# Masked Image Modeling with Local Multi-Scale Reconstruction

Haoqing Wang[1], Yehui Tang[1,2], Yunhe Wang[2]*, Jianyuan Guo[2], Zhi-Hong Deng[1]*, Kai Han[2]
[1]National Key Lab of General AI, School of Intelligence Science and Technology, Peking University
[2]Huawei Noah's Ark Lab
{wanghaoqing,yhtang,zhdeng}@pku.edu.cn, {kai.han, yunhe.wang}@huawei.com

## Abstract

*Masked Image Modeling (MIM) achieves outstanding success in self-supervised representation learning. Unfortunately, MIM models typically have huge computational burden and slow learning process, which is an inevitable obstacle for their industrial applications. Although the lower layers play the key role in MIM, existing MIM models conduct reconstruction task only at the top layer of encoder. The lower layers are not explicitly guided and the interaction among their patches is only used for calculating new activations. Considering the reconstruction task requires non-trivial inter-patch interactions to reason target signals, we apply it to multiple local layers including lower and upper layers. Further, since the multiple layers expect to learn the information of different scales, we design local multi-scale reconstruction, where the lower and upper layers reconstruct fine-scale and coarse-scale supervision signals respectively. This design not only accelerates the representation learning process by explicitly guiding multiple layers, but also facilitates multi-scale semantical understanding to the input. Extensive experiments show that with significantly less pre-training burden, our model achieves comparable or better performance on classification, detection and segmentation tasks than existing MIM models. Code is available with both MindSpore and PyTorch.*

## 1. Introduction

Recently, Masked Image Modeling (MIM) [2, 24, 60] achieves outstanding success in the field of self-supervised visual representation learning, which is inspired by the Masked Language Modeling (MLM) [4, 34] in natural language processing and benefits from the development of vision transformers [18, 39, 55]. MIM learns semantic representations by first masking some parts of the input and then predicting their signals based on the unmasked parts, e.g., normalized pixels [24, 60], discrete tokens [2, 17], HOG fea-

ture [57], deep features [1, 67] or frequencies [38, 59].

Despite superior performance on various downstream tasks, these models have huge computational burden and slow learning process [31]. They typically require thousands of GPU Hours for pre-training on ImageNet-1K to get generalizing representations. Since we expect to pre-train these models on more massive amount of unlabeled data (e.g., free Internet data) to obtain more generalizing representations in practice, the pre-training efficiency is an inevitable bottleneck limiting the industrial applications of MIM. How to accelerate the representation learning in MIM is an important topic. To this end, MAE [24] pioneered the asymmetric encoder-decoder strategy, where the costly encoder only operates few visible patches and the lightweight decoder takes all the patches as input for prediction. Further, GreenMIM [31] extends the asymmetric encoder-decoder strategy to hierarchical vision transformers (e.g., Swin [39]). Besides, [8, 22, 35] shrinks the input resolution to lessen the input patches, thereby reducing the computational burden. However, they all aim to accelerate the encoding process rather than the representation learning.

In MIM, the learning of upper layers depends on that of lower ones during pre-training, since the upper-layer features are calculated from the lower layers. Besides, during fine-tuning the upper layers are typically tuned quickly to adapt to the downstream task while the lower ones change more slowly and need to be well-learned [2, 29, 65]. Even fine-tuning only the several upper layers and freezing the others can obtain similar performance [24]. Therefore, the lower layers of encoder play the key role in MIM. However, all existing MIM models only conduct reconstruction task at the top layer of encoder and the lower ones are not explicitly guide, thus the interaction among their patches is only used for calculating the activations of the next layer. Considering the reconstruction task requires non-trivial inter-patch interactions to reason target signals, we apply it to both lower and upper layers to explicitly guide them and thus accelerate the overall learning process. Using tiny decoder is sufficient for each local reconstruction task and does not significantly increase the computational burden.
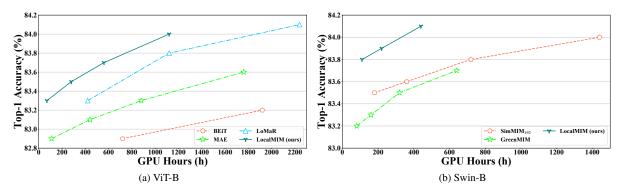
---

*Corresponding author.

Figure 1. Top-1 fine-tuning accuracy on ImageNet-1K vs. Pre-training duration. The duration is estimated on a machine with one Tesla V100-32G GPU, CUDA 10.2 and PyTorch 1.8. 'GPU Hours' is the running time on single GPU.

How to properly conduct reconstruction tasks at multiple local layers is a non-trivial problem. For example, applying the top-layer reconstruction task to carefully chosen local layers of ViT [18] can not achieve meaningful improvement. In general, the lower layers exploit low-level information and the upper ones learn high-level information [20, 44], so it is not appropriate to use the supervision signals of same scale for multiple local reconstruction tasks. Here 'scale' is the spatial size of the supervision signals calculated from the divided input regions, e.g., the signals from the $p \times p$ regions in an input of $H \times W$ resolution has the scale of $\frac{H}{p} \times \frac{W}{p}$. The fine-scale and coarse-scale supervisions typically contain low-level and high-level information of the input respectively, and these multi-scale supervisions from input are widely ignored by existing MIM models. To this end, we propose local multi-scale reconstruction where the lower and upper layers reconstruct fine-scale and coarse-scale supervisions respectively. This design not only accelerates the representation learning process, but also facilitates multi-scale semantic understanding to the input. When the decoded predictions have different scale with the supervisions (e.g., on ViT), we use the deconvolution/pool operations to rescale them. We also apply the asymmetric encoder-decoder strategy [24, 31] for quick encoding. Our model, dubbed as LocalMIM, are illustrated in Fig. 2 (a).

Overall, we summarize our contributions as follows.

- To the best of our knowledge, this is the first work in MIM to conduct local reconstructions and use multi-scale supervisions from the input.

- Our model is architecture-agnostic and can be used in both columnar and pyramid architectures.

- From extensive experiments, we find that 1) LocalMIM is more efficient than existing MIM models, as shown in Fig. 1 and Table 1. For example, LocalMIM achieves the best MAE result with $3.1\times$ acceleration on ViT-B and the best GreenMIM result with $6.4\times$ acceleration on Swin-B. 2) In terms of top-

1 fine-tuning accuracy on ImageNet-1K, LocalMIM achieves $84.0\%$ using ViT-B and $84.1\%$ using Swin-B with significantly less pre-training duration than existing MIM models. The obtained representations also achieve better generalization on detection and segmentation downstream tasks, as shown in Table 2 and 3.

## 2. Related Works

**Masked Image Modeling.** With the development of vision transformers [18, 23, 39, 50, 55], Masked Image Modeling (MIM) gradually replaces the dominant position of contrastive learning [10, 25, 54] in visual self-supervised representation learning due to its superior fine-tuning performance in various visual downstream tasks. Many target signals have been designed for the mask-prediction pretext task in MIM, such as normalized pixels [24, 60], discrete tokens [2, 17], HOG feature [57], deep features [1, 67] or frequencies [38, 59]. However, they are all only applied as single-scale supervisions for reconstruction. An inevitable bottleneck for the industrial applications of MIM is that these models typically require huge computational resources and long pre-training duration. To this end, some works accelerate the encoding process via the asymmetric encoder-decoder strategy [24, 31] or lessening the input patches [8, 35]. Only accelerating the encoding process sometimes doesn't really speed up the representation learning, like GreenMIM vs. SimMIM$_{192}$ in Fig. 1. In this work, we use local multi-scale reconstructions to explicitly guide multiple lower layers and thus accelerate the overall learning process. Our method is compatible with the above quick encoding approaches. ConvMAE [20] fuses the feature of local layers for the final reconstruction to explicitly guide them, but still applies single-scale supervision signals.

**Locally supervised learning.** Considering the biological brain learns mainly based on local information [6], some works [42, 43, 45, 56] used local error signals to greedily optimize individual blocks of the backbone without global back-propagation. This greedy training procedure signif-
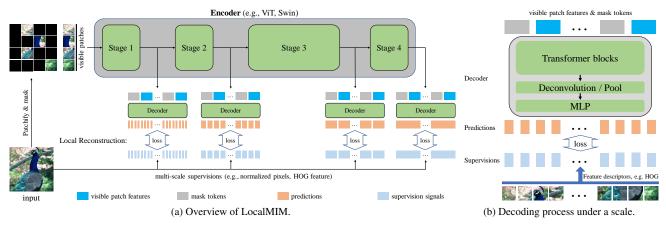
**Figure 2. Illustration of local multi-scale reconstruction.** a) We randomly mask some input patches and then predict their supervision signals of different scales at multiple local layers. The multi-scale supervisions are obtained by first dividing the input under different scales and then extracting signals with some feature descriptors. The lower layers are responsible for fine-scale reconstruction and the upper ones are responsible for coarse-scale reconstruction. We also use the asymmetric encoder-decoder strategy for quick encoding. b) The decoder consists of three parts: Transformer blocks, Deconvolution/Pool (optional) and Multilayer perceptron.

icantly reduces memory overhead, but can not accelerate model learning. In this work, we also optimize error at multiple local layers but still use global back-propagation for end-to-end training. Compared with existing MIM models, our local reconstruction is more biologically plausible since the human brain prefers local learning rules [3,14,16]. Feature distillation [28,47] also explicitly guides the local layers, but needs costly pre-trained or momentum teacher. In fact, the multi-scale supervisions from original input are sufficient to guide local layers and also readily available.

**Multi-scale property.** Biological visual perception is hierarchical and multi-scale [33]. Moreover, multi-scale features are also useful for many visual tasks [9,26,46]. To this end, many advanced vision transformers [19,39,53,55,64] hard-code the multi-scale property to the architectures and boost the performance. However, global single-scale reconstruction is not enough for guiding multiple local layers to learn multi-scale information. In this work, beyond the multi-scale features, we further introduce multi-scale supervisions to soft-code this property. This new design is not limited to specific architectures.

## 3. Model

### 3.1. Masked Image Modeling

In MIM, the models first mask some parts of the input image and then predict their information based on the unmasked observation. For quick encoding, MAE [24] and GreenMIM [31] use the asymmetric encoder-decoder strategy. The major designs are described below.

**Image presentation.** Generally, we present an image to the sequence of visual patches as the input to vision transformers [18,39]. The input image $x \in \mathbb{R}^{H \times W \times C}$ is reshaped

to $N = HW/p^2$ non-overlapping patches $x^p \in \mathbb{R}^{N \times (p^2 C)}$, where $p$ is the patch size, $(H, W)$ is the resolution of input image and $C$ is the number of channels. The patches $\{x_i^p\}_{i=1}^N$ are then linearly mapped to the patch embeddings. For retaining the positional information, the patches are typically added with positional embeddings.

**Vision transformers.** Unlike the grid operations in the convolutional neural networks, vision transformers [18,39,55] use the stacked multi-head self-attention modules. To obtain the multi-scale feature maps for dense prediction tasks, some advanced vision transformers [39,55] extend the columnar ViT [18] to the pyramidal structure, where the feature maps change from fine-scale to coarse-scale and are expected to capture the multi-scale information from the input. The patches $[x_1^p, \cdots, x_N^p]$ are fed to a vision transformer with $L$ layers and obtain the output feature at $l$-th layer $z^l = [z_1^l, \cdots, z_{N_l}^l], l = 1, \cdots, L$.

**Masking.** Given the patch sequence $\{x_i^p\}_{i=1}^N$, MIM constructs a random mask $m \in \{0, 1\}^N$ to indicate the masked patches that correspond to $m_i = 1$. There are two main masking strategies, *random masking* [24] and *block-wise masking* [2]. Specially, since the number of patches $N_l$ decreases through layers in the pyramidal architectures, some works [20,31] first construct the random mask $m^L$ with the size of $N_L$, and then up-sample it to size $N$.

**Encoder.** Only the visible patches $x^v = \{x_i^p | m_i = 0\}$ are fed to the encoder and mapped to the latent features $z_v^l, l = 1, \cdots, L$, which significantly reduces compute and memory. Specially, since the local window attention in the pyramidal architectures is not compatible with the incomplete input patches (i.e., only visible patches), GreenMIM [31] proposes the Optimal Grouping algorithm and the Group Window Attention scheme. After pre-training,

the encoder is used in various downstream tasks.

**Decoder.** The decoder takes as the input of both encoded visible patches $z_v^L$ and mask tokens $\{e_{[\mathbb{M}]}|m_i^L = 1\}$, where the mask token $e_{[\mathbb{M}]}$ is a shared and learnable vector. The positional embeddings are also added to them for providing the location information. Since the decoder is only used during pre-training to output the prediction $\hat{y}^L$, it can be any architectures which support the global information propagation among patches, e.g., a series of ViT or Swin blocks. Many works [24, 31, 35] show that using lightweight decoders is sufficient to learn generalizing representations.

**Global reconstruction.** All existing MIM models predict the supervision signal $y$ based on the final output feature $z_v^L$ of the encoder and minimize a global reconstruction loss

$$\mathcal{L}_{MIM} = -\sum_{i=1}^{N_L} m_i^L \cdot \ln P(y_i|\hat{y}_i^L) \qquad (1)$$

where the loss is calculated on the masked patches with $m_i^L = 1$, and $P(\cdot|\cdot)$ can be the Gaussian or Dirichlet distributions for regression or classification losses respectively.

### 3.2. Analysis and motivation

For pre-training, the upper-layer features are calculated from the lower layers, so the well-learned lower layers can propagate semantical knowledge to the upper ones and facilitate their learning. For fine-tuning, the upper layers typically adapt quickly to specific downstream tasks, while the lower ones change more slowly and need to be sufficient learned during pre-training [2,29,65]. Even fine-tuning only the several upper layers and freezing the others can obtain comparable performance [24]. Therefore, the lower layers of the encoder play the key role in MIM.

After patchification and linearly projection, the initial patch embeddings lose the inter-patch semantic relations. The self-attention mechanism in vision transformers is responsible for learning these relations by inter-patch interaction and build a better representation space than pixel space [5]. Further, since the self-attention mechanism has the computational complexity with a quadratic dependence on patch number $N$, it is difficult to learn the inter-patch interactions, especially for the lower layers of pyramidal architectures where the small patch size $p$ leads to huge $N$. However, under the global reconstruction loss, the inter-patch interaction at lower layers is not explicitly guided, and the simple task of calculating new activations is not sufficient to guide it. As a result, it is hard for the patches at lower layers to learn inter-patch relations. Concretely, we examine Normalized Mutual Information (NMI) [49] between query and key patches at each layer, where the high NMI value means the attention maps strongly depend on the query patches. Intuitively, the semantical representations
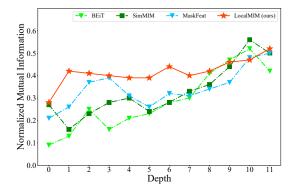


Figure 3. Normalized Mutual Information (NMI) between query and key patches at each layer of a pre-trained ViT-B.

should have highly query-adaptive attentions, i.e., the different query patches faithfully attend to their semantically related regions. This is an advantage of the self-attention mechanism. As shown in Fig. 3, existing MIM models with global loss have small NMI value at lower layers, which means their patches have less query-adaptive attentions.

The reconstruction task requires holistic reasoning among patches to predict the masked signals and thus obtains semantic understanding to the input. Since this challenging task facilitates non-trivial inter-patch interactions, we apply it at multiple local layers, including both lower and upper ones to explicitly guide them all.

### 3.3. Local multi-scale reconstruction

In MIM, the supervision signals for a reconstruction task are directly calculated from the input. Concretely, we evenly divide the input image $x \in \mathbb{R}^{H \times W \times C}$ into non-overlapping regions $\{x_i \in \mathbb{R}^{p \times p \times C}\}_{i=1}^{HW/p^2}$ and use some feature descriptor $\pi$ to extract the supervision signal $y_i = \pi(x_i)$. To learn generalizing representations, many feature descriptors have been designed, such as pixel normalization [24, 60], HOG [57] and the pre-trained or momentum teacher [1,2,17,67]. We define the scale of supervision $y$ as its spatial size $\frac{H}{p} \times \frac{W}{p}$. For a given input, the fine-scale supervisions from finely-divided input regions typically contain the low-level semantic information of the input, like corners, edges or textures. Relatively, the coarse-scale supervisions capture high-level semantical information of the input, like the shape of partial or whole object. Intuitively, multi-scale supervisions can guide representation learning better than the common single-scale ones due to their richer semantic information. In this work, we mainly consider the feature descriptors which are readily available without extra pre-training burden and costly forward inference of teacher networks, e.g., pixel normalization and HOG.

As shown in Table 4f, directly applying the top-layer reconstruction task at carefully chosen local layers of ViT can not achieve meaningful improvement, where each local task

| Model | Backbone | # Params | PT Epoch | GPU Hours/Ep. | Total GPU Hours | Acc |
|---|---|---|---|---|---|---|
| Scratch, ViT | ViT-B | 86M | 0 | 1.5 | - | 82.3 |
| Scratch, Swin | Swin-B | 88M | 0 | 2.4 | - | 83.5 |
| MoCo v3 [12] | ViT-B | 86M | 600 | - | - | 83.2 |
| DINO [7] | ViT-B | 86M | 300 | - | - | 82.8 |
| BEiT [2] | ViT-B | 86M | 800 | 2.4 | 1920 | 83.2 |
| iBOT [67] | ViT-B | 86M | 400 | 10.1 | 4040 | 83.8 |
| MAE [24] | ViT-B | 86M | 800 | 1.1 | 880 | 83.3 |
| MAE [24] | ViT-B | 86M | 1600 | 1.1 | 1760 | 83.6 |
| MAE [24] | ViT-L | 307M | 1600 | 1.7 | 2720 | 85.9 |
| MaskFeat [57] | ViT-B | 86M | 1600 | 3.9 | 6240 | 84.0 |
| CAE [11] | ViT-B | 86M | 800 | 2.8 | 2240 | 83.6 |
| LoMaR$^\dagger$ [8] | ViT-B | 86M | 1600 | 1.4 | 2240 | 84.1 |
| data2Vec$^\dagger$ [1] | ViT-B | 86M | 800 | 3.0 | 2400 | 84.2 |
| PeCo [17] | ViT-B | 86M | 800 | - | - | 84.5 |
| LocalMIM-HOG | ViT-B | 86M | 100 | 0.7 | 70 | 83.3 |
| LocalMIM-HOG | ViT-B | 86M | 1600 | 0.7 | 1120 | 84.0 |
| LocalMIM-HOG | ViT-L | 307M | 800 | 1.0 | 800 | 85.8 |
| SimMIM$_{192}$ [60] | Swin-B | 88M | 800 | 1.8 | 1440 | 84.0 |
| SimMIM$_{192}$ [60] | Swin-L | 197M | 800 | 3.0 | 2400 | 85.4 |
| GreenMIM [31] | Swin-B | 88M | 800 | 0.8 | 640 | 83.7 |
| GreenMIM [31] | Swin-L | 197M | 800 | 1.4 | 1120 | 85.1 |
| LocalMIM-Pixel | Swin-B | 88M | 100 | 1.0 | 100 | 83.7 |
| LocalMIM-HOG | Swin-B | 88M | 100 | 1.1 | 110 | 83.8 |
| LocalMIM-Pixel | Swin-B | 88M | 400 | 1.0 | 400 | 84.0 |
| LocalMIM-HOG | Swin-B | 88M | 400 | 1.1 | 440 | 84.1 |
| LocalMIM-HOG | Swin-L | 197M | 800 | 1.6 | 1280 | 85.6 |

Table 1. Top-1 fine-tuning accuracy on ImageNet-1K. All models are pre-trained and fine-tuned under $224 \times 224$ resolution except that SimMIM$_{192}$ uses $192 \times 192$ resolution for pre-training. $\dagger$ means using the relative positional encoding.

predicts the supervisions of same scale. In fact, the lower and upper layers expect to learning low-level and high-level information respectively [20, 44], so it is not appropriate to use the single-scale supervisions to guide multiple local layers, even for the columnar architectures with feature maps of the same scale at all layers. To this end, we make the lower layers to reconstruct fine-scale supervisions and the upper ones to reconstruct coarse-scale supervisions. Specially, for the pyramidal architectures which already hardcode the multi-scale property to the features by setting their spatial sizes, we use the supervisions with same scale as the feature maps at the chosen layers for compatibility.

The decoding process under a specific scale is illustrated in Fig. 2 (b). The decoder consists of three parts: Transformer blocks for reasoning, (optional) Deconvolution/Pool for rescaling and Multilayer perceptron for prediction. Concretely, based on the encoded visible patches $z_v^l$ from $l$-th layer and the mask tokens $\{e_{[\mathbb{M}]}^l | m_i^l = 1\}$, a decoder outputs prediction $\hat{y}^l$ that has the same scale as feature map $z_v^l$. When the supervision $y^l$ has different scale with the feature map $z_v^l$ (e.g., in the columnar architectures), the de-

coded prediction can not match the supervision. To this end, we use the deconvolution/pool operations to rescale the prediction $\hat{y}^l$ for matching the supervision $y^l$. For example, we can rescale the $14 \times 14$ prediction to the scale of $56 \times 56$ with twice deconvolution operations or to the scale of $7 \times 7$ with average pool. To avoid excessive computational overhead, we use tiny decoders containing one Transformer block with small embedding dimension.

The training loss is the weighted summation of the reconstruction losses at the chosen layers

$$\mathcal{L}_{LocalMIM} = -\sum_{l \in \mathcal{I}} w_l \cdot \sum_{i=1}^{N_l} m_i^l \cdot \ln P(y_i^l | \hat{y}_i^l) \quad (2)$$

where $\mathcal{I}$ is the set of chosen layers, $w_l$ is the coefficient of each local loss, and mask $m^l$ is calculated by up/downsampling the initial mask $m$. These local losses guide the patches at multiple chosen layers to conduct semantic interactions under different scales, which not only accelerates the learning of multiple layers but also facilitates multiscale semantical understanding to the input. As shown in Fig. 3, compared with existing models, our LocalMIM has

| Model | PT Epoch | PT Hours | mIoU |
|---|---|---|---|
| Supervised | - | - | 47.4 |
| MoCo v3 [12] | 300 | - | 47.3 |
| BEiT [2] | 800 | 1920 | 47.1 |
| MAE [24] | 1600 | 1760 | 48.1 |
| MaskFeat [57] | 1600 | 6240 | 48.8 |
| PeCo [17] | 800 | - | 48.5 |
| CAE [11] | 800 | 2240 | 48.8 |
| LocalMIM-HOG | 1600 | 1120 | 49.5 |

Table 2. Semantic segmentation on ADE20K using UperNet with ViT-B backbone. Our LocalMIM achieves better results than previous MIM models with less pre-training duration.

larger NMI values at lower layers, which means the attention maps depend more strongly on the query patches.

# 4. Experiments

In this section, we first evaluate our LocalMIM model on classification, detection and segmentation tasks, and then provide some ablation studies for deep understanding.

## 4.1. Classification on ImageNet-1K

**Settings.** Both pre-training and fine-tuning are conducted on ImageNet-1K [48] dataset under the $224 \times 224$ resolution. We mainly examine two representative architectures, columnar ViT [18] and pyramidal Swin [39]. The input images are patchified with patch size $p = 16$ for ViT and $p = 4$ for Swin, and the obtained patches are randomly masked with ratio $r = 0.75$ by default. We use both HOG feature [57] and normalized pixels [24] as the supervision signals. For Swin, we introduce the reconstruction task after each stage and the supervision signals have the same scale with the output feature maps, e.g., the chosen layers $\mathcal{I} = \{2, 4, 22, 24\}$ and the scales of supervision are $\{28^2, 14^2, 7^2, 7^2\}$ on Swin-B. Each decoder contains one transformer block with the embedding dimension of 128 and 4 attention heads. For ViT, the chosen layers $\mathcal{I} = \{2, 4, 4+n, 6+n\}$, e.g., $n = 6$ on ViT-B and $n = 18$ on ViT-L. The scales of supervision are $\{56^2, 28^2, 14^2, 7^2\}$. Each decoder contains one transformer block with the embedding dimension of 256 and 8 attention heads. The pre-training and fine-tuning schedules mostly follow [24, 31], and more detailed settings can be found in Appendix C.

**Results.** We compare our LocalMIM with existing MIM models and examine both pre-training efficiency and top-1 fine-tuning accuracy on ImageNet-1K. The results are provided in Fig. 1 and Table 1. For fair comparison, we estimate the pre-training duration of each model on the same machine with one Tesla V100-32G GPU, CUDA 10.2 and Pytorch 1.8. We report the running time on single GPU,

| Model | PT Epoch | PT Hours | $AP^b$ | $AP^m$ |
|---|---|---|---|---|
| Supervised | 300 | 840 | 48.5 | 43.2 |
| SimMIM$_{192}$ [60] | 800 | 1440 | 50.4 | 44.4 |
| GreenMIM [31] | 800 | 640 | 50.0 | 44.1 |
| LocalMIM-HOG | 400 | 440 | 50.7 | 44.9 |

Table 3. Object detection and instance segmentation on COCO. We fine-tune Mask R-CNN end-to-end with Swin-B backbone.

denoted as 'GPU Hours', see Appendix B for more details. On ViT-B, LocalMIM achieves the best results of MAE [24] and MaskFeat [57] with $3.1\times$ and $5.6\times$ acceleration respectively. On Swin-B, LocalMIM achieves those of SimMIM$_{192}$ [60] and GreenMIM [31] with $3.6\times$ and $6.4\times$ acceleration respectively. In terms of final top-1 fine-tuning accuracy, LocalMIM achieves 84.0% with ViT-B and 84.1% with Swin-B. Compared with previous best results, LocalMIM achieves comparable performance with significantly less pre-training duration. Note that PeCo [17] uses more advanced feature descriptor, a pre-trained codebook, which introduces additional pre-training burden. Besides, on Swin-B, with 100 epochs pre-training and 100 epochs supervised fine-tuning LocalMIM achieves 83.8% top-1 accuracy and takes about 350 GPU Hours, while 300 epochs supervised training from scratch only achieves 83.5% top-1 accuracy and takes about 720 GPU Hours. This means even for the high-level classification task, self-supervision learning is more efficient and effective than supervised learning. Similar results can also be found on ViT.

## 4.2. Downstream tasks

We transfer our pre-trained backbone to semantic segmentation on ADE20K [66] and object detection and segmentation on COCO [37].

**Semantic segmentation on ADE20K.** We conduct semantic segmentation on ADE20K using UperNet [58] and following the code in [2, 24]. See Appendix C for fine-tuning details. The results are shown in Table 2. With significantly less computation burden, LocalMIM outperforms the state-of-the-art result by 0.7. Besides, LocalMIM has the same top-1 fine-tuning accuracy with MaskFeat [57] on ImageNet-1K, but has better semantic segmentation performance, which means our local multi-scale reconstruction facilitates the learning of multi-scale semantic knowledge.

**Object detection and instance segmentation on COCO.** We fine-tune Mask R-CNN [26] on COCO with Swin-B backbone. Following [31], we also use the code base and $3\times$ fine-tuning schedule from the supervised Swin, where the model is fine-tuned for 36 epochs. See Appendix C for fine-tuning details. We report box AP ($AP^b$) for object detection and mask AP ($AP^m$) for instance segmentation. The results are shown in Table 3. With no labeling burden, our LocalMIM outperforms supervised pre-training by 2.2 $AP^b$

| target | ViT-B | Swin-B |
|---|---|---|
| baselines | 82.9 | 83.2 |
| Pixels | 83.0 | 83.7 |
| HOG | **83.3** | **83.8** |

(a) **Reconstruction target**. HOG is more suitable as multi-scale supervisions than pixels.

| ratio | ViT-B | Swin-B |
|---|---|---|
| 0.40 | 83.2 | 83.7 |
| 0.60 | 83.1 | 83.7 |
| 0.75 | **83.3** | **83.8** |
| 0.90 | 83.0 | 83.5 |

(b) **Mask ratio**. Masking 75% patches works the best for both ViT and Swin.

| decoder | ViT-B | Swin-B |
|---|---|---|
| 512D - 16H | 83.3 (1.0h) | 83.8 (1.6h) |
| 256D - 8H | **83.3 (0.7h)** | 83.7 (1.3h) |
| 128D - 4H | 83.2 (0.6h) | **83.8 (1.1h)** |

(c) **Decoder design**. A tiny decoder performs as well as the larger one but is more efficient.

| locations | backbone | acc |
|---|---|---|
| GreenMIM | Swin-B | 83.2 |
| [24] | | 83.3 |
| [22, 24] | | 83.4 |
| [4, 22, 24] | Swin-B | 83.6 |
| [2, 4, 22, 24] | | **83.8** |
| fusion | | 83.5 |

(d) **Locations for Swin**. Applying reconstructions at all stages mostly accelerates learning.

| locations | backbone | acc |
|---|---|---|
| MAE | ViT-B | 82.9 |
| [12] | | 83.0 |
| [3, 6, 9, 12] | ViT-B | 83.1 |
| [2, 4, 10, 12] | | **83.3** |
| [1, 2, 11, 12] | | 82.9 |

(e) **Locations for ViT**. The chosen layers for local losses obviously affect performance.

| scales | backbone | acc |
|---|---|---|
| MAE | ViT-B | 82.9 |
| $[14^2, 14^2, 14^2, 14^2]$ | | 83.0 |
| $[28^2, 14^2, 7^2, 7^2]$ | | 83.2 |
| $[56^2, 28^2, 14^2, 7^2]$ | ViT-B | **83.3** |
| $[7^2, 14^2, 28^2, 56^2]$ | | 83.1 |

(f) **Scales for ViT**. Using the supervisions from fine-scale to coarse-scale works the best.

Table 4. Ablation studies with ViT-B and Swin-B on ImageNet-1K. We report the top-1 fine-tuning accuracy (%) and the default setting is: the reconstruction target is HOG feature, the mask ratio is 75%, the decoder contains one Transformer block with dimension 256 and 8 heads for ViT and dimension 128 and 4 heads for Swin, the local reconstructions are applied at all stages of Swin and $[2, 4, 10, 12]$-th layer of ViT, the scales of supervision are $[28^2, 14^2, 7^2, 7^2]$ for Swin and $[56^2, 28^2, 14^2, 7^2]$ for ViT, and the pre-training length is 100 epochs.

and 1.7 AP$^m$. Compared to SimMIM$_{192}$ and GreenMIM, our LocalMIM achieves +0.3 and +0.7 gain respectively for detection, and +0.5 and +0.8 gain for segmentation with significantly less pre-training duration.

### 4.3. Ablation studies

We ablate our LocalMIM using the default setting in Table 4. MAE [24] and GreenMIM [31] are our main baselines on ViT and Swin respectively, since we all adopt the asymmetric encoder-decoder strategy.

**Reconstruction target.** We examine two types of signal as supervisions: HOG feature [57] and normalized pixel [24], and the results are shown in Table 4a. HOG feature is more suitable as multi-scale supervisions than normalized pixel on both ViT and Swin, since it contains more refined semantic information. Compared to the columnar ViT with single-scale feature maps, our multi-scale supervisions achieve more improvement on the pyramidal Swin. Since we need to rescale the predictions in the columnar architectures, it is harder to guide the learning of multi-scale information, e.g., we need more information-refined signals (HOG) and obtain relatively less improvement.

**Mask ratio.** We explore different mask ratios, ranging from 0.4 to 0.9, and the results are provided in Table 4b. Similar to [24, 31], our LocalMIM is robust to the mask ratio from 0.4 to 0.75, and excessively large one (0.9) harms the performance since it over-complicates the reconstruction tasks.

**Decoder design.** As shown in [24, 31], a small single-block decoder can achieve the similar fine-tuning accuracy to the heavier one with multiple blocks. However, since we need to use the decoder for each reconstruction task to process the signals of specified scale, their decoder which

has the embedding dimension of 512 and 16 attention heads is still heavy. To this end, we explore whether a smaller decoder can still handle the local reconstruction tasks for learning good representations. One Transformer block is the minimal requirement to propagate information from visible patches to mask tokens, so we still use one Transformer block but narrow the embedding dimension and reduce the number of attention heads. The results (top-1 fine-tuning accuracy and GPU Hours per epoch) are shown in Table 4c, where '$d$D - $h$H' represents the decoder with the embedding dimension of $d$ and $h$ attention heads. Surprisingly, a tiny decoder can achieve the same even better fine-tuning performance, but is significantly more efficient.

**Locations for local reconstructions.** An important design of our LocalMIM is the locations in the encoder where we introduce local reconstructions. Many pyramidal architectures like Swin already divide the entire backbone into multiple stages, so we introduce multi-scale reconstructions at their output locations. From the last stage, we gradually include the lower stages and the results are shown in Table 4d. Note that the 2-th, 4-th, 22-th and 24-th layer correspond to the output location of the four stages respectively. We also provide the result of fusion method [20] which fuses the output features from each stage for reconstruction. As we can see, introducing reconstructions at local stages can accelerate representation learning, and applying to all stages works the best. Our LocalMIM also outperforms the fusion method which only uses the single-scale reconstruction, and this verifies the importance of multi-scale reconstructions. Further, for the columnar ViT, we follow the experience from Swin and also select four local layers. Concretely, we consider the uniform division $[3, 6, 9, 12]$, Swin-style divi-
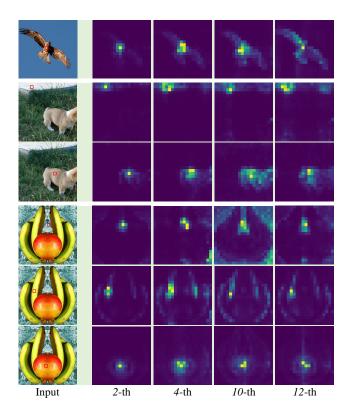
Figure 4. Visualization of the attention maps for different query points, marked with red boxes. LocalMIM can distinguish semantical regions using self-attention mechanism.

sion $[2, 4, 10, 12]$ and a extreme division $[1, 2, 11, 12]$. The results are shown in Table 4e. The performance on columnar ViT is sensitive to the chosen local layers, and the Swin-style division $[2, 4, 4 + n, 6 + n]$ works the best.

**Scales of supervisions.** Since the pyramidal architectures hard-code the multi-scale property of the feature maps, we directly use the supervisions of same scale as the feature maps for compatibility, i.e., we use $[28^2, 14^2, 7^2, 7^2]$ for Swin. For columnar ViT, we consider single-scale supervisions $[14^2, 14^2, 14^2, 14^2]$, fine-to-coarse multi-scale supervisions $[28^2, 14^2, 7^2, 7^2]$ and $[56^2, 28^2, 14^2, 7^2]$, and coarse-to-fine multi-scale supervisions $[7^2, 14^2, 28^2, 56^2]$. The results are shown in Table 4f. The multi-scale supervisions achieve better fine-tuning accuracy than the single-scale ones, which means the multiple local layers require the supervision of different scales. Further, the fine-to-coarse supervisions outperform the coarse-to-fine ones, which is consistent with the design priors for visual architectures [39] and biological visual processing [33].

**Analysis of self-attention mechanism.** To qualitatively illustrate the behavior of self-attention mechanism in LocalMIM, we visualize the attention maps for different query positions within an image, as shown in Fig. 4. We use each chosen patch as query and show the patches it attends to.

| model | backbone | GPU Hours/Ep. | acc |
|---|---|---|---|
| LocalMIM | ViT-B | 0.7 | 83.3 |
| w/ isolated grad | | 0.7 | 83.0 |
| LocalMIM | Swin-B | 1.1 | 83.8 |
| w/ isolated grad | | 1.1 | 83.7 |

Table 5. Training LocalMIM with isolated gradients achieves similar performance with global back-propagation.

We select the attention maps at 2-th, 4-th, 10-th and 12-th layers of a pre-trained ViT-B/16 backbone to show their changes during forward inference. For object-centric images, LocalMIM can distinguish the foreground object from the background. For more complex multi-object images, LocalMIM can effectively separate different objects without any task-specific supervision, which means the attention maps are query-adaptive. On the other hand, the patches at lower layers typically more focus on their neighboring regions and capture low-level information, while those at upper layers attend to a wide range of semantically related regions and capture high-level shape information.

**Gradient-isolated training.** Inspired by locally supervised learning [42, 43], we remove global back-propagation and stop the gradients after each chosen local layer used for reconstruction. The entire backbone is thus divided into multiple gradient-isolated parts. The results are shown in Table 5. Surprisingly, the gradient-isolated training achieves similar performance to global back-propagation, which further verifies the effectiveness of our multi-scale reconstructions for guiding the local layers. It even requires no gradient information from the upper layers. This observation also shows the promise of our local multi-scale reconstruction for the decoupled training of neural networks, which allows for training very deep networks without memory concern and reduces explosive or vanishing gradients.

## 5. Conclusions

We present a novel and efficient pretext task, *local multi-scale reconstruction*, where the lower layers and upper layers reconstruct the fine-scale and coarse-scale supervisions from the input respectively. We obtain the multi-scale supervisions by first dividing the input under different scales and then extracting supervisions with appropriate feature descriptors. Our model needs no extra pre-trained codebook and no costly forward inference of teacher networks during pre-training. The asymmetric encoder-decoder strategy with tiny encoders also have small computational burden, our model thus can be trained quickly. The novel pretext task further accelerates the representation learning, especially for the pyramidal architectures. All these designs allow our model to achieve comparable performance to existing models but with significantly less pre-training burden. The gradient-isolated training also verifies the effectiveness

8

of our novel task design in guiding the local layers.

# References

[1] 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. *arXiv preprint arXiv:2202.03555*, 2022. 1, 2, 4, 5

[2] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers. In *International Conference on Learning Representations*, 2022. 1, 2, 3, 4, 5, 6, 13

[3] Yoshua Bengio, Dong-Hyun Lee, Jorg Bornschein, Thomas Mesnard, and Zhouhan Lin. Towards biologically plausible deep learning. *arXiv preprint arXiv:1502.04156*, 2015. 3

[4] 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. 1

[5] Shuhao Cao, Peng Xu, and David A Clifton. How to understand masked autoencoders. *arXiv preprint arXiv:2202.03670*, 2022. 4

[6] Natalia Caporale, Yang Dan, et al. Spike timing-dependent plasticity: a hebbian learning rule. *Annual review of neuroscience*, 31(1):25–46, 2008. 2

[7] 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. 5

[8] Jun Chen, Ming Hu, Boyang Li, and Mohamed Elhoseiny. Efficient self-supervised vision pretraining with local masked reconstruction. *arXiv preprint arXiv:2206.00790*, 2022. 1, 2, 5

[9] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017. 3

[10] 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. 2

[11] Xiaokang Chen, Mingyu Ding, Xiaodi Wang, Ying Xin, Shentong Mo, Yunhao Wang, Shumin Han, Ping Luo, Gang Zeng, and Jingdong Wang. Context autoencoder for self-supervised representation learning. *arXiv preprint arXiv:2202.03026*, 2022. 5, 6

[12] 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. 5, 6

[13] Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. ELECTRA: pre-training text encoders as discriminators rather than generators. In *ICLR*, 2020. 13

[14] Francis Crick. The recent excitement about neural networks. *Nature*, 337(6203):129–132, 1989. 3

[15] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 702–703, 2020. 13

[16] Yang Dan and Mu-ming Poo. Spike timing-dependent plasticity of neural circuits. *Neuron*, 44(1):23–30, 2004. 3

[17] 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. 1, 2, 4, 5, 6

[18] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. 1, 2, 3, 6, 12, 13

[19] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. Multiscale vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6824–6835, 2021. 3

[20] Peng Gao, Teli Ma, Hongsheng Li, Jifeng Dai, and Yu Qiao. Convmae: Masked convolution meets masked autoencoders. *arXiv preprint arXiv:2205.03892*, 2022. 2, 3, 5, 7, 12

[21] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017. 13

[22] Jianyuan Guo, Kai Han, Han Wu, Yehui Tang, Yunhe Wang, and Chang Xu. Fastmim: Expediting masked image modeling pre-training for vision. *arXiv preprint arXiv:2212.06593*, 2022. 1

[23] Kai Han, An Xiao, Enhua Wu, Jianyuan Guo, Chunjing Xu, and Yunhe Wang. Transformer in transformer. *Advances in Neural Information Processing Systems*, 34:15908–15919, 2021. 2

[24] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16000–16009, 2022. 1, 2, 3, 4, 5, 6, 7, 13

[25] 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. 2

[26] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international*

*conference on computer vision*, pages 2961–2969, 2017. 3, 6, 14

[27] 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. 12

[28] Byeongho Heo, Jeesoo Kim, Sangdoo Yun, Hyojin Park, Nojun Kwak, and Jin Young Choi. A comprehensive overhaul of feature distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1921–1930, 2019. 3

[29] Jeremy Howard and Sebastian Ruder. Universal language model fine-tuning for text classification. In *ACL*, 2018. 1, 4

[30] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q Weinberger. Deep networks with stochastic depth. In *European conference on computer vision*, pages 646–661. Springer, 2016. 13

[31] Lang Huang, Shan You, Mingkai Zheng, Fei Wang, Chen Qian, and Toshihiko Yamasaki. Green hierarchical vision transformer for masked image modeling. *arXiv preprint arXiv:2205.13515*, 2022. 1, 2, 3, 4, 5, 6, 7, 13, 14

[32] Huawei. Mindspore. https://www.mindspore.cn/, 2020. 9

[33] DH Hubel and TN Wiesel. Receptive fields of optic nerve fibres in the spider monkey. *The Journal of physiology*, 154(3):572, 1960. 3, 8

[34] 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. 1

[35] Xiang Li, Wenhai Wang, Lingfeng Yang, and Jian Yang. Uniform masking: Enabling mae pre-training for pyramid-based vision transformers with locality. *arXiv preprint arXiv:2205.10063*, 2022. 1, 2, 4

[36] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017. 14

[37] 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, 14

[38] Hao Liu, Xinghua Jiang, Xin Li, Antai Guo, Deqiang Jiang, and Bo Ren. The devil is in the frequency: Geminated gestalt autoencoder for self-supervised visual pre-training. *arXiv preprint arXiv:2204.08227*, 2022. 1, 2

[39] 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. 1, 2, 3, 6, 8, 12, 14

[40] Ilya Loshchilov and Frank Hutter. SGDR: stochastic gradient descent with warm restarts. In *ICLR*, 2017. 13

[41] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2019. 13

[42] Sindy Löwe, Peter O'Connor, and Bastiaan Veeling. Putting an end to end-to-end: Gradient-isolated learning of representations. *Advances in neural information processing systems*, 32, 2019. 2, 8

[43] Arild Nøkland and Lars Hiller Eidnes. Training neural networks with local error signals. In *International conference on machine learning*, pages 4839–4850. PMLR, 2019. 2, 8

[44] Namuk Park and Songkuk Kim. How do vision transformers work? In *International Conference on Learning Representations*, 2022. 2, 5

[45] Myeongjang Pyeon, Jihwan Moon, Taeyoung Hahn, and Gunhee Kim. Sedona: Search for decoupled neural networks toward greedy block-wise learning. In *International Conference on Learning Representations*, 2021. 2

[46] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28, 2015. 3

[47] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. In *ICLR*, 2015. 3

[48] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015. 6

[49] Alexander Strehl and Joydeep Ghosh. Cluster ensembles—a knowledge reuse framework for combining multiple partitions. *Journal of machine learning research*, 3(Dec):583–617, 2002. 4, 12

[50] Zhiqing Sun, Zhuohan Li, Haoqing Wang, Di He, Zi Lin, and Zhihong Deng. Fast structured decoding for sequence models. *Advances in Neural Information Processing Systems*, 32, 2019. 2

[51] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826, 2016. 13

[52] 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. 12

[53] Haoqing Wang and Zhi-Hong Deng. Cross-domain few-shot classification via adversarial task augmentation. In Zhi-Hua Zhou, editor, *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*, pages 1075–1081. International Joint Conferences on Artificial Intelligence Organization, 2021. Main Track. 3

[54] Haoqing Wang, Xun Guo, Zhi-Hong Deng, and Yan Lu. Rethinking minimal sufficient representation in contrastive learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16041–16050, 2022. 2

[55] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense

prediction without convolutions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 568–578, 2021. 1, 2, 3

[56] Yulin Wang, Zanlin Ni, Shiji Song, Le Yang, and Gao Huang. Revisiting locally supervised learning: an alternative to end-to-end training. In *International Conference on Learning Representations*, 2021. 2

[57] Chen Wei, Haoqi Fan, Saining Xie, Chao-Yuan Wu, Alan Yuille, and Christoph Feichtenhofer. Masked feature prediction for self-supervised visual pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14668–14678, 2022. 1, 2, 4, 5, 6, 7

[58] 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. 6, 13

[59] Jiahao Xie, Wei Li, Xiaohang Zhan, Ziwei Liu, Yew Soon Ong, and Chen Change Loy. Masked frequency modeling for self-supervised visual pre-training. *arXiv preprint arXiv:2206.07706*, 2022. 1, 2

[60] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han Hu. Simmim: A simple framework for masked image modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9653–9663, 2022. 1, 2, 4, 5, 6

[61] Weihao Yu, Mi Luo, Pan Zhou, Chenyang Si, Yichen Zhou, Xinchao Wang, Jiashi Feng, and Shuicheng Yan. Metaformer is actually what you need for vision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10819–10829, 2022. 12

[62] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6023–6032, 2019. 13

[63] Hongyi Zhang, Moustapha Cissé, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *ICLR*, 2018. 13

[64] Pengchuan Zhang, Xiyang Dai, Jianwei Yang, Bin Xiao, Lu Yuan, Lei Zhang, and Jianfeng Gao. Multi-scale vision longformer: A new vision transformer for high-resolution image encoding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2998–3008, 2021. 3

[65] Tianyi Zhang, Felix Wu, Arzoo Katiyar, Kilian Q Weinberger, and Yoav Artzi. Revisiting few-sample bert finetuning. In *International Conference on Learning Representations*, 2021. 1, 4

[66] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 633–641, 2017. 6

[67] Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong. Image bert pre-training with online tokenizer. In *International Conference on Learning Representations*, 2022. 1, 2, 4, 5

|  | single-scale | multi-scale |  | single-scale | multi-scale |
|---|---|---|---|---|---|
| global | 83.0 | 82.8 | global | 83.3 | 82.9 |
| local | 83.0 | 83.3 | local | 83.6 | 83.8 |
| (a) ViT-B |  |  | (b) Swin-B |  |  |

Table 6. Decoupling local reconstruction and multi-scale supervisions. We report the top-1 fine-tuning accuracy on ImageNet-1K.

| method | acc | method | acc |
|---|---|---|---|
| global | 83.0 | global | 83.3 |
| global with fusion | 83.0 | global with fusion | 83.5 |
| local | 83.3 | local | 83.8 |
| (a) ViT-B |  | (b) Swin-B |  |

Table 7. Top-1 fine-tuning accuracy on ImageNet-1K. We compare the global reconstruction, the global reconstruction with feature fusion and our local reconstruction.

# A. More experiments

In this section, we provide more experiments to support our work.

## A.1. Decoupling local and multi-scale

To effectively guide the local layers, we propose multi-scale supervisions for multiple local reconstruction tasks. Here we decouple the local (or global) reconstruction and multi-scale (or single-scale) supervisions to further understand their relations. For global multi-scale reconstruction, we conduct at the top layer of encoder and use separate decoders to predict multiple supervisions of different scales. When the predictions have different scale with supervisions, we use deconvolution/pooling options to rescale them for matching supervisions. The results are shown in Table 6 and the pre-training length is 100 epochs. As we can see, global reconstruction prefers to single-scale supervisions, and using the supervisions of different scales to guide the same layer could make confusion. Conversely, local reconstruction prefers to multi-scale supervisions, and multiple local layers expect to learn the information of different scales. Local reconstruction can achieve better performance than the global one in most cases, and the gain increases when using multi-scale supervisions.

## A.2. Comparison with feature fusion

To explicitly guide the lower layers, we conduct reconstruction task at multiple chosen local layers. The other method is fusing the features of multiple local layers to the top layer for global reconstruction [20]. It uses single-scale supervision for avoiding confusion. We compare our local reconstruction with this feature fusion method, and the results are shown in Table 7. Local reconstruction achieves consistently better performance than feature fusion on both columnar ViT [18] and pyramidal Swin [39]. For further exploration, we examine the gradient norm of each layer in the encoder during training process. Concretely, we load the checkpoint (state) of the median epoch in a complete training schedule and then calculate the gradient norm of parameters in each layer under this state. The results are shown in Fig. 5. For other middle epochs, we observe the same results. The lower layers have larger gradient norm than the upper ones due to the skip-connections in vision transformers. The skip-connections allow the lower layers

to learn more quickly than the upper ones, which may be one reason for its significant effectiveness in various architectures [27,52,61]. Our local reconstruction can strengthen this characteristic and thus obtain better performance. Feature fusion essentially has the similar effect with the skip-connections. Besides, another advantage of our local construction is that it is compatible with multi-scale supervisions and thus can take advantage of richer information.

## A.3. Query-adaptive attention

In the main text, we use Normalized Mutual Information (NMI) [49] between query and key patches to examine how much the attention map depends on the query patch. Here we use another metric, the Kullback-Leibler divergence between the attention distributions of different query patches. Intuitively, when the attention map strongly depends on the query patch, the attention distributions of a pair of query patches should have large KL divergence. We calculate the average on all pairs of query patches at each layer and the results are shown in Fig. 6. As we expect, existing MIM models with global loss have small KL divergence at lower layers, which means the patches there have less query-adaptive attention. Relatively, the lower layers in our LocalMIM have larger KL divergence and the attention maps depend more strongly on the query patches.

# B. GPU Hours

'GPU Hours' denotes the running time on single Tesla V100-32G GPU. For fair comparison, we estimate that of each model at the same machine with one Tesla V100-32G GPU, CUDA 10.2 and PyTorch 1.8. We pre-train each model for 10 epochs using its official released codes and default hyper-parameters, and then calculate the average running time per epoch. We find that each epoch takes similar time with each other during estimation, so pre-training 10 epochs is enough to estimate the GPU Hours per epoch. The batch size is an important factor that affects the running time, and we choose it from $\{32, 48, 64, 128, 256\}$ to take full advantage of GPU memory and computing capability. This estimation method avoids the interference of the communication time among multiple GPUs.
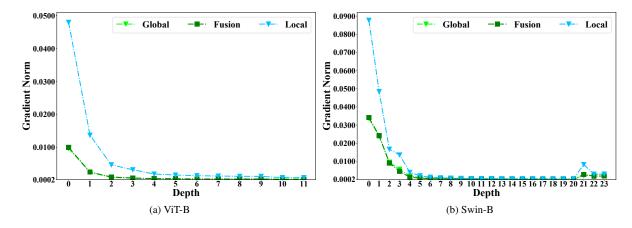
(a) ViT-B



(b) Swin-B

Figure 5. Gradient norm of each layer in the encoder. We compare the global reconstruction, the global reconstruction with feature fusion and our local reconstruction, which are denoted as 'Global', 'Fusion' and 'Local' respectively.
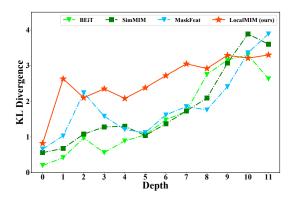


Figure 6. The KL divergence between attention distributions of different query patches at each layer of a pre-trained ViT-B backbone, averaged on all pairs of query patches.

## C. Implementation details

For ViT [18], we use the standard architecture with the sine-cosine positional embeddings and do not use relative positional encoding or layer scaling. For HOG feature, we set the number of orientation bins $\#bins = 18$ and the cell size is the same as the divided regions. We set the same weight to each local loss for simplicity. The pre-training and fine-tuning schedules mostly follow [24, 31].

**Pre-training.** The default setting is shown in Table 8. We use the simple data augmentation and do not use drop path or gradient clip. We use the linear learning rate scaling rule [21]: $lr = base\_lr \times batch\_size/256$. The warmup epoch [21] is set to 10 for pre-training 100 epochs, 40 for pre-training 400, 800 and 1600 epochs.

**Fine-tuning on ImageNet-1K.** The default fine-tuning setting is shown in Table 9. Most of the hyper-parameters are shared, except the peak learning rate, layer-wise learning rate decay and drop path rate, which are influenced by the

| config | ViT | Swin |
|---|---|---|
| optimizer | AdamW [41] | |
| base learning rate | $2e^{-4}$ | $1e^{-4}$ |
| weight decay | 0.05 | |
| optimizer momentum | $\beta_1, \beta_2 = 0.9, 0.95$ | |
| batch size | 2048 (B) / 4096 (L) | |
| learning rate schedule | cosine decay [40] | |
| augmentation | RandomResizedCrop | |
| input resolution | $224 \times 224$ | |

Table 8. Pre-training setting on ImageNet-1K.

| config | ViT | | Swin | |
|---|---|---|---|---|
| | ViT-B | ViT-L | Swin-B | Swin-L |
| optimizer | AdamW | | | |
| peak learning rate | $\{2e^{-3}, 3e^{-3}, 4e^{-3}\}$ | | $\{3e^{-3}, 4e^{-3}, 5e^{-3}\}$ | |
| weight decay | 0.05 | | | |
| optimizer momentum | $\beta_1, \beta_2 = 0.9, 0.999$ | | | |
| layer-wise lr decay [13] | $\{0.65, 0.75\}$ | | $\{0.80, 0.90\}$ | |
| batch size | 1024 (B) / 4096 (L) | | | |
| learning rate schedule | cosine decay | | | |
| fine-tuning epochs | 100 | 50 | 100 | 100 |
| warmup epochs | 20 | 5 | 20 | 20 |
| drop path [30] | 0.1 | 0.2 | 0.1 | 0.3 |
| augmentation | RandAug (9, 0.5) [15] | | | |
| label smoothing [51] | 0.1 | | | |
| mixup [63] | 0.8 | | | |
| cutmix [62] | 1.0 | | | |
| input resolution | $224 \times 224$ | | | |

Table 9. Fine-tuning setting on ImageNet-1K.

backbones and the number of pre-training epochs.

**Semantic segmentation on ADE20K.** We use UperNet [58] with ViT-B backbone and follow the semantic segmentation code of [2, 24]. Concretely, we fine-tune end-to-end for 160K iterations using AdamW optimizer with the peak learning rate of $4e^{-4}$, weight decay of 0.05 and batch size of 16. The learning rate warmups with 1500 iterations and then decays with linear strategy. The model is trained with

input resolution of $512 \times 512$ and uses bilinear positional embedding interpolate. We choose the out indices of feature maps as $[2, 4, 10, 12]$ and use FPN [36] to rescale them.

**Object detection and segmentation on COCO.** We fine-tune Mask R-CNN [26] on COCO [37] with Swin-B backbone. Following [31], we also use the code base and schedule from [39]. Concretely, the model is fine-tuned on COCO 2017 train split and evaluated on 2017 val split. We adopt the $3\times$ fine-tuning schedule which trains the model for 36 epochs in total and decays the learning rate at the 27-th and 33-th epoch by a factor of 10. We use AdamW optimizer with the learning rate of $1e^{-4}$ and weight decay of $0.05$.