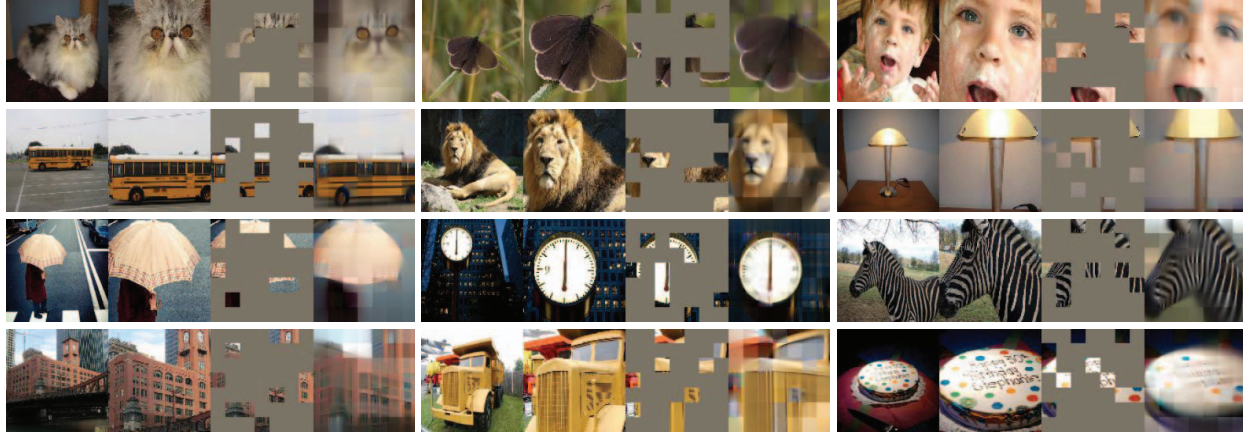


Figure 5. Example results on ImageNet (upper two rows) and COCO (lower two rows) validation images. We mask 80% patches out and reconstruct them with our pretrained model. For each image reconstruction figure, we split them into 4 parts: 1) the left-most is the original image. 2) the second-left is the sampled window ( $7 \times 7$  patches). 3) The second-right is the masked image. 4) The right-most is our reconstructed image.

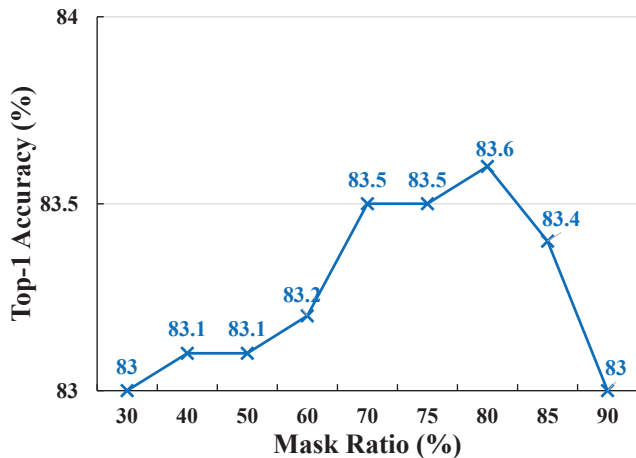


Figure 6. **Mask ratio ablations:** It compares the LoMaR under different mask ratios from 30% to 90%

## 5. Discussion and Limitation

Self-supervised learning (SSL) can benefit from training with a massive amount of unlabeled data, which has brought many promising results [6, 11, 20, 30, 36, 49, 50]. However, their high computational demands remain a significant concern under large-scale pretraining. In our study, we observe that the local masked reconstruction (LoMaR) for generative SSL is more efficient than the global version used by the influential works of MAE [30] and BEiT [6]. LoMaR demonstrates good generalization in image classification, instance segmentation, and object detection; it can be easily incorporated into both MAE and BEiT. LoMaR holds the promise to scale up SSL to even bigger datasets and higher resolution [48, 56] datasets. LoMaR can also be extended

to video analysis, where the computation problem is even more severe.

Another advantage of LoMaR is its efficiency gains when the number of image patches increases in high-resolution images such as  $384 \times 384$  and  $448 \times 448$  or even larger. The primary reason is that LoMaR restricts self-attention within a small region, and its computational complexity grows linearly with the number of sampled regions per image. This characteristic enables efficient pretraining under high image resolution, which would be prohibitively expensive for other SSL methods. It can benefit many vision tasks such as object detection or instance segmentation, which require dense prediction at the pixel level.

Despite the high pretraining efficiency gain of LoMaR over other baselines for high-resolution images, one limitation is that LoMaR underperforms in linear probing (see results in Supplementary), which is mainly due to two reasons: 1) There is a discrepancy between training and inference. During pretraining, we feed only a small region of patches and masked tokens to the network. The input contains all image patches without masked tokens during linear probing, resulting in a shift of input distribution and damaging linear probing performance. 2) LoMaR applies a much shallower decoder than MAE. A deep decoder improves linear probing performance because the last few layers in an autoencoder are specialized for reconstruction and not very helpful for recognition; MAE removes these layers during linear probing. However, as shown in Table 2, fine-tuning the entire model can easily mitigate this limitation. We hope the local masked reconstruction idea, pioneered by LoMaR, can lead to further research on efficient self-supervised learning.

## References

- [1] Mahmoud Assran, Mathilde Caron, Ishan Misra, Piotr Bojanowski, Florian Bordes, Pascal Vincent, Armand Joulin, Michael Rabbat, and Nicolas Ballas. Masked siamese networks for label-efficient learning, 2022. 5, 6
- [2] Philip Bachman, R Devon Hjelm, and William Buchwalter. Learning representations by maximizing mutual information across views. *Advances in neural information processing systems*, 32, 2019. 1
- [3] Roman Bachmann, David Mizrahi, Andrei Atanov, and Amir Zamir. Multimaec: Multi-modal multi-task masked autoencoders, 2022. 3
- [4] Roman Bachmann, David Mizrahi, Andrei Atanov, and Amir Zamir. Multimaec: Multi-modal multi-task masked autoencoders, 2022. 5, 6
- [5] Alexei Baevski, Wei-Ning Hsu, Qiantong Xu, Arun Babu, Jiatao Gu, and Michael Auli. data2vec: A general framework for self-supervised learning in speech, vision and language. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 1298–1312. PMLR, 17–23 Jul 2022. 3
- [6] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. BEit: BERT pre-training of image transformers. In *International Conference on Learning Representations*, 2022. 1, 3, 4, 5, 6, 7, 8
- [7] Adrien Bardes, Jean Ponce, and Yann LeCun. Vicregl: Self-supervised learning of local visual features. *arXiv preprint arXiv:2210.01571*, 2022. 3
- [8] Suzanna Becker and Geoffrey E Hinton. Self-organizing neural network that discovers surfaces in random-dot stereograms. *Nature*, 355(6356):161–163, 1992. 2
- [9] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013. 3
- [10] Yoshua Bengio, Pascal Lamblin, Dan Popovici, and Hugo Larochelle. Greedy layer-wise training of deep networks. In *Proceedings of the 19th International Conference on Neural Information Processing Systems*, NeurIPS’06, page 153–160. Cambridge, MA, USA, 2006. MIT Press. 3
- [11] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020. 8
- [12] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in Neural Information Processing Systems*, 33:9912–9924, 2020. 3
- [13] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9650–9660, 2021. 1, 3, 5
- [14] Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. Generative pre-training from pixels. In *International Conference on Machine Learning*, pages 1691–1703. PMLR, 2020. 1, 3
- [15] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020. 3
- [16] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey E Hinton. Big self-supervised models are strong semi-supervised learners. *Advances in neural information processing systems*, 33:22243–22255, 2020. 3
- [17] X. Chen, M. Ding, and X. et al. Wang. Context autoencoder for self-supervised representation learning. *International Journal of Computer Vision*, 132(2):208–223, 2024. 5, 6
- [18] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020. 3
- [19] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15750–15758, 2021. 3
- [20] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9640–9649, 2021. 1, 3, 5, 8
- [21] Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers. *URL <https://openai.com/blog/sparse-transformers>*, 2019. 2
- [22] Jia Deng, Wei Dong, R Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, 2009. 2, 4, 5, 7
- [23] Carl Doersch, Abhinav Gupta, and Alexei A Efros. Unsupervised visual representation learning by context prediction. In *Proceedings of the IEEE international conference on computer vision*, pages 1422–1430, 2015. 3
- [24] Xiaoyi Dong, Jianmin Bao, Ting Zhang, Dongdong Chen, Weiming Zhang, Lu Yuan, Dong Chen, Fang Wen, and Nenghai Yu. Peco: Perceptual codebook for bert pre-training of vision transformers. *arXiv preprint arXiv:2111.12710*, 2021. 3
- [25] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 5, 6
- [26] Yuxin Fang, Li Dong, Hangbo Bao, Xinggang Wang, and Furu Wei. Corrupted image modeling for self-supervised visual pre-training. *arXiv preprint arXiv:2202.03382*, 2022. 3
- [27] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations. In *International Conference on Learning Representations*, 2018. 2

- [28] Jean-Bastien Grill, Florian Strub, Florent Althé, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent—a new approach to self-supervised learning. *Advances in Neural Information Processing Systems*, 33:21271–21284, 2020. 3
- [29] Tengda Han, Weidi Xie, and Andrew Zisserman. Video representation learning by dense predictive coding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 0–0, 2019. 3
- [30] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. *arXiv preprint arXiv:2111.06377*, 2021. 1, 2, 3, 4, 5, 6, 8
- [31] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020. 3
- [32] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 4
- [33] G. E. Hinton and R. R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504–507, 2006. 3
- [34] R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman, Adam Trischler, and Yoshua Bengio. Learning deep representations by mutual information estimation and maximization. In *International Conference on Learning Representations*, 2018. 1
- [35] Elad Hoffer, Itay Hubara, and Daniel Soudry. Train longer, generalize better: Closing the generalization gap in large batch training of neural networks. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*, page 1729–1739, Red Hook, NY, USA, 2017. Curran Associates Inc. 6
- [36] Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, pages 4171–4186, 2019. 8
- [37] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 3
- [38] Gustav Larsson, Michael Maire, and Gregory Shakhnarovich. Learning representations for automatic colorization. In *European conference on computer vision*, pages 577–593. Springer, 2016. 3
- [39] Gang Li, Heliang Zheng, Daqing Liu, Chaoyue Wang, Bing Su, and Changwen Zheng. Semmae: Semantic-guided masking for learning masked autoencoders. *arXiv preprint arXiv:2206.10207*, 2022. 5, 6
- [40] Yanghao Li, Hanzi Mao, Ross Girshick, and Kaiming He. Exploring plain vision transformer backbones for object detection. *arXiv preprint arXiv:2203.16527*, 2022. 2, 6
- [41] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 6, 7
- [42] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021. 2
- [43] Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. Adversarial autoencoders. *arXiv preprint arXiv:1511.05644*, 2015. 3
- [44] Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *European conference on computer vision*, pages 69–84. Springer, 2016. 3
- [45] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 1, 3
- [46] Deepak Pathak, Ross Girshick, Piotr Dollár, Trevor Darrell, and Bharath Hariharan. Learning features by watching objects move. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2701–2710, 2017. 2
- [47] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2536–2544, 2016. 3
- [48] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. 8
- [49] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. *Technical report, OpenAI*, 2018. 8
- [50] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019. 8
- [51] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021. 3, 5
- [52] Marcaurelio Ranzato, Christopher Paultney, Sumit Chopra, and Yann LeCun. Efficient learning of sparse representations with an energy-based model. In *Advances in Neural Information Processing Systems*, volume 19, 2006. 3
- [53] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In Eric P. Xing and Tony Jebara, editors, *Proceedings of the 31st International Conference on Machine Learning Research*, pages 1278–1286, Beijing, China, 22–24 Jun 2014. PMLR. 3
- [54] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 4

- [55] Sainbayar Sukhbaatar, Édouard Grave, Piotr Bojanowski, and Armand Joulin. Adaptive attention span in transformers. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 331–335, 2019. 2
- [56] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable effectiveness of data in deep learning era. In *Proceedings of the IEEE international conference on computer vision*, pages 843–852, 2017. 8
- [57] Grant Van Horn, Oisín Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. The inaturalist species classification and detection dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8769–8778, 2018. 4, 5
- [58] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 3
- [59] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th International Conference on Machine Learning, ICML '08*, page 1096–1103, New York, NY, USA, 2008. Association for Computing Machinery. 3
- [60] Kan Wu, Houwen Peng, Minghao Chen, Jianlong Fu, and Hongyang Chao. Rethinking and improving relative position encoding for vision transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10033–10041, 2021. 4, 7
- [61] Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3733–3742, 2018. 1, 3
- [62] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 418–434, 2018. 2, 6
- [63] Jianwei Yang, Chunyuan Li, Pengchuan Zhang, Xiyang Dai, Bin Xiao, Lu Yuan, and Jianfeng Gao. Focal self-attention for local-global interactions in vision transformers. *arXiv preprint arXiv:2107.00641*, 2021. 2
- [64] Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. Big bird: Transformers for longer sequences. *Advances in Neural Information Processing Systems*, 33:17283–17297, 2020. 2
- [65] Qiming Zhang, Yufei Xu, Jing Zhang, and Dacheng Tao. Vitaev2: Vision transformer advanced by exploring inductive bias for image recognition and beyond. *arXiv preprint arXiv:2202.10108*, 2022. 6
- [66] Richard Zhang, Phillip Isola, and Alexei A Efros. Colorful image colorization. In *European conference on computer vision*, pages 649–666. Springer, 2016. 3
- [67] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ade20k dataset. *International Journal of Computer Vision*, 127(3):302–321, 2019. 6