# Revealing the Dark Secrets of Masked Image Modeling

**Zhenda Xie**[*13], **Zigang Geng**[*23], **Jingcheng Hu**[13], **Zheng Zhang**[3], **Han Hu**[3], **Yue Cao**[3†]

[1]Tsinghua University
[2]University of Science and Technology of China
[3]Microsoft Research Asia

## Abstract

Masked image modeling (MIM) as pre-training is shown to be effective for numerous vision downstream tasks, but how and where MIM works remain unclear. In this paper, we compare MIM with the long-dominant supervised pre-trained models from two perspectives, the visualizations and the experiments, to uncover their key representational differences. From the visualizations, we find that MIM brings locality inductive bias to all layers of the trained models, but supervised models tend to focus locally at lower layers but more globally at higher layers. That may be the reason why MIM helps Vision Transformers that have a very large receptive field to optimize. Using MIM, the model can maintain a large diversity on attention heads in all layers. But for supervised models, the diversity on attention heads almost disappears from the last three layers and less diversity harms the fine-tuning performance. From the experiments, we find that MIM models can perform significantly better on geometric and motion tasks with weak semantics or fine-grained classification tasks, than their supervised counterparts. Without bells and whistles, a standard MIM pre-trained SwinV2-L could achieve state-of-the-art performance on pose estimation (78.9 AP on COCO test-dev and 78.0 AP on CrowdPose), depth estimation (0.287 RMSE on NYUv2 and 1.966 RMSE on KITTI), and video object tracking (70.7 SUC on LaSOT). For the semantic understanding datasets where the categories are sufficiently covered by the supervised pre-training, MIM models can still achieve highly competitive transfer performance. With a deeper understanding of MIM, we hope that our work can inspire new and solid research in this direction.

## 1 Introduction

Pre-training of effective and general representations applicable to a wide range of tasks in a domain is the key to the success of deep learning. In computer vision, supervised classification on ImageNet [13] has long been the dominant pre-training task which is manifested to be effective on a wide range of vision tasks, especially on the semantic understanding tasks, such as image classification [16, 42, 40, 17, 55], object detection [67, 25, 65, 32], semantic segmentation [57, 74], video action recognition [69, 72, 6, 56] and so on. Over the past several years, "masked signal modeling", which masks a portion of input signals and tries to predict these masked signals, serves as a universal and effective self-supervised pre-training task for various domains, including language, vision, and speech. After (masked) language modeling repainted the NLP field [14, 53], recently, such task has also been shown to be a competitive challenger to the supervised pre-training in computer vision [7, 17, 2, 30, 85, 77]. That is, masked image modeling (MIM) pre-trained models achieve very high fine-tuning accuracy on a wide range of vision tasks of different nature and complexity.

However, there still remain several questions:

---

* Equal Contribution. The work is done when Zhenda Xie, Zigang Geng, and Jingcheng Hu are long-term interns at Microsoft Research Asia. † Contact person.

1. What are the key mechanisms that contribute to the excellent performance of MIM?
2. How transferable are MIM and supervised models across different types of tasks, such as semantic understanding, geometric and motion tasks?

To investigate these questions, we compare MIM with supervised models from two perspectives, the visualization perspective and the experimental perspective, trying to uncover key representational differences between these two pre-training tasks and deeper understand the behaviors of MIM pre-training.

We start with studying the attention maps of the pre-trained models. Firstly, we visualize the averaged attention distance in MIM models, and we find that **masked image modeling brings locality inductive bias to the trained model, that the models tend to aggregate near pixels in part of the attention heads,** and the locality strength is highly correlated with the masking ratio and masked patch size in the pre-training stage. But the supervised models tend to focus locally at lower layers but more globally at higher layers.

We next probe how differently the attention heads in MIM trained Transformer behave. We find that **different attention heads tend to aggregate different tokens on all layers in MIM models**, according to the large KL-divergence on attention maps of different heads. But for supervised models, the diversity on attention heads diminishes as the layer goes deeper and almost disappears in the last three layers. We drop the last several layers for supervised pre-trained models during fine-tuning and find that it benefits the fine-tuning performance on downstream tasks, however this phenomenon is not observed for MIM models. That is, **less diversity on attention heads would somewhat harm the performance on downstream tasks**.

Then we examine the representation structures in the deep networks of MIM and supervised models via the similarity metric of Centered Kernel Alignment (CKA) [41]. We surprisingly find that **in MIM models, the feature representations of different layers are of high similarity, that their CKA values are all very large (e.g., [0.9, 1.0]).** But for supervised models, as in [63], different layers learn different representation structures, that their CKA similarities vary greatly (e.g., [0.5,1.0]). To further verify this, we load the pre-trained weights of randomly shuffled layers during fine-tuning and find that supervised pre-trained models suffer more than the MIM models.

From the experimental perspective, a fundamental pretraining task should be able to benefit a wide range of tasks, or at least it is important to know for which types of tasks MIM models work better than the supervised counterparts. To this end, we conduct a large-scale study by comparing the fine-tuning performance of MIM and supervised pre-trained models, on three types of tasks, semantic understanding tasks, geometric and motion tasks, and the combined tasks which simultaneously perform both.

For semantic understanding tasks, we select several representative and diverse image classification benchmarks, including Concept Generalization (CoG) benchmark [66], the widely-used 12-dataset benchmark [42], as well as a fine-grained classification dataset iNaturalist-18 [73]. For the classification datasets whose categories are sufficiently covered by ImageNet categories (e.g. CIFAR-10/100), supervised models can achieve better performance than MIM models. However, for other datasets, such as fine-grained classification datasets (e.g., Food, Birdsnap, iNaturalist), or datasets with different output categories (e.g., CoG), most of the representation power in supervised models is difficult to transfer, thus MIM models remarkably outperform supervised counterparts.

For geometric and motion tasks that require weaker semantics and high-resolution object localization capabilities, such as pose estimation on COCO [52] and CrowdPose [48], depth estimation on NYUv2 [68] and KITTI [22], and video object tracking on GOT10k [36], TrackingNet [59], and LaSOT [20], MIM models outperform supervised counterparts by large margins. Note that, without bells and whistles, Swin-L with MIM pre-training could achieve state-of-the-art performance on these benchmarks, e.g., 80.5 AP on COCO $val$, 78.9 AP on COCO $test$-$dev$, and 78.0 AP on CrowdPose of pose estimation, 0.287 RMSE on NYUv2 and 1.966 RMSE on KITTI of depth estimation, and 70.7 SUC on LaSOT of video object tracking.

We select object detection on COCO as the combined task which simultaneously performs both semantic understanding and geometric learning. For object detection on COCO, MIM models would outperform supervised counterparts. Via investigating the training losses of object classification and localization, we find that MIM models help localization task converge faster, and supervised models benefit more for object classification, that categories of COCO are fully covered by ImageNet.

In general, MIM models can perform significantly better on geometric/motion tasks with weak semantics or fine-grained classification tasks, than the supervised counterparts. For tasks/datasets where supervised models are good at transfer, MIM models can still achieve highly competitive transfer performance. It seems time to embrace masked image modeling as a general-purpose pre-trained model. We hope our paper can drive this belief deeper in the community and inspire new and solid research in this direction.

## 2    Background

**Masked Image Modeling.** Masked image modeling (MIM) is a sub-task of masked signal prediction, that masks a portion of input images, and lets the deep networks predict the masked signals conditioned on the visible ones. We use SimMIM [77], a simple framework for masked image modeling, as the exampled framework of pre-trained image models in our visualizations and experiments, because it is simple, effective, and generally applicable. SimMIM consists of four major components with simple designs: 1) random masking with a moderately large masked patch size (e.g., 32); 2) the masked tokens and image tokens are fed together to the encoder; 3) the prediction head is as light as a linear layer: 4) directly predicting raw pixels of RGB values as the target with the $\ell_1$ loss of direct regression. With these simple designs, SimMIM can achieve state-of-the-art performance on ImageNet-1K classification, COCO object detection, and ADE-20K semantic segmentation. Note that, the SimMIM framework could be directly applied to different types of backbone architectures, such as Vision Transformer (ViT) [17], Swin Transformer [55], and ConvNets [33, 15]. This property enables us to study the characteristics of MIM under different types of backbone architectures, as well as in multiple types of downstream tasks.

**Backbone Architectures.** Masked image modeling is mostly studied in the Transformer architectures, thus the major understandings and experiments in this paper are performed on Vision Transformers (ViT) [17] and Swin Transformers [55, 54]. Due to the simple and clear architecture designs of ViT, most of the visualizations are performed on ViT, shown in Section 3. Due to the general-purpose property of Swin Transformer, most of the experiments on different downstream tasks are conducted on Swin Transformer, shown in Section 4.

## 3    Visualizations

### 3.1    Revealing the Properties of Attention Maps

Attention mechanism [1] has been an exceptional component in deep networks. It is naturally interpretable since attention weights have a clear meaning: how much each token is weighted when determining the output representation of the current token. Fortunately, most MIM pre-trained models [17, 2, 30, 85, 77] are established upon the Vision Transformers, where self-attention block is its major component. Here we start with studying the attention maps of the pre-trained models from three angles: (a) averaged attention distance to measure whether it is local attention or global attention; (b) entropy of attention distribution to measure whether it is focused attention or broad attention; (c) KL divergence of different attention heads to investigate that attention heads are attending different tokens or similar ones.

#### 3.1.1    Local Attention or Global Attention?

Images are observed to exhibit strong locality: pixels near each other tend to be highly correlated [37], motivating the use of local priors in a wide range of visual perception architectures [21, 46, 44, 33, 55]. In the era of Vision Transformers, the usefulness of local priors has still undergone rich discussions and trials [17, 55, 49]. Thus it is valuable to investigate whether MIM models bring the locality inductive bias to the models. We do this by computing averaged attention distance in each attention head of each layer.

Results of the averaged attention distance in different attention heads (dots) w.r.t the layer number, on supervised model (DeiT), contrastive learning model (MoCo v3) and SimMIM model with ViT-B as backbone are shown in Figure 1. We find that the supervised model tends to focus locally at lower layers but more globally at higher layers, which well matches the observations in ViT [17]. Surprisingly, the contrastive learning model acts very similarly to the supervised counterpart. This
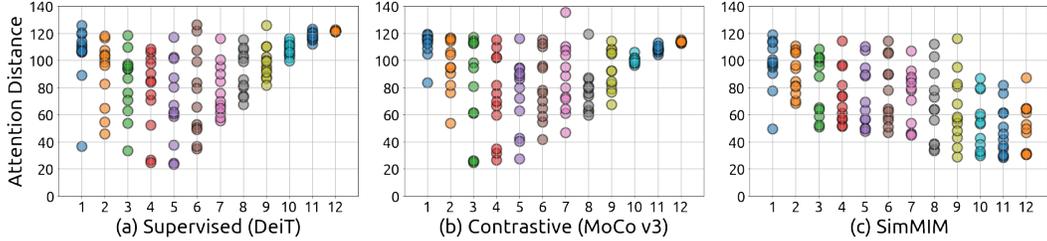
Figure 1: The averaged attention distance in different attention heads (dots) w.r.t the layer number on supervised model (a), contrastive learning model (b), and SimMIM model (c) with ViT-B as the backbone architecture.
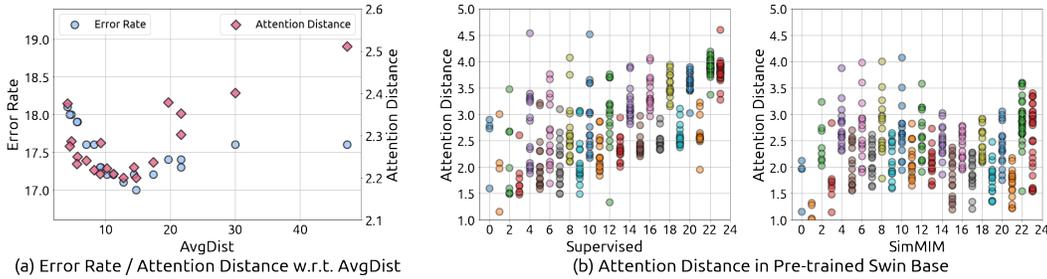


Figure 2: (a) The error rate of fine-tuning on ImageNet-1K (blue circle ∘) and averaged attention distance (red diamond ◇) w.r.t AvgDist (averaged distance of masked pixels to the nearest visible pixels) with Swin-B as the backbone. Points (◇ or ∘) denote the SimMIM models with different masking ratios and masked patch sizes. (b) The averaged attention distance in different attention heads (dots) w.r.t the layer number on supervised model (b1) and SimMIM model (b2) with Swin-B as the backbone.

may also be understandable, since MoCo v3 has a very high linear evaluation accuracy on ImageNet-1K (76.7% of top-1 accuracy), which indicates that the features of the last layer of MoCo v3 are very similar to that of the supervised counterpart. But for the model trained by SimMIM, its behavior is significantly different to supervised and contrastive learning models. Each layer has diverse attention heads that tend to aggregate both local and global pixels, and the average attention distance is similar to the lower layers of the supervised model. As the number of layers gets deeper, the averaged attention distance becomes even slightly smaller. That is, MIM brings locality inductive bias to the trained model, that the models tend to aggregate near pixels in part of the attention heads. Also, a similar observation could be observed with Swin-B as the backbone, as shown in Figure 2(b).

SimMIM [77] designed a new metric, AvgDist, which measures the averaged Euclidean distance of masked pixels to the nearest visible ones and indicates the task difficulty and effectiveness of MIM depending on the masking ratio and masked patch size. As shown in Figure 2(a), AvgDist is a good indicator that the entries of high fine-tuning accuracy roughly distribute in a range of [10, 20] of AvgDist, while entries with smaller or higher AvgDist perform worse. Interestingly, in the range of [10, 20] of AvgDist, we can also observe a small averaged attention distance. That is, a moderate prediction distance in MIM will bring a greater strength of locality and incur a better fine-tuning performance.

### 3.1.2 Focused Attention or Broad Attention?

We then measure the attention maps on whether attention heads focus on a few tokens or attend broadly over many tokens, via averaging the entropy of each head's attention distribution. Results of entropy values w.r.t different layers of three pre-trained models, supervised model (DeiT), contrastive learning model (MoCo v3), and MIM model (SimMIM) with ViT-B as the backbone, are shown in Figure 3. For supervised models, we find that some attention heads in lower layers have very focused attention, but in higher layers, most attention heads focus very broadly. The contrastive model still behaves very similarly to the supervised model. But for the MIM model, the entropy
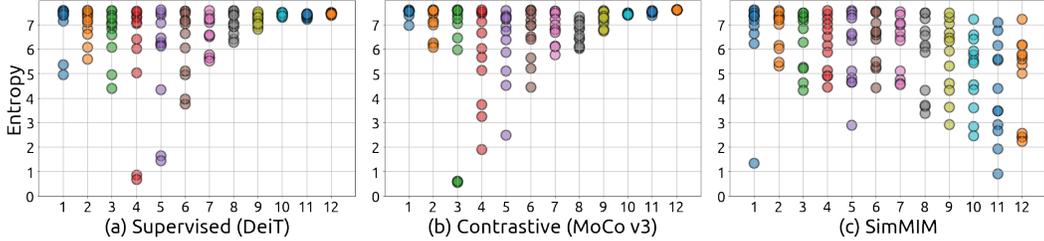
Figure 3: The entropy of each head's attention distribution w.r.t the layer number on (a) supervised model, (b) contrastive learning model, and (c) SimMIM model with ViT-B as the backbone.
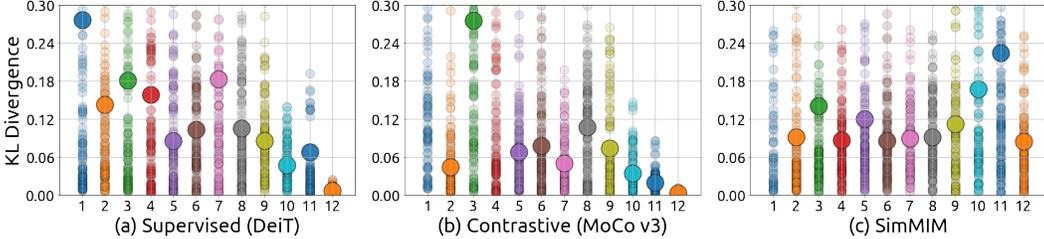


Figure 4: The KL divergence between attention distributions of different heads (small dots) and the averaged KL divergence (large dots) in each layer w.r.t the layer number on (a) supervised model, (b) contrastive learning model, and (c) SimMIM model with ViT-B as the backbone architecture.

values in different attention heads are diverse in all layers, that some attention heads are more focused and some heads have very broad attention.

### 3.1.3 Diversity on Attention Heads

From the previous two sub-sections, we observe a similar phenomenon, that is, for the supervised model, the attention distance or entropy of attention heads in the last few layers seem to be similar, while for the MIM model, different heads in all layers behave more diversely. Therefore, we want to further explore whether the different heads pay attention to different/similar tokens, via computing the Kullback–Leibler (KL) divergence [45] between the attention maps of different heads in each layer.

Results of KL divergence between attention distributions of different heads w.r.t different layers of three pre-trained models, supervised model (DeiT), contrastive learning model (MoCo v3), and MIM model (SimMIM) with ViT-B as the backbone, are shown in Figure 4. As we expect, different attention heads tend to aggregate different tokens on all layers in MIM models, according to the large KL-divergence on attention maps of different heads. But for supervised models and contrastive learning models, the diversity on attention heads becomes smaller as the layer goes deeper and almost disappears from the last three layers.

Intuitively, losing diversity across different attention heads may limit the capacity of the model. To investigate whether the loss of diversity on attention heads has any adverse effect, we gradually drop layers from the end, and only load previous layers when fine-tuning the model for the downstream tasks of COCO $val2017$ pose estimation and NYUv2 depth estimation. From Figure 5, we can observe that when we drop two to eight layers, although the model becomes smaller, the performance of the supervised pre-trained model on COCO $val2017$ pose estimation is better than the baseline, and the performance on NYUv2 depth estimation is comparable with the baseline. This shows that in the supervised pre-trained model, the last layers with small diversity on attention heads indeed affect the performance of downstream tasks. The detailed setup of this experiment is in the Appendix.

### 3.2 Investigating the Representation Structures via CKA similarity

Studying the behaviors of attention mechanisms is analyzing inside the block, from a micro perspective. Next, we hope to study from a macro perspective of deep networks, such as studying the
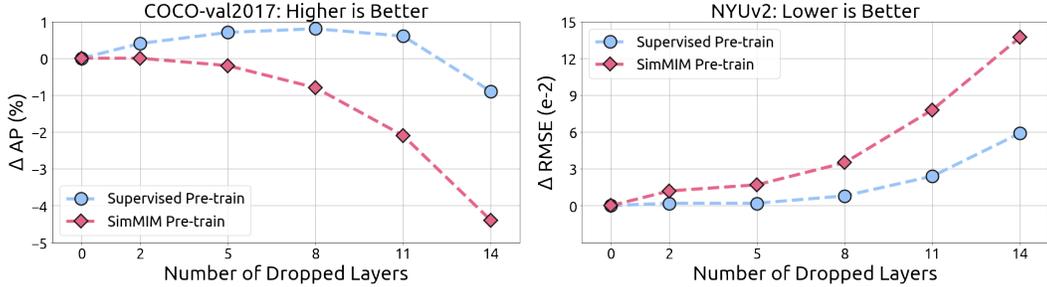
Figure 5: The performance of the COCO $val2017$ pose estimation (left) and NYUv2 depth estimation (right) when we drop several last layers of the SwinV2-B backbone. When the model becomes smaller, the performance of the supervised pre-trained model increases on the pose estimation and keeps the same on the depth estimation. The last layers in the supervised pre-trained model lose diversity across different attention heads and are harmful to the downstream tasks.
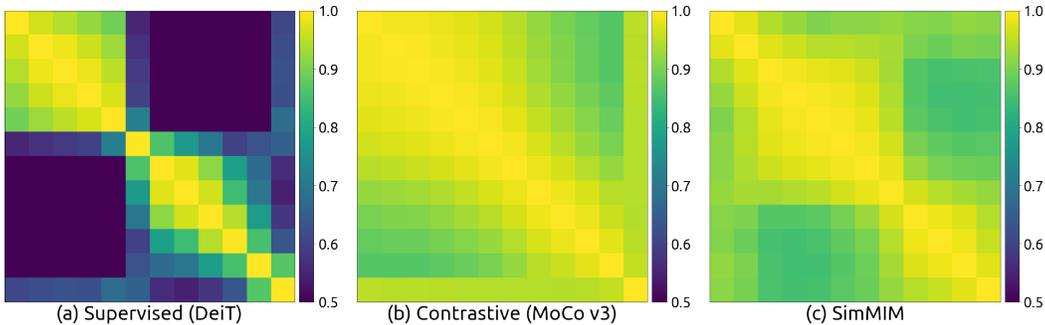


Figure 6: The CKA heatmap between the feature maps of different layers of (a) supervised model, (b) contrastive learning model, and (c) SimMIM model with ViT-B as the backbone architecture.

similarity between feature maps across different layers via the CKA similarity [41]. Results of CKA similarity between feature representations of different layers of three pre-trained models, supervised model (DeiT), contrastive learning model (MoCo v3), and MIM model (SimMIM) with ViT-B as the backbone, are shown in Figure 6. We surprisingly find that in MIM models, the representation structures of different layers are almost the same, that their CKA similarities are all very large (e.g., [0.9, 1.0]). But for supervised models, as in [63], different layers learn different representation structures, that their CKA similarities vary greatly (e.g., [0.5,1.0]). Different from previous visualizations, MoCo v3 behaves similarly to SimMIM in this case.

To further verify this observation, we load the pre-trained weights of randomly shuffled layers and fine-tune the model for the downstream tasks of COCO pose estimation and NYUv2 depth estimation. We observe that by loading the models with the randomly sampled layers, the performance on 1K-MIM drops from 75.5 to 75.2 (-0.3) on pose estimation and 0.382 to 0.434 (-0.052) on depth estimation. But supervised pre-trained models suffer more than the MIM models, which drops from 75.8 to 74.9 (-0.9) on pose estimation, and 0.376 to 0.443 (-0.067) on depth estimation. The detailed setup of this experiment is in the Appendix.

## 4   Experimental Analysis on Three Types of Downstream Tasks

In this section, we conduct a large-scale study by comparing the fine-tuning performance of MIM and supervised pre-trained models, on three types of tasks, semantic understanding tasks (e.g., image classification in different domains), geometric and motion tasks (e.g., pose/depth estimation, and video object tracking), and the combined tasks which simultaneously perform both types of tasks (e.g., object detection). We use 8 NVIDIA V100 GPUs for our experiments.

| pre-train | Concept Generalization (CoG) | | | | | Kornlith et al's 12 datasets (K12) | | | | | iNat18 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | $L_1$ | $L_2$ | $L_3$ | $L_4$ | $L_5$ | Food | Birdsnap | Cars | Aircraft | Average (7) | |
| 1K-SUP | 79.4 | 76.2 | 72.7 | 72.5 | 68.4 | 93.2 | 81.8 | 88.6 | 83.0 | 89.7 | 77.7 |
| 1K-MIM | 79.6 | 77.1 | 73.6 | 73.0 | 69.1 | 94.2 | 83.7 | 89.2 | 83.5 | 86.1 | 79.6 |

Table 1: Comparisons of MIM and supervised (SUP) pre-trained models on semantic understanding tasks with SwinV2-B as the backbone. We follow [42] to report top-1 accuracy ($\uparrow$) and mean per-class accuracy ($\uparrow$) for specific datasets. Results on the multi-label dataset Pascal Voc 2007 are not included, whose evaluation metric is not compatible with others.

## 4.1 Semantic Understanding Tasks

For semantic understanding tasks, we select several representative and diverse image classification benchmarks, including Concept Generalization (CoG) benchmark [66], the widely-used 12-dataset benchmark [42], as well as a fine-grained classification dataset iNaturalist-18 [73].

**Setup.** The CoG benchmark consists of five 1k-category datasets split from ImageNet-21K, which has an increasing semantic gap with ImageNet-1K, from $L_1$ to $L_5$. On the CoG dataset, we search for the best hyper-parameters based on the top-1 accuracy of the $L_1$ validation set and then apply the best setting to CoG $L_2$ to $L_5$ to report the top-1 accuracy. On the K12 dataset, we adopt standard splits of train/val/test sets as in [42]. We use the training set to fine-tune the models, use the validation set to search for the best hyper-parameters, and then train the models on the merged training and validation sets using the best setting. Following [42], we report mean-per-class accuracy for Aircraft, Pets, Caltech-101, Oxford 102 Flowers and top-1 accuracy for other datasets. The iNat18 dataset includes 437,513 training images and 24,426 validation images, with more than 8,000 categories. We fine-tune the pre-trained models using the training set and report the top-1 accuracy on the validation set. For all datasets, we choose learning rate, weight decay, layer decay, and DropPath [34] on the valid set respectively for the MIM pre-trained model and the supervised pre-trained model. We use the AdamW optimizer [58] and cosine learning rate schedule. We train the model for 100 epochs with 20 warm-up epochs. The input image size is $224 \times 224$. Other detailed setups of these datasets are in the Appendix.

**Results.** Results of different semantic understanding tasks are shown in Table 1. For the classification datasets whose categories are sufficiently covered by ImageNet categories (e.g. CIFAR-10/100), supervised models can achieve better performance than MIM models as pre-training. However, for other datasets, such as fine-grained classification datasets (e.g., Food, Birdsnap, iNaturalist), or datasets with different output categories (e.g., CoG), most of the representation power in supervised models is difficult to transfer; thus MIM models remarkably outperform supervised counterparts.

## 4.2 Geometric and Motion Tasks

We study how MIM models perform on the geometric and motion tasks that require the ability to localize the objects and are less dependent on semantic information. We select several benchmarks, such as pose estimation on COCO [52] and CrowdPose [48], depth estimation on NYUv2 [68] and KITTI [22], and video object tracking on GOT10k [36], TrackingNet [59], and LaSOT [20].

**Setup.** For pose estimation on COCO and Crowdpose, we use the standard splits for training and evaluation and report the AP based on OKS as the evaluation metric. We use the standard person detection results from [75]. We follow Simple Baseline [75], which upsamples the last feature of the backbone by deconvolutions and predicts the heatmaps at $4\times$ resolution. The data augmentations include random flipping, half body transformation, random scale, random rotation, grid dropout, and color jittering. The input image size is $256 \times 256$ by default. We use the AdamW [58] optimizer with the base learning rate $5e$-4 and the weight decay $5e$-2. The learning rate is dropped to $5e$-5 at the $120th$ epoch. We train the models for 150 epochs. We use a layer decay of 0.9/0.85 for Swin-B/L and the DropPath [34] of 0.3/0.5 for Swin-B/L.

For depth estimation on NYUv2 and KITTI, we use the standard splits and report the RMSE (Root Mean Square Error) as the evaluation metric. To compare with the previous works [64, 38], we set the maximum range as 10m/80m for NYUv2/KITTI. The head of the depth estimation is the same as that of the pose estimation and is comprised of deconvolutions. Similar to the GLPDepth [38], we use the

| backbone | pre-train | Pose Estimation | | | Depth Estimation | | Video Object Tracking | | |
|---|---|---|---|---|---|---|---|---|---|
| | | COCO *val* | COCO *test* | Crowd-Pose | NYUv2 | KITTI | GOT10k *test* | Track-Net | LaSOT |
| SwinV2-B | 1K-SUP | 75.2 | 74.5 | 70.7 | 0.352 | 2.313 | 70.1 | 81.5 | 69.4 |
| | 22K-SUP | 75.9 | 75.1 | 72.2 | 0.335 | 2.240 | 69.9 | 81.0 | 67.8 |
| | 1K-MIM | **77.6** | **76.7** | **74.9** | **0.304** | **2.050** | **70.8** | **82.0** | **70.0** |
| SwinV2-L | 22K-SUP | 76.5 | 75.7 | 72.7 | 0.334 | 2.150 | 71.1 | 81.5 | 69.2 |
| | 1K-MIM | **78.1** | **77.2** | **75.5** | **0.287** | **1.966** | **72.9** | **82.5** | **70.7** |
| Representative methods | | HRFormer [79] | | | BinsFormer [50] | | MixFormer [12] | | |
| | | 77.2 | 76.2 | 72.5 | 0.330 | 2.098 | 75.6 | 83.9 | 70.1 |

Table 2: Comparisons of MIM and supervised (SUP) pre-trained models on the geometric and motion tasks. We report the AP (↑) for the pose estimation tasks, RMSE (↓) for the monocular depth estimation tasks, AO (↑) for the GOT10K dataset, and SUC (↑) for the TrackingNet dataset and LaSOT tracking dataset. The best results among the different pre-trained models are shown in the **bold** text. We provide the best results of the representative methods for reference.

following data augmentations: random horizontal flip, random brightness/gamma/hue/saturation/value and random vertical CutDepth. We randomly crop the images to $480 \times 480$ / $352 \times 352$ size for NYUv2/KITTI dataset. The optimizer, layer decay, and DropPath is the same as the pose estimation. The learning rate is scheduled via polynomial strategy with a factor of $0.9$ with a minimal value of $3e\text{-}5$ and a maximum value of $5e\text{-}4$. The total number of epochs is $25$. We use the flip testing and sliding window test.

Following the previous methods [51, 12], we train the models on the train splits of four datasets GOT10k [36], TrackingNet [59], LaSOT [20], and COCO [52] and report the success score (SUC) for the TrackingNet dataset and LaSOT dataset, and the average overlap (AO) for GOT10k. We use the SwinTrack [51] to train and evaluate our pre-trained models with the same data augmentations, training, and inference settings. We sample $131072$ pairs per epoch and train the models for $300$ epochs. We use the AdamW optimizer with a learning rate of $5e\text{-}4$ for the head, a learning rate of $5e\text{-}5$ for the backbone, and a weight decay of $1e\text{-}4$. We decrease the learning rate by a ratio of $0.1$ at the 210th epoch. We set the sizes of search images and templates as $224 \times 224$ and $112 \times 112$.

**Results.** From Table 2, for the pose estimation, MIM models pre-trained with ImageNet-1K surpass supervised counterparts by large margins, $2.4$ AP on COCO *val*, $2.2$ AP on COCO *test-dev*, and $4.2$ AP on CrowdPose dataset which contains more crowded scenes. Even if the supervised models are pre-trained with ImageNet-22K, the performances are still worse than MIM models pre-trained with ImageNet-1K. The observation of the SwinV2-L is similar to that of the SwinV2-B. With a larger image size $384 \times 384$, MIM pre-trained SwinV2-L reaches $78.4$ on COCO *test-dev*, and $77.1$ on the challenging CrowdPose dataset. Using a stronger detection result from BigDetection [3], we obtain $80.5$ AP on COCO *val*, $78.9$ AP on COCO *test-dev*, and $78.0$ AP on CrowdPose.

For the depth estimation, using a simple deconvolution head, SwinV2-B with MIM pre-training with ImageNet-1K achieves $0.304$ RMSE on NYUv2 and $2.050$ RMSE on KITTI, outperforming the previous SOTA method BinsFormer-L [50]. The MIM pre-training does improve the performance of SwinV2-B by $0.03$ RMSE compared with the supervised pre-training with ImageNet-22K. Note that with supervised pre-training, a larger model SwinV2-L shows no gain for the NYUv2 dataset, while with MIM pre-training, SwinV2-L leads to about $0.02$ RMSE gain over SwinV2-B.

For the video object tracking, MIM models also show a stronger transfer ability over supervised pre-trained models. On the long-term dataset LaSOT, SwinTrack [51] with MIM pre-trained SwinV2-B backbone achieves comparable result with the SOTA MixFormer-L [12] with a larger image size $320 \times 320$. We obtain the best SUC of $70.7$ on the LaSOT with SwinV2-L backbone with the input image size $224 \times 224$ and template size $112 \times 112$.

## 4.3 Combined Task of Object Detection

We select object detection on COCO as the combined task which simultaneously performs both semantic understanding and geometric learning. For object detection, a Mask-RCNN[32] framework is adopted and trained with a $3\times$ schedule (36 epochs). We utilize an AdamW [39] optimizer with a
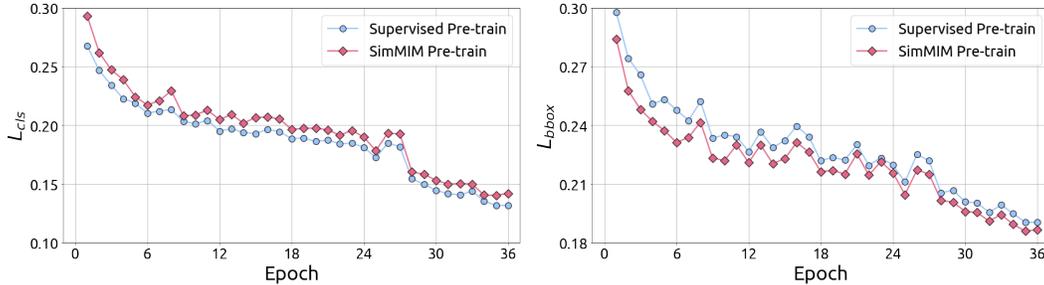
Figure 7: Loss curves of $L_{cls}$ and $L_{bbox}$ w.r.t the epoch number using supervised and MIM models with SwinV2-B as the backbone architecture.

learning rate of 6e-5/8e-5 for supervised/MIM model, a weight decay of 0.05, and a batch size of 32. We employ a large jittering augmentation ($1024 \times 1024$ resolution, scale range [0.1, 2.0]).

On COCO, we could clearly observe that MIM model outperforms its supervised counterpart (52.9/46.7 v.s. 51.9/45.7 of box/mask AP) with SwinV2-B as the backbone. We also plot the loss curves of object classification $L_{cls}$ and localization $L_{bbox}$, as shown in Figure 7. We find that MIM model helps localization task converge faster and better, and the supervised model benefits more for object classification. This also matches our previous observations, that MIM model can perform better on geometric and motion tasks, and on par or slightly worse on the tasks that its categories are sufficiently covered by ImageNet like COCO.

## 5 Related Work

**Visual Pre-training.** Throughout the deep learning era, supervised classification on ImageNet [13] has been the dominant pretraining task. It is found to deliver strong finetuning performance on numerous semantic understanding tasks [16, 42, 40, 17, 67, 25, 55, 57, 69, 6]. Over the past several years, self-supervised pretraining has attracted more and more attention, and achieved finetuning performance on par with the supervised counterparts on several representative downstream tasks [31, 8], including two representative ones, contrastive learning [18, 31, 8, 5, 27] and masked image modeling [7, 2, 30, 77]. In our work, we focus on understanding the different behaviors of supervised and emergent MIM pre-training.

**Understanding Pre-training.** There are some outstanding works [78, 42, 81, 61, 60, 63, 62] trying to understand the pre-training procedure and inspire a lot of following works in a wide range. [78] reveals how features of different layers are transferable in deep neural networks. [42] performs a sufficient experimental study on different backbones and tries to answer whether better ImageNet models transfer better. Some works [63, 84, 63] try to understand the behaviors of ViT, with CKA [41], loss landscape [47] and Fourier analysis. In NLP, after BERT [14] pre-training came out, there is also a lot of works [43, 28, 10, 29] trying to understand it. Most of them focus on the only interpretable component of Transformer, self-attention block, to give some detailed understanding.

## 6 Conclusion

In this work, we present a sufficient and sound analysis on masked image modeling, to reveal how and where MIM models work well. From visualizations, our most interesting finding is that the MIM pre-training brings locality to the trained model with sufficient diversity on the attention heads. This reveals why MIM is very helpful to the Vision Transformers (ViT, Swin, etc), because the Vision Transformer has a much larger receptive field, and to optimize it to a solution with strong generalization ability is difficult. In experiments, our most interesting finding is that MIM pre-training can perform very well on the geometric and motion tasks with weak semantics. This finding helps the model to achieve state-of-the-art performance on those benchmarks without bells and whistles.

It seems time to embrace masked image modeling as a general-purpose pre-trained model. We hope our paper can drive this belief deeper in the community and inspire new and solid research in this direction. The best destination for an understanding paper would be to appear in the motivation of future technologies.

# References

[1] Bahdanau, D., Cho, K., and Bengio, Y. (2014). Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.

[2] Bao, H., Dong, L., and Wei, F. (2021). Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*.

[3] Cai, L., Zhang, Z., Zhu, Y., Zhang, L., Mu, L., and Xue, X. (2022). Bigdetection: A large-scale benchmark for improved object detector pre-training. *arXiv preprint arXiv:2203.13249*.

[4] Cai, Z. and Vasconcelos, N. (2019). Cascade r-cnn: High quality object detection and instance segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.

[5] Cao, Y., Xie, Z., Liu, B., Lin, Y., Zhang, Z., and Hu, H. (2020). Parametric instance classification for unsupervised visual feature learning. *Advances in Neural Information Processing Systems*, 33.

[6] Carreira, J. and Zisserman, A. (2017). Quo vadis, action recognition? a new model and the kinetics dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6299–6308.

[7] Chen, M., Radford, A., Child, R., Wu, J., and Jun, H. (2020a). Generative pretraining from pixels. *Advances in Neural Information Processing Systems*.

[8] Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. (2020b). A simple framework for contrastive learning of visual representations. *ICML*.

[9] Cheng, B., Misra, I., Schwing, A. G., Kirillov, A., and Girdhar, R. (2021). Masked-attention mask transformer for universal image segmentation. *arXiv preprint arXiv:2112.01527*.

[10] Clark, K., Khandelwal, U., Levy, O., and Manning, C. D. (2019). What does bert look at? an analysis of bert's attention. *arXiv preprint arXiv:1906.04341*.

[11] Cubuk, E. D., Zoph, B., Mane, D., Vasudevan, V., and Le, Q. V. (2018). Autoaugment: Learning augmentation policies from data. *arXiv preprint arXiv:1805.09501*.

[12] Cui, Y., Jiang, C., Wang, L., and Wu, G. (2022). Mixformer: End-to-end tracking with iterative mixed attention. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*.

[13] Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. (2009). Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255. Ieee.

[14] Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

[15] Ding, X., Zhang, X., Zhou, Y., Han, J., Ding, G., and Sun, J. (2022). Scaling up your kernels to 31x31: Revisiting large kernel design in cnns. *arXiv preprint arXiv:2203.06717*.

[16] Donahue, J., Jia, Y., Vinyals, O., Hoffman, J., Zhang, N., Tzeng, E., and Darrell, T. (2014). Decaf: A deep convolutional activation feature for generic visual recognition. In *International conference on machine learning*, pages 647–655. PMLR.

[17] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Houlsby, N. (2021). An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*.

[18] Dosovitskiy, A., Springenberg, J. T., Riedmiller, M., and Brox, T. (2014). Discriminative unsupervised feature learning with convolutional neural networks. In *Advances in neural information processing systems*, pages 766–774.

[19] Eigen, D., Puhrsch, C., and Fergus, R. (2014). Depth map prediction from a single image using a multi-scale deep network. In *Advances in neural information processing systems*, pages 2366–2374.

[20] Fan, H., Lin, L., Yang, F., Chu, P., Deng, G., Yu, S., Bai, H., Xu, Y., Liao, C., and Ling, H. (2019). Lasot: A high-quality benchmark for large-scale single object tracking. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5374–5383.

[21] Fukushima, K. (1975). Cognitron: A self-organizing multilayered neural network. *Biological cybernetics*, 20(3):121–136.

[22] Geiger, A., Lenz, P., Stiller, C., and Urtasun, R. (2013). Vision meets robotics: The kitti dataset. *International Journal of Robotics Research*, 32(11):1231–1237.

[23] Geng, Z., Sun, K., Xiao, B., Zhang, Z., and Wang, J. (2021). Bottom-up human pose estimation via disentangled keypoint regression. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14676–14686.

[24] Ghiasi, G., Cui, Y., Srinivas, A., Qian, R., Lin, T.-Y., Cubuk, E. D., Le, Q. V., and Zoph, B. (2021). Simple copy-paste is a strong data augmentation method for instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2918–2928.

[25] Girshick, R., Donahue, J., Darrell, T., and Malik, J. (2014). Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 580–587.

[26] Gretton, A., Bousquet, O., Smola, A., and Schölkopf, B. (2005). Measuring statistical dependence with hilbert-schmidt norms. In *International conference on algorithmic learning theory*, pages 63–77. Springer.

[27] Grill, J.-B., Strub, F., Altché, F., Tallec, C., Richemond, P., Buchatskaya, E., Doersch, C., Avila Pires, B., Guo, Z., Gheshlaghi Azar, M., et al. (2020). Bootstrap your own latent-a new approach to self-supervised learning. *Advances in Neural Information Processing Systems*, 33.

[28] Hao, Y., Dong, L., Wei, F., and Xu, K. (2019). Visualizing and understanding the effectiveness of bert. *arXiv preprint arXiv:1908.05620*.

[29] Hao, Y., Dong, L., Wei, F., and Xu, K. (2020). Self-attention attribution: Interpreting information interactions inside transformer. *arXiv preprint arXiv:2004.11207*, 2.

[30] He, K., Chen, X., Xie, S., Li, Y., Dollár, P., and Girshick, R. (2021). Masked autoencoders are scalable vision learners. *arXiv preprint arXiv:2111.06377*.

[31] He, K., Fan, H., Wu, Y., Xie, S., and Girshick, R. (2020). Momentum contrast for unsupervised visual representation learning. *CVPR*.

[32] He, K., Gkioxari, G., Dollár, P., and Girshick, R. (2017). Mask r-cnn. In *ICCV*, pages 2961–2969.

[33] He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *CVPR*, pages 770–778.

[34] Huang, G., Sun, Y., Liu, Z., Sedra, D., and Weinberger, K. Q. (2016). Deep networks with stochastic depth. In *European conference on computer vision*, pages 646–661. Springer.

[35] Huang, J., Zhu, Z., Guo, F., and Huang, G. (2020). The devil is in the details: Delving into unbiased data processing for human pose estimation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5699–5708.

[36] Huang, L., Zhao, X., and Huang, K. (2021). Got-10k: A large high-diversity benchmark for generic object tracking in the wild. *IEEE transactions on pattern analysis and machine intelligence*, 43(5):1562–1577.

[37] Hubel, D. H. and Wiesel, T. N. (1962). Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. *The Journal of physiology*, 160(1):106–154.

[38] Kim, D., Ga, W., Ahn, P., Joo, D., Chun, S., and Kim, J. (2022). Global-local path networks for monocular depth estimation with vertical cutdepth. *arXiv preprint arXiv:2201.07436*.

[39] Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

[40] Kolesnikov, A., Beyer, L., Zhai, X., Puigcerver, J., Yung, J., Gelly, S., and Houlsby, N. (2019). Big transfer (bit): General visual representation learning.

[41] Kornblith, S., Norouzi, M., Lee, H., and Hinton, G. (2019a). Similarity of neural network representations revisited. In *International Conference on Machine Learning*, pages 3519–3529. PMLR.

[42] Kornblith, S., Shlens, J., and Le, Q. V. (2019b). Do better imagenet models transfer better? In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2661–2671.

[43] Kovaleva, O., Romanov, A., Rogers, A., and Rumshisky, A. (2019). Revealing the dark secrets of bert. *arXiv preprint arXiv:1908.08593*.

[44] Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 1097–1105.

[45] Kullback, S. and Leibler, R. A. (1951). On information and sufficiency. *The annals of mathematical statistics*, 22(1):79–86.

[46] LeCun, Y., Haffner, P., Bottou, L., and Bengio, Y. (1999). Object recognition with gradient-based learning. In *Shape, contour and grouping in computer vision*, pages 319–345. Springer.

[47] Li, H., Xu, Z., Taylor, G., Studer, C., and Goldstein, T. (2018). Visualizing the loss landscape of neural nets. *Advances in neural information processing systems*, 31.

[48] Li, J., Wang, C., Zhu, H., Mao, Y., Fang, H., and Lu, C. (2019). Crowdpose: Efficient crowded scenes pose estimation and a new benchmark. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10863–10872.

[49] Li, Y., Zhang, K., Cao, J., Timofte, R., and Van Gool, L. (2021). Localvit: Bringing locality to vision transformers. *arXiv preprint arXiv:2104.05707*.

[50] Li, Z., Wang, X., Liu, X., and Jiang, J. (2022). Binsformer: Revisiting adaptive bins for monocular depth estimation. *arXiv preprint arXiv:2204.00987*.

[51] Lin, L., Fan, H., Xu, Y., and Ling, H. (2021). Swintrack: A simple and strong baseline for transformer tracking. *arXiv preprint arXiv:2112.00995*.

[52] Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., and Zitnick, C. L. (2014). Microsoft coco: Common objects in context. In *ECCV*, pages 740–755. Springer.

[53] Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., and Stoyanov, V. (2019). Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.

[54] Liu, Z., Hu, H., Lin, Y., Yao, Z., Xie, Z., Wei, Y., Ning, J., Cao, Y., Zhang, Z., Dong, L., et al. (2021a). Swin transformer v2: Scaling up capacity and resolution. *arXiv preprint arXiv:2111.09883*.

[55] Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., and Guo, B. (2021b). Swin transformer: Hierarchical vision transformer using shifted windows. *arXiv preprint arXiv:2103.14030*.

[56] Liu, Z., Ning, J., Cao, Y., Wei, Y., Zhang, Z., Lin, S., and Hu, H. (2021c). Video swin transformer.

[57] Long, J., Shelhamer, E., and Darrell, T. (2015). Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440.

[58] Loshchilov, I. and Hutter, F. (2019). Decoupled weight decay regularization. In *International Conference on Learning Representations*.

[59] Müller, M., Bibi, A., Giancola, S., Al-Subaihi, S., and Ghanem, B. (2018). Trackingnet: A large-scale dataset and benchmark for object tracking in the wild. In *Proceedings of the European Conference on Computer Vision*, pages 310–327. Springer.

[60] Neyshabur, B., Sedghi, H., and Zhang, C. (2020). What is being transferred in transfer learning? *Advances in neural information processing systems*, 33:512–523.

[61] Nguyen, T., Raghu, M., and Kornblith, S. (2020). Do wide and deep networks learn the same things? uncovering how neural network representations vary with width and depth. *arXiv preprint arXiv:2010.15327*.

[62] Park, N. and Kim, S. (2022). How do vision transformers work? *arXiv preprint arXiv:2202.06709*.

[63] Raghu, M., Unterthiner, T., Kornblith, S., Zhang, C., and Dosovitskiy, A. (2021). Do vision transformers see like convolutional neural networks? *Advances in Neural Information Processing Systems*, 34.

[64] Ranftl, R., Bochkovskiy, A., and Koltun, V. (2021). Vision transformers for dense prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12159–12168.

[65] Ren, S., He, K., Girshick, R., and Sun, J. (2015). Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99.

[66] Sariyildiz, M. B., Kalantidis, Y., Larlus, D., and Alahari, K. (2021). Concept generalization in visual representation learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9629–9639.

[67] Sermanet, P., Eigen, D., Zhang, X., Mathieu, M., Fergus, R., and LeCun, Y. (2013). Overfeat: Integrated recognition, localization and detection using convolutional networks. *arXiv preprint arXiv:1312.6229*.

[68] Silberman, N., Hoiem, D., Kohli, P., and Fergus, R. (2012). Indoor segmentation and support inference from rgbd images. In *Proceedings of the European Conference on Computer Vision*, pages 746–760. Springer.

[69] Simonyan, K. and Zisserman, A. (2014). Two-stream convolutional networks for action recognition in videos. In *Advances in neural information processing systems*, pages 568–576.

[70] Sun, K., Xiao, B., Liu, D., and Wang, J. (2019). Deep high-resolution representation learning for human pose estimation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5693–5703.

[71] Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., and Wojna, Z. (2016). Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826.

[72] Tran, D., Bourdev, L., Fergus, R., Torresani, L., and Paluri, M. (2015). Learning spatiotemporal features with 3d convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 4489–4497.

[73] Van Horn, G., Mac Aodha, O., Song, Y., Cui, Y., Sun, C., Shepard, A., Adam, H., Perona, P., and Belongie, S. (2018). The inaturalist species classification and detection dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8769–8778.

[74] Wang, X., Girshick, R., Gupta, A., and He, K. (2018). Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7794–7803.

[75] Xiao, B., Wu, H., and Wei, Y. (2018a). Simple baselines for human pose estimation and tracking. In *Proceedings of the European Conference on Computer Vision*, pages 472–487. Springer.

[76] Xiao, T., Liu, Y., Zhou, B., Jiang, Y., and Sun, J. (2018b). Unified perceptual parsing for scene understanding. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 418–434.

[77] Xie, Z., Zhang, Z., Cao, Y., Lin, Y., Bao, J., Yao, Z., Dai, Q., and Hu, H. (2021). Simmim: A simple framework for masked image modeling. *arXiv preprint arXiv:2111.09886*.

[78] Yosinski, J., Clune, J., Bengio, Y., and Lipson, H. (2014). How transferable are features in deep neural networks? *Advances in neural information processing systems*, 27.

[79] Yuan, Y., Fu, R., Huang, L., Lin, W., Zhang, C., Chen, X., and Wang, J. (2021). Hrformer: High-resolution vision transformer for dense predict. In *Advances in neural information processing systems*, pages 7281–7293.

[80] Yun, S., Han, D., Oh, S. J., Chun, S., Choe, J., and Yoo, Y. (2019). Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6023–6032.

[81] Zhai, X., Puigcerver, J., Kolesnikov, A., Ruyssen, P., Riquelme, C., Lucic, M., Djolonga, J., Pinto, A. S., Neumann, M., Dosovitskiy, A., et al. (2019). A large-scale study of representation learning with the visual task adaptation benchmark. *arXiv preprint arXiv:1910.04867*.

[82] Zhang, H., Cisse, M., Dauphin, Y. N., and Lopez-Paz, D. (2017). mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*.

[83] Zhong, Z., Zheng, L., Kang, G., Li, S., and Yang, Y. (2020). Random erasing data augmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 13001–13008.

[84] Zhou, H.-Y., Lu, C., Yang, S., and Yu, Y. (2021a). Convnets vs. transformers: Whose visual representations are more transferable? In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2230–2238.

[85] Zhou, J., Wei, C., Wang, H., Shen, W., Xie, C., Yuille, A., and Kong, T. (2021b). ibot: Image bert pre-training with online tokenizer. *arXiv preprint arXiv:2111.07832*.

# A    Visualizations on Swin Transformer

It is crucial to know whether our observations in visualizations are general across different backbone architectures. Thanks to the general applicability of SimMIM [77], we further perform the visualizations on SwinV2-B [54] (in Section A) and RepLKNet [15] (in Section B). Fortunately, we find that most of the observations could be transferred across architectures, ViT-B, SwinV2-B, and RepLKNet.
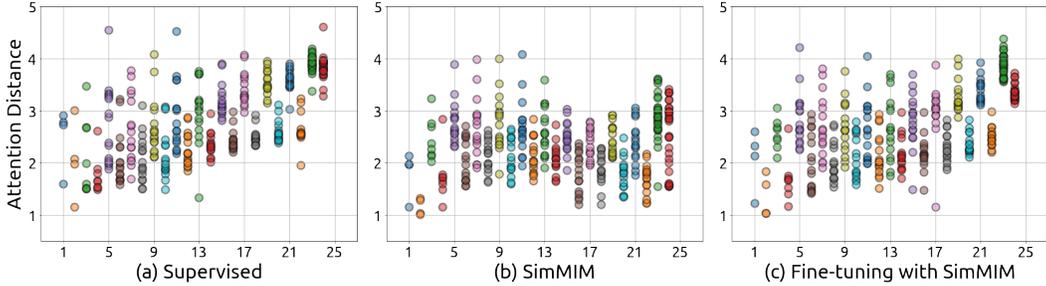
## A.1    Visualizations on Attention Maps



Figure 8: The averaged attention distance in different attention heads (dots) w.r.t the layer number on (a) supervised model, (b) SimMIM model, and (c) supervised fine-tuned model with SimMIM pre-training with SwinV2-B as the backbone.

**Local Attention or Global Attention?**    Results are shown in Figure 8. First, we can have a similar observation as in ViT-B that the supervised model (a) tends to focus locally at lower layers but more globally at higher layers, and the SimMIM model (b) tends to aggregate both local and global pixels in all layers, and the average attention distance of SimMIM model is similar to the lower layers of the supervised counterpart. The supervised fine-tuned model (c) with SimMIM pre-training behaves very similarly to the supervised model trained from scratch, but still maintains some good properties in SimMIM pre-training (a larger diversity on the last several layers). Also, we find that the averaged aggregated distances in two consecutive layers are one high and one low. This is due to the shifted windowing scheme in Swin Transformer, that is, the ranges that each pixel can aggregate in two consecutive layers are different.
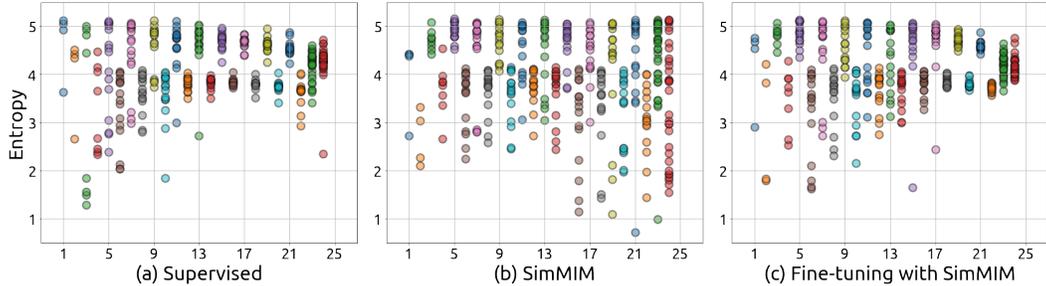


Figure 9: The entropy of each head's attention distribution in different attention heads (dots) w.r.t the layer number on (a) supervised model, (b) SimMIM model, and (c) supervised fine-tuned model with SimMIM pre-training with SwinV2-B as the backbone.

**Focused Attention or Broad Attention?**    A similar observation could be found with Swin-B as the backbone as using ViT-B as the backbone in the main paper, as shown in Figure 9.

**Diversity on Attention Heads**    As shown in Figure 10, similar to ViT-B, in SimMIM models (b), different attention heads tend to aggregate different tokens on all layers. But for supervised models (a), the diversity on attention heads becomes smaller as the layer goes deeper. Interestingly, after supervised fine-tuning the SimMIM model on ImageNet-1K, the model (c) behaves much more similarly to the supervised model (a) trained from scratch, but maintains an advantage of the SimMIM model, that is, a larger diversity on attention heads of the last two layers.
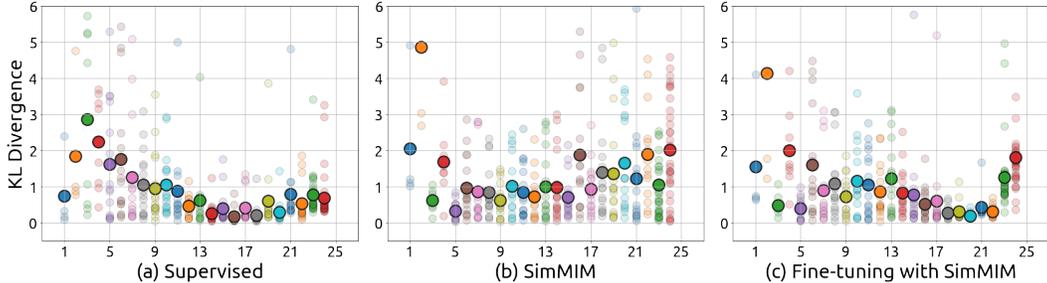
Figure 10: The KL divergence between attention distributions of different heads (small dots) and the averaged KL divergence (large dots) in each layer w.r.t the layer number on (a) supervised model, (b) SimMIM model, and (c) supervised fine-tuned model with SimMIM pre-training with SwinV2-B as the backbone.

## A.2 Investigating the Representation Structures via CKA Similarity
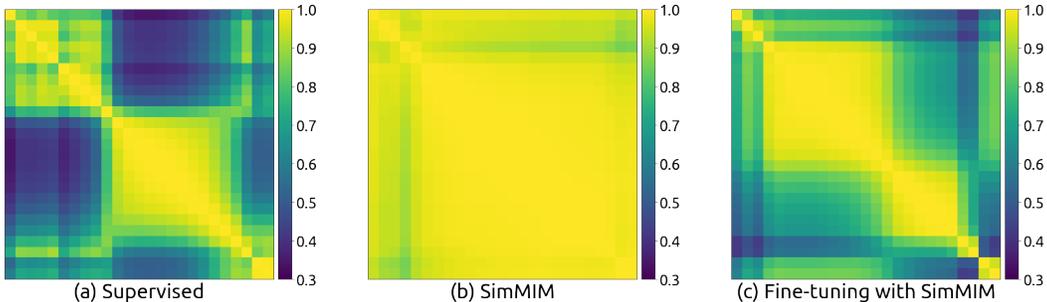


Figure 11: The CKA heatmap between the feature maps of different layers of (a) supervised model, (b) SimMIM model, and (c) supervised fine-tuned model with SimMIM pre-training with SwinV2-B as the backbone.

It is challenging to analyze and compare the layer representations of deep networks, because their features are high-dimensional and with different dimensions. Centered kernel alignment (CKA) [41] is defined to address this challenge, and enables quantitative comparisons of feature representations within and across networks. Given two inputs of $X \in \mathbb{R}^{N \times D_1}$ and $Y \in \mathbb{R}^{N \times D_2}$, where $N$ denotes number of examples and $D_1$ and $D_2$ denote the dimension. Then the Gram matrices are computed as $K = XX^T$ and $L = YY^T$. CKA is then defined as

$$\text{CKA}(K, L) = \frac{\text{HSIC}(K, L)}{\sqrt{\text{HSIC}(K, K)\text{HSIC}(L, L)}}, \tag{1}$$

where $\text{HSIC}(\cdot, \cdot)$ denotes the Hilbert-Schmidt independence criterion [26]. Note that, CKA is invariant to the orthogonal transformation and isotropic scaling, which enables valuable and effective comparison and analysis on hidden representations of deep networks.

Results of CKA similarity between feature representations of different layers on (a) supervised model, (b) SimMIM model, and (c) supervised fine-tuned model with SimMIM pre-training with SwinV2-B as the backbone, are shown in Figure 11. We still have a similar observation as in ViT-B, that the representation structures of different layers in SimMIM models are almost the same, and supervised models trained from scratch learn different representation structures in different layers. With the help of the SimMIM pre-training, the representation structures of different layers in supervised model are not as different as that in the scratch supervised models.

## B Investigations on Large-kernel ConvNets (RepLKNet [15])

From the previous visualizations on Vision Transformers (ViT) and Swin Transformers, we find that the MIM pre-training brings the locality inductive bias and larger diversity on attention heads to the trained models comparing to the supervised counterpart, which may benefit the optimization of the

trained models on downstream tasks. This reminds us that large-kernel ConvNets [15] without special designs still face the optimization issue, and need the re-parametrization trick with small kernels to bring the locality back and help them optimize. Thus it is valuable to know whether the masked image modeling (MIM) as pre-training could help the large-kernel ConvNets to optimize without the re-parametrization trick. Thanks to the general applicability of SimMIM [77], we could also perform experiments and visualizations on large-kernel ConvNets [15] with the MIM pre-training.

## B.1  Experimental Results

| backbone | pre-train | ImageNet-1K | Pose Estimation | | |
| --- | --- | --- | --- | --- | --- |
| | | | COCO *val* | COCO *test* | Crowd-Pose |
| RepLKNet-31B | 1K-SUP w/ Reparam. | 83.5 | 74.6 | 73.9 | 70.2 |
| RepLKNet-31B | 1K-MIM w/o Reparam. | 83.3 | 76.5 | 75.8 | 72.4 |

Table 3: Detailed comparisons of pre-trained RepLKNet models on the classification and the pose estimation tasks. We report the top-1 accuracy ($\uparrow$) for the ImageNet-1K dataset and the AP ($\uparrow$) for the pose estimation tasks.

**Setup**  For MIM pretraining, we utilize the RepLKNet-31B [15] without the specially designed re-parametrization trick. Before the stem of the RepLKNet, using a normal $1 \times 1$ convolution, we map the 3-dimension space of the image into a high-dimensional space where we randomly mask out some patches. Following SimMIM [77], the image size is $192 \times 192$, we divide it into $6 \times 6$ patches and randomly mask out $60\%$ patches. The decoder contains a linear projection layer and an upsample layer. We use $\ell_1$-loss to supervise the reconstruction of the masked pixels.

We use the ImageNet-1k for MIM pre-training and augment the data using the random resize cropping (scale range $[0.67, 1]$ and aspect ratio range $[3/4, 4/3]$), and random flipping. The optimizer is the AdamW[58] optimizer with a weight decay of $5e$-2 and a base learning rate of $4e$-4. We use warm-up for 10 epochs, drop the learning rate to $4e$-5 at 260th epoch, and train for 300 epochs in total. The batch size is 2048. We use the DropPath of $0.1$ for RepLKNet-31B and gradient clipping.

We report the top-1 accuracy of the supervised pre-trained model on ImageNet-1k in the original paper [15]. For fine-tuning of MIM pre-trained model on ImageNet-1k, we follow the setting of SimMIM [77] and use the AdamW optimizer with a weight decay of $5e$-2, a base learning rate of $5e$-3 with a layer decay of $0.8$. The learning rate is scheduled via cosine strategy and we use 20 epochs for warm-up and train for 100 epochs in total. The batch size is 2048. We adopt the DropPath of $0.1$ and gradient clipping. The data augmentations contain AutoAug [11], Mixup [82], CutMix [80], color jitter, random erasing [83], and label smoothing [71]. The settings of the pose estimation are the same as the details in Section E.

**Results**  As shown in Table 3, the MIM pre-training can help the large-kernel convnets to address the optimization issue to some extent and achieve on par performance on ImageNet-1K compared with the supervised model with the re-parametrization trick. Note that, on pose estimation, MIM models still surpass supervised counterparts with the re-parametrization trick by large margins, which indicates that the benefit of MIM pre-training on geometric and motion tasks is general across different backbone architectures.

## B.2  Visualizations

To further understand whether the behaviors of large-kernel ConvNets with MIM pre-training are similar to those of Vision/Swin Transformers, we visualize the convolutional kernels with similar tools used in visualizing the attention maps. As the basic component in RepLKNet is the depth-wise convolution with the kernel dimension of $C \times H \times W$, we normalize each channel of the depth-wise convolutional kernels (on the dimension of $H \times W$) to make them as a similar role of attention map, and regard different channels ($C$ channels) of the depth-wise convolutional kernels as the attention heads. Then we could directly apply the previous tools on attention maps for visualizations.
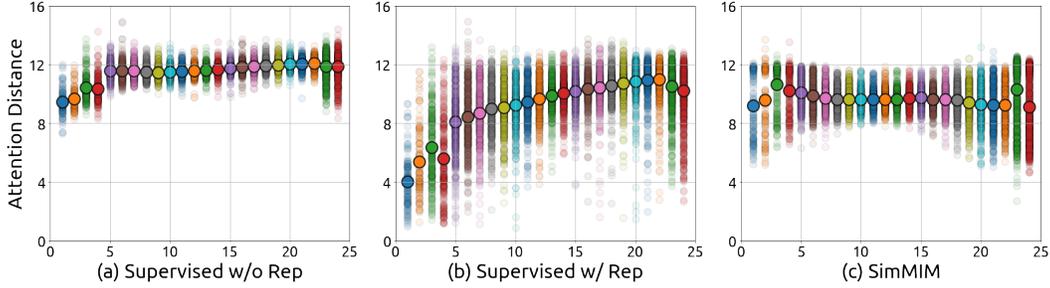
Figure 12: The aggregated distance in different channels (small dots) and the averaged aggregated distance (large dots) w.r.t the layer number on (a) supervised model without the re-parametrization trick, (b) supervised model with the re-parametrization trick, and (c) SimMIM model, with RepLKNet-31B as the backbone architecture.

**Local Kernels or Global Kernels?**  As shown in Figure 12, with the re-parametrization trick, the RepLKNet-31B model (b) with supervised training focuses much more locally in all layers. Similar to previous supervised trained models, RepLKNet-31B models with supervised training still tend to focus locally at lower layers but more globally at higher layers. But for the model trained by SimMIM (c), each layer has diverse kernels that tend to aggregate both local and global pixels, and the average aggregated distance is much smaller than the supervised trained model without the re-parametrization trick (a), indicating that MIM still brings locality inductive bias to the large-kernel ConvNets with a similar role of the re-parametrization trick but less strength.
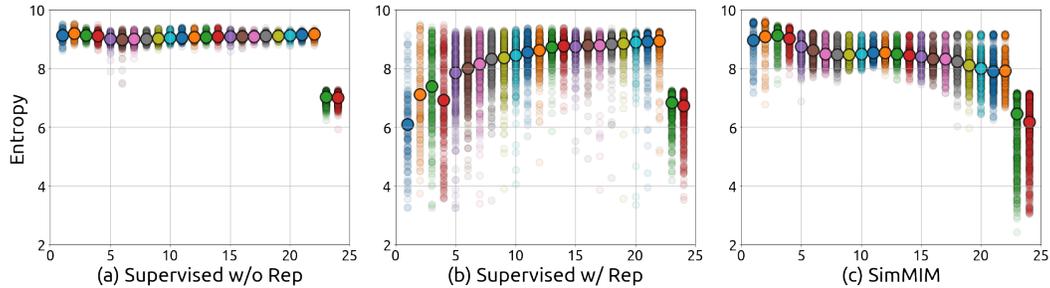


Figure 13: The entropy values in different channels (small dots) and the averaged entropy values (large dots) w.r.t the layer number on (a) supervised model without the re-parametrization trick, (b) supervised model with the re-parametrization trick, and (c) SimMIM model, with RepLKNet-31B as the backbone architecture.

**Focused Kernels or Broad Kernels?**  As shown in Figure 13, with the re-parametrization trick, the supervised RepLKNet-31B model (b) has very focused attention in lower layers, but broader attention in higher layers. But for the MIM model (c), the entropy values in different kernels focus diversely in all layers, that some kernels are more focused and some kernels have very broad attention. These observations well match that in the Vision/Swin Transformers.

**Diversity across Different Kernels**  Interestingly, in Figure 12, it seems that the different kernels in both supervised model with the re-parametrization trick and SimMIM model have diverse averaged aggregated distance. But in Figure 14, we could clearly observe that the diversity on different convolution kernels of SimMIM model (c) is remarkably larger than that of supervised counterparts (b), especially for the deeper layers.

## C   Detailed Results on Semantic Understanding Tasks

Detailed comparisons of Kornblith 12-dataset classification benchmark [42] and Concept Generalization (CoG) benchmark [66] with a fine-grained classification dataset iNaturalist-18 [73] using
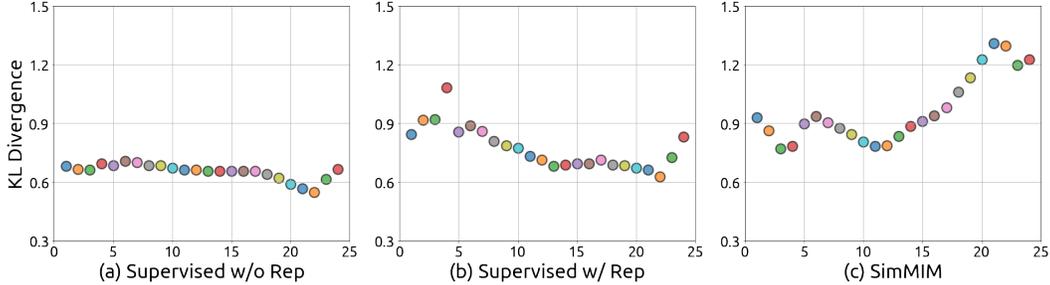
Figure 14: The averaged KL divergence in each layer w.r.t the layer number on (a) supervised model without the re-parametrization trick, (b) supervised model with the re-parametrization trick, and (c) SimMIM model, with RepLKNet-31B as the backbone architecture.

| Methods | Food101 | Birdsnap | Stanford Cars | FGVC Aircraft | Oxford Pets | Caltech101 | Flowers102 | DTD | SUN397 | CIFAR10 | CIFAR100 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1K-SUP | 93.2 | 81.7 | 88.6 | 83.0 | 95.9 | 91.9 | 97.7 | 80.3 | 72.3 | 99.1 | 91.0 |
| 1K-MIM | 94.2 | 83.7 | 89.2 | 83.5 | 90.9 | 85.5 | 91.4 | 73.4 | 70.8 | 99.2 | 91.4 |

Table 4: Detailed comparisons of MIM and supervised (SUP) pre-trained models on Kornblith 12-dataset classification benchmark [42] with SwinV2-B as the backbone. We follow [42] to report top-1 accuracy ($\uparrow$) and mean per-class accuracy ($\uparrow$) for specific datasets. Results on the multi-label dataset Pascal Voc 2007 are not included, whose evaluation metric is not compatible with others.

SwinV2-B as the backbone, are shown in Table 4 and 5, respectively. These results are already discussed in Section 4.1 of the main paper.

## D   Comparisons on Combined Task of Semantic Segmentation

We further select semantic segmentation on ADE-20K as another combine task which simultaneously performs both semantic understanding and geometric learning. For this task, we select two different frameworks, UperNet [76] and Mask2former [9] for evaluation. The detailed settings are shown in Section E.

Results are shown in Table 6. Different to COCO, we find that the supervised pre-trained model slightly outperforms the MIM counterpart on ADE-20K semantic segmentation. Therefore, for the combined tasks, it may be difficult to predict which pretrained model will perform better. But if the model gets larger, MIM models still have the unique advantage that MIM tasks are harder to be overfitted than supervised tasks [30, 77], which is beyond the scope of this paper. Also, we can observe that the performance gap between supervised and MIM models on Mask2former is smaller than that of UperNet ($-1.6$ v.s. $-0.6$). This may be due to that Mask2former decomposes the semantic segmentation task into object localization and recognition tasks, while MIM is better at object localization tasks, as shown in Figure 7 of the main paper.

## E   Detailed Settings

**Concept Generalization benchmark (CoG).** The Concept Generalization benchmark (CoG) consists of five 1k-category datasets splitted from ImageNet-22K, which have increasing semantic gaps with ImageNet-1K, from $L_1$ to $L_5$. On the CoG datasets, for a fair comparison, we first fine-tune the models on the CoG $L_1$ training set and search for the best hyper-parameter based on the validation top-1 accuracy of CoG $L_1$, and then directly apply the searched setting to CoG $L_2$ to $L_5$ and report the top-1 accuracy. The detailed hyperparameters are shown in Table 7.

| pre-train | Concept Generalization (CoG) | | | | | iNat18 |
|---|---|---|---|---|---|---|
| | $L_1$ | $L_2$ | $L_3$ | $L_4$ | $L_5$ | |
| RAND | 79.4 | 76.7 | 73.1 | 72.7 | 68.5 | 76.5 |
| 1K-SUP | 79.4 | 76.2 | 72.7 | 72.5 | 68.4 | 77.7 |
| 1K-MIM | 79.6 | 77.1 | 73.6 | 73.0 | 69.1 | 79.6 |

Table 5: Detailed comparisons of randomly initialized model (RAND), MIM and supervised (SUP) pre-trained models on Concept Generalization (CoG) benchmark [66] and a fine-grained classification dataset iNaturalist-18 [73] with SwinV2-B as the backbone. Top-1 accuracy ($\uparrow$) is reported.

| backbone | pre-train | Object Det. (COCO) | | Semantic Seg. (ADE-20K) | |
|---|---|---|---|---|---|
| | | Mask R-CNN | | UperNet | Mask2former |
| | | $AP^{box}$ | $AP^{mask}$ | mIoU | mIoU |
| SwinV2-B | 1K-SUP | 51.9 | 45.7 | 50.9 | 52.3 |
| | 1K-MIM | 52.9 | 46.7 | 49.3 $(-1.6)$ | 51.7 $(-0.6)$ |

Table 6: Comparisons of MIM and supervised (SUP) pre-trained models on the combined tasks of object detection and semantic segmentation. We report the $AP^{box}$ ($\uparrow$) and $AP^{mask}$ ($\uparrow$) for the object detection and instance segmentation tasks, mIoU ($\uparrow$) for the semantic segmentation task.

**Kornlith et al's 12-dataset benchmark (K12) and iNaturalist-18 (iNat18).** On the K12 dataset, we follow the previous standard settings [42] to use training set and validation set to search for the best hyper-parameters, and then merge the training and validation sets as the final training set with the searched best hyper-parameters, and evaluate the final trained models on the test set. And we adopt standard splits of train/val/test sets as in [42]. For Aircraft, Pets, Caltech-101, Oxford 102 Flowers, the mean-per-class accuracy metric is adopted, for other datasets, the top-1 accuracy is adopted. For K12, we follow [42] to select the optimal learning rate, weight decay, layer decay, and drop path rate. In pilot experiments, we find that for 1K-SUP pre-trained models, the drop path rate can be fixed as $0.2$, and for 1K-MIM pre-trained models, on smaller datasets like Stanford Cars, FGVC Aircraft, DTD, Caltech101, Flowers102, and Oxford Pets, drop path rate is first fixed as $0.0$ and fixed as $0.2$ for other datasets. And the weight decay can be fixed as $0.05$. Then we do a grid search on learning rate and layer decay. For 1K-MIM pre-trained models, our grid consists of 5 approximately logarithmically spaced learning rates between $1.25e$-$4$ and $2.5e$-$3$ and 3 equally spaced layer decay between $0.75$ and $0.95$. For 1K-SUP pre-trained models, our grid consists of 5 approximately logarithmically spaced learning rates between $2.5e$-$5$ and $5e$-$4$ and 3 equally spaced layer decay between $0.75$ and $0.95$. Then we adjust the learning rate, layer decay, and drop path rate in the neighborhood of the best setting in the grid search to get the final results.

The iNat18 dataset includes 437,513 training images and 24,426 validation images, with more than 8,000 categories. The detailed hyperparameters of iNat18 are shown in Table 7.

**Pose estimation.** We compare the performance of MIM and supervised pre-trained models on the COCO [52] and CrowdPose [48] dataset. For the COCO dataset, We train the models on the $train2017$ set ($57K$ training images) and report the performance of the COCO $val2017$ split ($5K$ images), COCO $test$-$dev2017$ split ($20K$ images). For the CrowdPose dataset, following the DEKR [23], we train the models on the CrowdPose train and val sets ($12K$ training images) and evaluate on the test split ($8K$ images). The standard average precision based on OKS is adopted as the evaluation metric for all datasets.

We adopt the heatmap-based top-down pipeline. We upsample the last feature of the backbone by deconvolutions and predict the heatmaps at $4\times$ resolution like Simple Baseline [75].

In the ablation study on the number of the dropped layers in the section 3.1.3 of the main paper, we feed the feature at the different layers in the third stage of SwinV2-B into the pose head. We observe that when we use the feature at the ninth layer, the downstream performances of the supervised pre-trained model and MIM pre-trained model are almost comparable, so we use this model as the baseline of the experiments of randomly sampling pre-trained weights in the section 3.2 of the main paper. In the experiments of randomly sampling pre-trained weights, we randomly sample the weights of nine layers from the weights of the eighteen pre-trained layers in the third stage and then load them to the first nine layers.

| Hyperparameters | RAND | | 1K-SUP | | 1K-MIM | |
|---|---|---|---|---|---|---|
| | CoG (1-5) | iNat18 | CoG (1-5) | iNat18 | CoG (1-5) | iNat18 |
| Input size | | | 224 | | | |
| Window size | | | 14 | | | |
| Patch size | | | 4 | | | |
| Training epochs | 300 | 300 | 100 | 100 | 100 | 100 |
| Warm-up epochs | | | 20 | | | |
| Layer decay | 1.0 | 1.0 | 0.85 | 0.9 | 0.8 | 0.75 |
| Batch size | | | 2048 | | | |
| Optimizer | | | AdamW | | | |
| Base learning rate | 2e-3 | 4e-3 | 2e-4 | 1.6e-3 | 5e-3 | 1.6e-2 |
| Weight decay | 0.05 | 0.1 | 0.05 | 0.1 | 0.05 | 0.1 |
| Adam $\epsilon$ | | | 1e-8 | | | |
| Adam $\beta$ | | | (0.9, 0.999) | | | |
| Learning rate scheduler | | | Cosine | | | |
| Gradient clipping | | | 5.0 | | | |
| Stochastic depth | 0.5 | 0.5 | 0.2 | 0.2 | 0.2 | 0.2 |
| Label smoothing | | | 0.1 | | | |
| Rand crop scale | | | (0.08, 1) | | | |
| Rand resize ratio | | | (3. / 4., 4. / 3.) | | | |
| Rand horizontal flip | | | 0.5 | | | |
| Color jitter | | | 0.4 | | | |
| Rand augment | | | 9 / 0.5 | | | |
| Rand erasing prob. | | | 0.25 | | | |
| Mixup prob. | | | 0.8 | | | |
| Cutmix prob. | | | 1.0 | | | |

Table 7: Detailed settings and hyperparameters for fine-tuning on CoG (1-5) and iNat18 with supervised and MIM pre-trained models.

The data augmentations include random flipping, half body transformation, random scale $(0.5, 1.5)$, random rotation $(-40°, 40°)$, grid dropout and color jitterring (h=0.2, s=0.4, c=0.4, b=0.4). The input image size is $256 \times 256$ by default. We use the AdamW [58] optimizer with the base learning rate $5e$-4 and the weight decay $5e$-2. The learning rate is dropped to $5e$-5 at the $120th$ epoch. We totally train the models for $150$ epochs. We use a layer decay of $0.9/0.85$ for Swin-B/L and the DropPath [34] of $0.3/0.5$ for Swin-B/L. The batch size is $512$.

For the COCO dataset, we use the person detection results from the previous methods [70, 75] for a fair comparison. For the CrowdPose dataset, we use a cascade mask-rcnn [4] with Swin-B backbone trained on the COCO detection dataset to generate the person detection results. We use the UDP [35] to reduce the quantization errors brought by the heatmaps and use flip testing by averaging the heatmaps predicted by the original and flipped images during the inference.

**Depth estimation.** We evaluate the performance of MIM and supervised pre-trained models on the NYUv2 [68] and KITTI [22] monocular depth estimation datasets. The NYUv2 dataset includes $464$ indoor scenes captured by a Microsoft Kinect camera. The official training split ($24K$ images) is used for training and we report the RMSE (Root Mean Square Error) on the $654$ testing images from $215$ indoor scenes. The KITTI dataset contains various driving scenes. The Eigen split [19] contains $23K$ training images and $697$ testing images. To compare with the previous approaches [64, 38], we set the maximum range as 10m for NYUv2 and 80m for KITTI.

The head of the depth estimation is the same as the head of the pose estimation and is comprised of three deconvolutions (with BN and ReLU) and a normal convolution. The kernel and filter of the deconvolution are 2 and 32, respectively.

Similar to the GLPDepth [38], we use the following data augmentations: random horizontal flip, random brightness (-0.2, 0.2), random gamma (-0.2, 0.2), random hue (-20, 20), random saturation (-30, 30), random value (-20, 20) and random vertical CutDepth. We randomly crop the images to

$480 \times 480$ size for NYUv2 dataset and $352 \times 352$ size for KITTI dataset. The optimizer, layer decay, and DropPath is the same as the pose estimation. The learning rate is scheduled via polynomial strategy with a factor of $0.9$. The minimal learning rate and the maximal learning rate are $3e$-$5$ and $5e$-$4$, respectively. The batch size is 24. The total number of epochs is 25. We use the flip testing and sliding window test for the SwinV2 backbone. We average the prediction of the two square windows for NYUv2 dataset and the sixteen square windows for KITTI dataset.

**Video Object Tracking.** Following the previous arts, we train the models on the train splits of four datasets GOT10k [36], TrackingNet [59], LaSOT [20], and COCO [52] and report the success score (SUC) for the TrackingNet dataset and LaSOT dataset. For the GOT10k test set, we report the average overlap as the evaluation metric. The GOT10k and the TrackingNet are two short-term large-scale benchmarks, the GOT10K test set contains 180 video sequences, and the TrackingNet test set contains 511 video sequences. The LaSOT is a long-term tracking benchmark and has 280 video sequences with an average length of about 2500 frames.

We use the SwinTrack [51] to evaluate our pre-trained models. The data augmentations and the training settings Strictly follow SwinTrack [51]. We sample 131072 pairs per epoch and train the models for 300 epochs. We use the AdamW optimizer with a learning rate of $5e$-$4$ for the head, a learning rate of $5e$-$5$ for the backbone, and a weight decay of $1e$-$4$. We decrease the learning rate by a ratio of $0.1$ at the 210th epoch. We set the sizes of search images and templates as $224 \times 224$ and $112 \times 112$. The batch size is 160. The inference process is the same as the SwinTrack [51].

**Object Detection.** Following [77], we adopt a Mask-RCNN [32] framework to evaluate the pre-trained models on COCO object detection. All models are trained with a $3\times$ schedule (36 epochs). We utilize an AdamW [39] optimizer with a learning rate of 6e-5 for supervised model and a learning rate of 8e-5 for MIM model, a weight decay of 0.05 and a batch size of 32 for both models. Following [24, 55], we employ a large jittering augmentation ($1024 \times 1024$ resolution, scale range $[0.1, 2.0]$). The window size is set to 14 for both models and drop path rate is set to 0.3 for supervised model and 0.1 for MIM model. $\text{AP}^{box}$ and $\text{AP}^{mask}$ are reported for comparison.

**Semantic Segmentation.** Following [55], an UPerNet [76] framework is used for ADE-20K semantic segmentation. We use an AdamW [39] optimizer with a learning rate of 8e-5 for supervised model and a learning rate of 1e-4 for MIM model, a weight decay of 0.05 and a batch size of 32 for both models. Both models utilize a layer-wise learning rate decay of 0.95. All models are trained for 80K iterations with an input resolution of $640 \times 640$ and a window size of 20. The drop path rate is set to 0.3 for supervised model and 0.1 for MIM model. In inference, a single-scale test using resolution of $2560 \times 640$ is employed.

Besides, we also adopt Mask2Former [9] to evaluate the pre-trained models on ADE-20K semantic segmentation. We use an AdamW [39] optimizer with a base learning rate of 1e-4 for supervised model and a base learning rate of 3e-4 for MIM model, a weight decay of 0.05 and a patch size of 16 for both supervised and MIM models. The learning rate of backbone is multiplied by a factor of 0.1. All models are trained for 160K iterations with an input resolution of $512 \times 512$, a scale ratio range from 0.5 to 2, a window size of 8, and a drop path rate of 0.3. In inference, the input resolution will be set to $2048 \times 512$. mIoU is reported for comparison for both UPerNet and Mask2Former.