
On Data Scaling in Masked Image Modeling

Zhenda Xie^{*13}, Zheng Zhang^{*3}, Yue Cao^{*3}, Yutong Lin²³, Yixuan Wei¹³, Qi Dai³, Han Hu³

¹Tsinghua University

²Xi'an Jiaotong University

³Microsoft Research Asia

Abstract

An important goal of self-supervised learning is to enable model pre-training to benefit from almost unlimited data. However, one method that has recently become popular, namely masked image modeling (MIM), is suspected to be unable to benefit from larger data. In this work, we break this misconception through extensive experiments, with data scales ranging from 10% of ImageNet-1K to full ImageNet-22K, model sizes ranging from 49 million to 1 billion, and training lengths ranging from 125K iterations to 500K iterations. Our study reveals that: (i) Masked image modeling is also demanding on larger data. We observed that very large models got over-fitted with relatively small data; (ii) The length of training matters. Large models trained with masked image modeling can benefit from more data with longer training; (iii) The validation loss in pre-training is a good indicator to measure how well the model performs for fine-tuning on multiple tasks. This observation allows us to pre-evaluate pre-trained models in advance without having to make costly trial-and-error assessments of downstream tasks. We hope that our findings will advance the understanding of masked image modeling in terms of scaling ability.

1 Introduction

In natural language processing, scaling model capacity and data size has been an important driving force for the remarkable improvements of language models over the past few years [17, 26, 27, 23, 2, 13, 24]. Behind the success is a self-supervised pre-training approach, masked language modeling (MLM) [9], that can take advantage of and benefit from almost unlimited data. As the same time, the relevant research in the field of computer vision has also been intensifying. However, due to the lack of effective self-supervision methods, most previous works are based on image classification tasks [29, 18, 38, 7], where the huge labeling cost and low information contained in the labels limit broader exploration of scaling visual models, or the models being scaled up further, thus leaving progress in computer vision largely behind the NLP field.

Recently, a self-supervision visual pre-training method named masked image modeling (MIM) [1, 15, 36] has become popular due to its impressive fine-tuning performance on a variety of downstream computer vision tasks. Given its high analogy with MLM [9], the dominant pre-training approach in NLP, we expect masked image modeling to advance the scaling performance of visual models. Specifically, we are concerned with two aspects of scaling ability: model scaling and data scaling. While the masked image modeling approach is shown to be good at scaling up model capacity [15, 20], like NLP models, its ability to benefit from larger data is unclear or even a bit negative. For example, [11, 30] show that using a small amount of training data in masked image modeling can achieve comparable performance than that using larger data. The data scaling capability is critical, as an

* Equal Contribution. The work is done when Zhenda Xie, Yutong Lin, and Yixuan Wei are interns at Microsoft Research Asia.

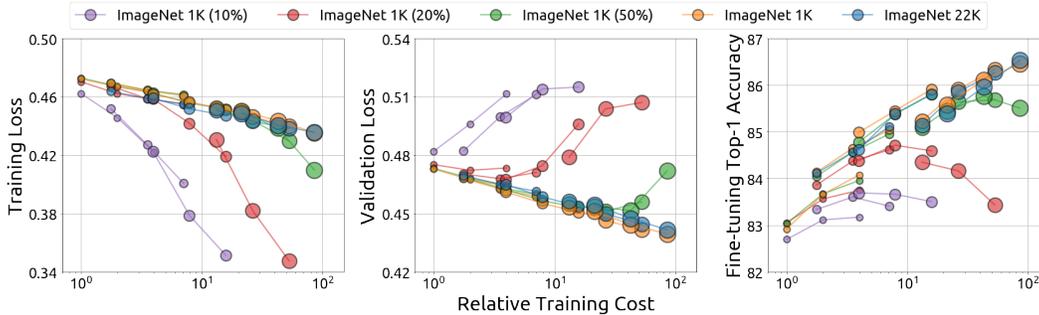


Figure 1: The curves of training loss, validation loss of pre-training, and fine-tuning accuracy on ImageNet-1K of different model sizes, data sizes and training lengths, w.r.t. the relative training cost. We set the training cost of SwinV2-S for 125K iterations as the value of 1. Bigger circles indicate larger models. *Best viewed in color.*

important hallmark of self-supervised learning is the ability to leverage almost unlimited data, and failure to benefit from larger data may hinder the future potential of the masked image modeling.

In this paper, we systematically investigate the data scaling capability of masked image modeling at different model sizes and training lengths. We use Swin Transformer V2 [20] as the visual encoder because of its proven trainability for large models and its applicability to a wide range of vision tasks, and adopt SimMIM [36] for masked image modeling pre-training because it has no restrictions on the encoder architectures. With extensive experiments, we find that:

(i) *Masked image modeling is demanding for large data.* We observed large models overfitted with relatively small data, as reflected by the increased validation losses with longer training when a large model while relatively small data is used (see Figure 1 center). The overfitting issue will result in degraded fine-tuning performance, as shown in Figure 1 right.

(ii) *Training length matters. Large models trained with masked image modeling can benefit from more data at a longer training length.* When the training length is short, the difference in performance between using large and small datasets is not significant. However, with sufficient training, more data shows better performance. In addition, as the data size increases, the fine-tuning performance of large models saturates more slowly than that of small models.

(iii) *The validation loss is a good proxy indicator for fine-tuning performance.* We observe a strong correlation between validation loss and fine-tuning performance on multiple tasks. This finding suggests that the validation loss can be used as a good indicator of how well the model is trained, which can reduce the overhead of evaluation by direct fine-tuning on downstream tasks.

These findings suggest that masked image modeling (MIM) is not only a *model* scalable learner, but also a *data* scalable learner. Particularly, our revealing of data scaling capability of masked image modeling breaks the misconception of previous studies that suspected masked image modeling could not benefit from more data. We hope these findings will deepen the understanding of masked image modeling.

2 Background and Experimental Setup

2.1 Masked Image Modeling

Masked image modeling is used to train the vision model by taking a corrupted image as input and predicting the content of the masked region as the target. In this study, we use SimMIM [36] as the default masked image modeling approach because of its simplicity and lack of restrictions on the architecture of the vision encoder. SimMIM consists of a visual encoder and an extremely lightweight prediction head of a linear layer for predicting the raw pixels of the corrupted images via ℓ_2 regression loss. To facilitate the implementation of the vision transformer, SimMIM adopts the patch-wise mask strategy with the masked patch size of 32×32 and mask ratio of 0.6. To further alleviate the local dependency of raw pixels, we improved the SimMIM by normalizing the predicted target according to [12] with a sliding window of 47^2 . As the result, a slight performance improvement is observed.

Model	Base Channel	Depth	Head	Window Size		Backbone Params
				pre-train	fine-tune	
SwinV2-S	96	{2, 2, 18, 2}	{3, 6, 12, 24}	12	14	49M
SwinV2-B	128	{2, 2, 18, 2}	{4, 8, 16, 32}	12	14	87M
SwinV2-L	192	{2, 2, 18, 2}	{6, 12, 24, 48}	12	14	195M
SwinV2-H	352	{2, 2, 18, 2}	{11, 22, 44, 88}	12	14	655M
SwinV2-g	448	{2, 2, 18, 2}	{14, 28, 56, 112}	12	14	1061M

Table 1: Detailed architecture specifications. SwinV2-g (giant) is a new variant to those in [20], with number of parameters between SwinV2-L and the 3-billion-parameter SwinV2-G (Giant).

	IN1K (10%)	IN1K (20%)	IN1K (50%)	IN100	IN1K (100%)	IN22K (100%)
# Classes	1×10^3	1×10^3	1×10^3	1×10^2	1×10^3	2.18×10^4
# Images	1.28×10^5	2.56×10^5	6.41×10^5	1.27×10^5	1.28×10^6	1.42×10^7

Table 2: Detailed dataset specifications used in the pre-training of masked image modeling.

2.2 Architecture Specifications

We use Swin Transformer V2 [20] as the vision encoder in this study. Thanks to its generality and scalability, we evaluate a series of SwinV2 models with a wide range of model sizes (the number of parameters ranges from $\sim 50\text{M}$ to $\sim 1\text{B}$, and FLOPs range from $\sim 9\text{G}$ to $\sim 190\text{G}$) on multiple downstream tasks. The detailed model specifications are shown in Table 1. We use a new variant SwinV2-g (giant), with number of parameters between SwinV2-L and the 3-billion-parameter SwinV2-G (Giant) used in [20].

2.3 Pre-training Datasets

To study the effect of data size on masked image modeling, we build datasets with different sizes. We use the training set of ImageNet-1K and ImageNet-22K as two large-scale datasets, and randomly sample 10%, 20%, 50% of images in the ImageNet-1K training set as smaller datasets. By default, the images are uniformly sampled from each category. We also consider the sampling strategies could perform differently. To this end, we randomly sample 100 classes from ImageNet-1K as ImageNet-100, and compare it with ImageNet-1K (10%) but find their training loss and fine-tuning performance are almost the same. The details and statistics of all pre-training datasets used in our study are shown in Table 2.

2.4 Pre-training Details

To better compare the performance of models with different amounts of data under the same pre-training length, we use training iterations rather than training epochs and adopt the same hyper-parameters for all models with different sizes during pre-training. The total number of training iterations is in {125K, 250K, 500K} and the batch size is set as 2048 for all experiments. In pre-training stage, we use the same hyper-parameters for all models, and the training details and hyper-parameters of pre-training are summarized in Table 10. Because of the excessive amount of experiments, we follow SimMIM [36] and also use the following two techniques for reducing the experimental overheads: First, we use the step learning rate scheduler in pre-training for sharing the first training step among experiments with different training lengths. The first 7/8 training iterations are the first step and the last 1/8 training iterations are the second step with the learning rate ratio of 0.1 (*i.e.* learning rate is divided by 10 in the second step). Second, we adopt the input image size of 192^2 and set the window size of 12. We improve the SimMIM by normalizing the predicted target according to [12] with a sliding window of 47^2 and observe an improvement of 0.3 on top-1 accuracy of ImageNet-1K for the SwinV2-Large model. The same light data augmentation strategy as SimMIM is used: random resize cropping with a scale range of [0.67, 1], an aspect ratio range of [3/4, 4/3] and a random flipping with probability 0.5.

2.5 Fine-tuning Tasks

To extensively and accurately evaluate the performance of pre-trained models under different pre-training schedulers and datasets, a series of diverse and representative tasks including fine-tuning on ImageNet-1K, fine-grained image classification, object detection, instance segmentation, and semantic segmentation are selected for evaluation.

ImageNet-1K We follow [1] to evaluate the quality of learnt representations by fine-tuning the pre-trained models on ImageNet-1K [8] image classification task, which is the most commonly used scenario and evaluation criterion for pre-trained models [15, 36]. The setting details and fine-tuning hyper-parameters for ImageNet-1K image classification are summarized in Table 11. Different from pre-training, We adopt the image size with 224^2 with window size of 14 in fine-tuning. The AdamW with batch size of 2048, base learning rate of $5e-3$, weight decay of 0.05, β_1 of 0.9 and β_2 of 0.999 are used, and we adopt cosine learning rate scheduler. As larger models are more prone to overfitting, we fine-tune SwinV2-S/B/L for 100 epochs with 20 warm-up epochs and SwinV2-H/g for 50 epochs with 10 warm-up epochs, and decrease the layer decay as the model size increases. In addition, gradient clipping, stochastic depth, label smoothing and data augmentations (*e.g.* random crop, rand erasing [40], rand augment [6], mixup [39], cutmix [37], *etc.*) are also used by following [36].

iNaturalist-18 iNaturalist [32] 2018 is a long-tailed fine-grained image classification dataset. The details and fine-tuning hyper-parameters for iNaturalist 2018 are summarized in Table 12. As fine-tuning in ImageNet-1K, we also use the input image size of 224^2 , window size of 14 and patch size of 4 in iNaturalist 2018. We fine-tune all models for 100 epochs with 20 warm-up epochs, and set layer decay to 0.8, 0.75 and 0.7 for SwinV2-S/B/L, respectively. The AdamW optimizer with cosine learning rate scheduler, batch size of 2048, base learning rate of $1.6e-2$, weight decay of 0.1, β_1 of 0.9 and β_2 of 0.999 are used. In addition, we also adopt stochastic depth, label smoothing, gradient clipping and data augmentations in fine-tuning.

COCO Object Detection and Instance Segmentation [19] The details and fine-tuning hyper-parameters for COCO dataset are summarized in Table 13. We use Mask R-CNN [16]² for evaluation. We set the window size to 14 and patch size to 4. The AdamW optimizer with batch size of 32, base learning rate of $8e-5$, weight decay of 0.05, β_1 of 0.9, β_2 of 0.999 and a step learning rate scheduler (step learning rate ratio of 0.1, step epochs are 27 and 33) are used. In training, the random cropping with crop size of [1024, 1024], large scale jittering with a range of [0.1, 2.0], random horizontal flip with probability 0.5, and stochastic depth regularization are used. In testing, all images are resized to (800, 1333) and keeping the aspect ratio unchanged.

ADE20K Semantic Segmentation [41] The details and fine-tuning hyper-parameters for ADE20K dataset are summarized in Table 14. Following [21], we use UPerNet [35] for evaluation. We set the window size to 20 and the patch size to 4. The AdamW optimizer with batch size of 32, base learning rate searched in a range of [$1e-4$, $3e-4$], weight decay of 0.05, β_1 of 0.9, β_2 of 0.999 and a linear learning rate scheduler with a total of 80K iterations are used. Also, we use the layer decay of 0.95, 0.95, 0.9 for SwinV2-S/B/L, respectively. In training, the random cropping with crop size of [640, 640], scale jittering with a range of [0.5, 2.0], random horizontal flip with probability 0.5, random photometric distortion and stochastic depth regularization of 0.1 are used. In testing, all images are evaluated by sliding window manner, and use the test image size of (2560, 640) and set sliding window stride to 426, following [21, 36].

3 Results and Findings

3.1 Training Length, Data Size and Model Size

We train numerous models with different training lengths, data sizes, and model sizes, and study how these factors affect the performance of masked image modeling. Figure 1 and Figure 2 illustrate the relationship between the training loss, and the validation loss of pre-training³, and the fine-tuning

²Our implementation based on MMDetection [3].

³The validation loss of pre-training is measured on the validation set of ImageNet-1K for all experiments.

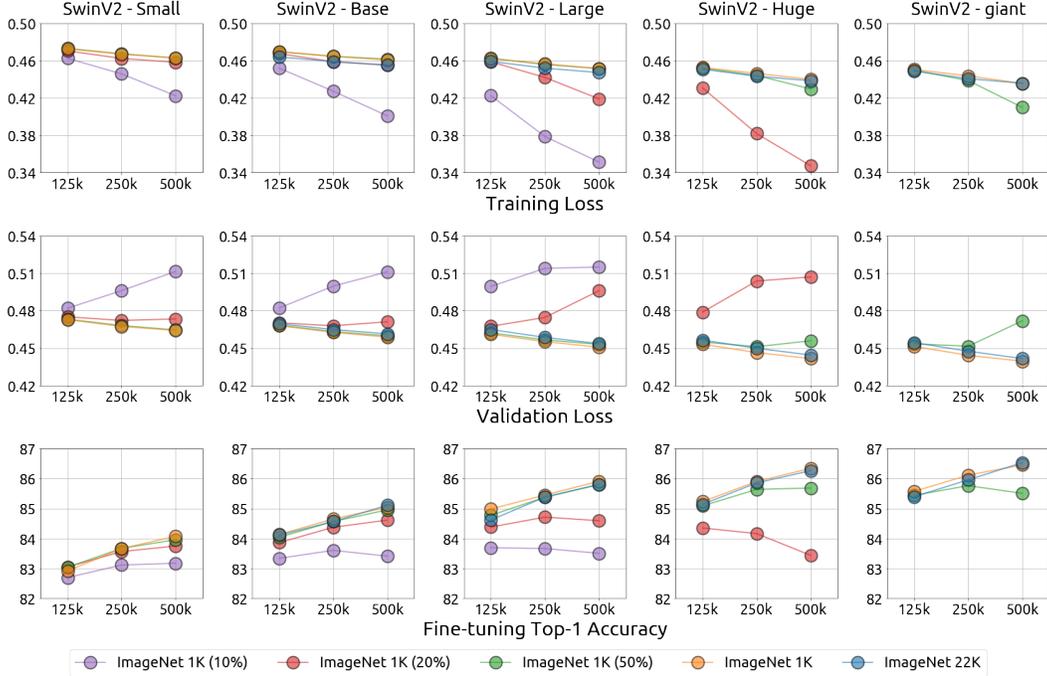


Figure 2: Relationship among training loss, validation loss of pre-training, and fine-tuning performance of ImageNet-1K measured by top-1 accuracy, w.r.t. the training length. *Best viewed in color.*

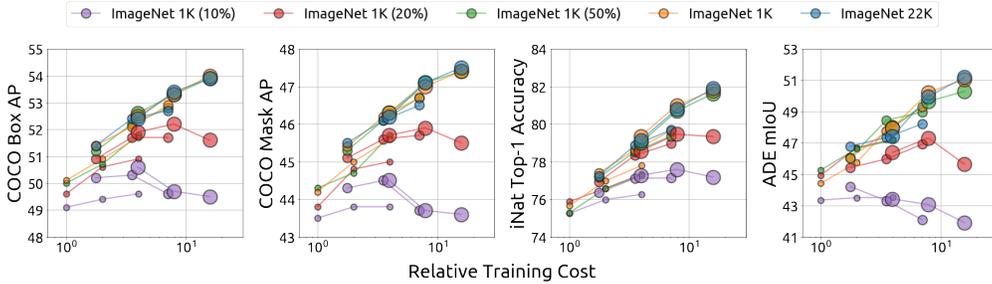


Figure 3: The curves of performances on COCO object detection (a), COCO instance segmentation (b), iNaturalist-18 (c), and ADE20K semantic segmentation (d) w.r.t. the relative training cost. Note that the training cost indicates the pre-training cost. We set the training cost of SwinV2-S for 125K iterations as 1. Bigger circles indicate larger models. *Best viewed in color.*

top-1 accuracy of ImageNet-1K. Based on these extensive experiments, we make the following observations:

Masked image modeling remains demanding for large dataset. When with the high masking rate (*e.g.*, 60% in our work), the masked image modeling is considered a very challenging training objective and has been found to be data efficient by previous literature [11, 22], *i.e.*, a comparable performance can be achieved with small datasets as with large datasets. However, Figure 1 shows that as the training cost increases, the training loss of some models drops significantly, and their validation loss rises significantly, even on using 50% images of ImageNet-1K (*i.e.*, IN1K (50%)), indicating the *overfitting* phenomenon exists. And significant decrease to the fine-tuning performance caused by overfitting could be observed in Figure 2. Moreover, we measure the best fine-tuning performance of each model trained by different training schedulers in Table 3. We find the large models perform even worse than smaller models when small dataset is used for training. For example, the best top-1 accuracy of SwinV2-H with IN1K (20%) is 84.4, worse than the best performance of SwinV2-L by

Model	IN1K (10%)	IN1K (20%)	IN1K (50%)	IN1K (100%)	IN22K (100%)
SwinV2-S	83.2	83.7	84.0	84.1	-
SwinV2-B	83.6	84.6	85.0	85.0	85.1
SwinV2-L	83.7	84.7	85.8	85.9	85.8
SwinV2-H	-	84.4	85.7	86.3	86.3
SwinV2-g	-	-	85.8	86.5	86.5

Table 3: The best fine-tuning performance (top-1 accuracy) of each model with different scales of data on ImageNet-1K image classification.

Model	Iter	IN1K (10%)	IN1K (20%)	IN1K (50%)	IN1K (100%)	IN22K (100%)
SwinV2-S	125K	75.2	75.9	75.3	75.6	-
	250K	76.0	76.6	76.6	77.0	-
	500K	76.3	77.3	77.3	77.8	-
SwinV2-B	125K	76.4	76.9	77.2	77.4	77.2
	250K	77.1	78.3	78.5	78.6	78.9
	500K	77.1	79.0	79.4	79.6	79.7
SwinV2-L	125K	77.3	78.6	79.0	79.4	79.1
	250K	77.6	79.5	80.7	81.0	80.8
	500K	77.2	79.3	81.6	81.8	81.9

Table 4: Results of top-1 accuracy on iNaturalist-18 fine-grained image classification.

0.3. In addition, by comparing the best performance that can be obtained using different sizes of dataset, we find that using more data results in better performance. These observations suggest that masked image modeling does not alleviate the demands of large dataset.

The training length matters. Larger models can benefit from more data at a longer training length. By comparing the performance of models pre-trained by different data sizes (3rd row of Figure 2), we find that the fine-tuning performance of the large models saturates more slowly with the increasing data size compared to the smaller models. For example, the SwinV2-S model pre-trained on IN1K (50%) has a very similar fine-tuning performance to the model pre-trained on IN1K (100%). In comparison, the performance difference between the SwinV2-H model pre-trained on IN1K (50%) and IN1K (100%) is near 0.5, which is a significant gap for ImageNet-1K classification.

Furthermore, a comprehensive observation reveals that the improvements from using more data are not significant under short training lengths. For example, while there is a noticeable performance gap between SwinV2-H trained on IN1K (50%) and IN1K (100%) at a training length of 500K iterations, the gap is less than 0.1 at a training length of 125K iterations. This observation suggests that while larger models can benefit from more data, the training length must also increase at the same time.

Evaluation on more tasks. In addition to ImageNet-1K image classification, we also evaluate the MIM pre-trained SwinV2-S, SwinV2-B and SwinV2-L on iNaturalist-18 fine-grained image classification, ADE20K semantic segmentation, and COCO object detection/segmentation. Figure 3 shows a similar pattern with ImageNet-1K (Figure 1 (right)) that as the training cost increases, some

Model	Iter	IN1K (10%)	IN1K (20%)	IN1K (50%)	IN1K (100%)	IN22K (100%)
SwinV2-S	125K	43.4	44.9	45.3	44.2	-
	250K	43.5	46.7	46.6	45.8	-
	500K	43.5	47.2	47.2	48.3	-
SwinV2-B	125K	44.2	45.4	46.1	46.0	46.8
	250K	43.3	46.0	48.5	47.7	47.3
	500K	42.1	46.9	49.0	49.3	48.2
SwinV2-L	125K	43.4	46.4	48.0	48.0	47.4
	250K	43.1	47.3	49.6	50.2	50.0
	500K	41.9	45.6	50.3	51.1	51.2

Table 5: Results (mIoU) on validation set of ADE20K semantic segmentation.

Model	Iter	IN1K (10%)		IN1K (20%)		IN1K (50%)		IN1K (100%)		IN22K (100%)	
		AP ^{box}	AP ^{mask}								
SwinV2-S	125K	49.1	43.5	49.6	43.8	50.0	44.3	50.1	44.2	-	-
	250K	49.4	43.8	50.6	44.8	50.7	44.7	50.9	45.0	-	-
	500K	49.6	43.8	50.9	44.8	51.8	45.7	51.7	45.6	-	-
SwinV2-B	125K	50.2	44.3	50.9	45.1	51.2	45.3	51.4	45.4	51.4	45.5
	250K	50.3	44.5	51.7	45.6	52.2	46.1	52.1	46.2	52.4	46.1
	500K	49.6	43.7	51.7	45.7	52.8	46.7	52.9	46.7	52.7	46.5
SwinV2-L	125K	50.6	44.5	51.9	45.7	52.6	46.3	52.5	46.3	52.4	46.2
	250K	49.7	43.7	52.2	45.9	53.3	47.1	53.3	47.0	53.4	47.1
	500K	49.5	43.6	51.6	45.5	53.9	47.4	54.0	47.4	53.9	47.5

Table 6: Results of box/mask AP on the validation set of COCO object detection and instance segmentation.

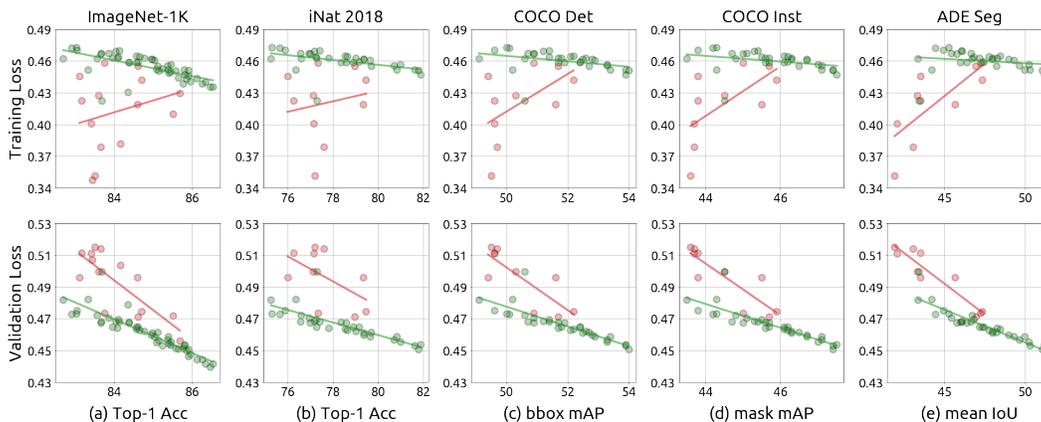


Figure 4: The correlations between pre-training losses (training and validation losses) and the fine-tuning performances. (a) ImageNet-1K image classification; (b) iNat 2018 fine-grained classification; (c) COCO object detection; (d) COCO instance segmentation; (e) ADE20K semantic segmentation. Pre-training losses are highly correlated with fine-tuning performance on all tasks. Red circles are the overfitting models and green circles are non-overfitting models. *Best viewed in color.*

models have significantly performance drop. In addition, as shown in Table 4, 5, and 6, the smaller models rapidly reach saturation as the amount of data increases, while larger models can continuously benefit from more data after sufficient training. These results suggest that the conclusions drawn on ImageNet-1K are broadly applicable to other vision tasks.

3.2 Reconstruction Results of Overfitting and Non-overfitting Models

To better understand the difference between overfitting and non-overfitting models, we visualize the reconstruction results of SwinV2-L that pre-trained on ImageNet1K (10%) and ImageNet1K (100%). Figure 5(a) shows the reconstruction results on the training images from ImageNet1K (10%) dataset, and Figure 5(b) shows the reconstruction results on the images from ImageNet-1K validation set. Based on the reconstruction results on the training images, we observed the overfitting model (*i.e.* SwinV2-L pre-trained on ImageNet1K (10%)) is more like "*remembering*" the masked regions, while the non-overfitting model (*i.e.* SwinV2-L pre-trained on ImageNet1K (100%)) is more like performing "*reasoning*" on the masked regions. For example, the results on the left of the first row in Figure 5(a) show that the overfitting model even predicts the black hair of the dog, but the seen regions only indicate that the dog is white. And the non-overfitting model only predicts the dog with the white hair. Furthermore, we observe that the overfitting model seems to lack the "*reasoning*" ability and has a poorer prediction quality on the images of the validation set compared to the non-overfitting model. For example, the results on the left of the first row in Figure 5(b) show the overfitting model even fails to predict the eyes of the dog.

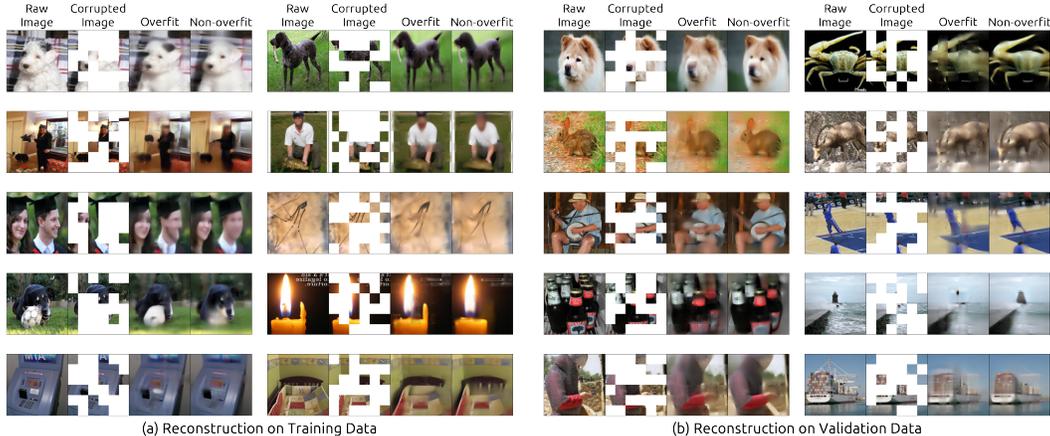


Figure 5: We visualize the reconstruction results of overfitting model (SwinV2-L pre-trained on ImageNet-1K(10%)) and non-overfitting model (SwinV2-L pre-trained on ImageNet-1K(100%)). (a) shows the reconstruction results on the training images from **ImageNet-1K(10%)** dataset, which are jointly contained by the training set of two models. (b) shows the reconstruction results on the validation images from **ImageNet-1K validation set**. Each group contains 4 images from left to right are: the original image, the corrupted images, reconstructed image of overfitting model, and reconstructed image of non-overfitting model.

3.3 Correlation between Pre-training Losses and the Fine-tuning Performance

Evaluating a pre-trained model by its fine-tuned performance on downstream tasks is costly. In supervised pre-training, the validation accuracy is used as the proxy metric to evaluate the quality of the pre-trained models. While in previous studies [5] on other self-supervised learning approaches (*e.g.*, contrastive learning), such a proxy metric is lacking. In this study, we would like to explore whether the pre-training loss in the training of masked image modeling is a good indicator of its fine-tuning performance. We collect all pre-trained models and plot their training and validation loss curves on Figure 4. Interestingly, the correlations between pre-training losses and the fine-tuning performance on multiple tasks could be observed with a *phase transition* around overfitting.

Specifically, the correlation between training loss and fine-tuning performance is negative for the overfitting model (green circles) and positive for the non-overfitting model (red circles). The correlation between validation loss and fine-tuning performance is always negative, but the slope of their linear fit lines ⁴ is significantly different.

In addition, we further analyze the Pearson correlation coefficient between training loss and fine-tuning performance (Table 7), and find the validation loss has stronger *linear correlation* with fine-tuned performance than train loss for all cases, especially for non-overfitting models.

3.4 Effects of Different Sizes of Decoders

We have studied the effects of encoder size from the data scaling perspective. Here, the effects of decoder size are further studied. We pre-train SwinV2-B models with decoder heads of different sizes on IN1K (20%), and Table 8 shows the results. Interestingly, although we find that the heavier decoder has lower training loss and higher validation loss than the linear decoder, indicating a more severe overfitting issue. But there is no decrease in its fine-tuning performance on ImageNet-1K than the linear decoder. This experiment shows that the decoder behaves very differently from the encoder, and we speculate that this is because the decoder "blocks" the damage to the encoder from overfitting.

3.5 Impact of Different Dataset Sampling Strategies

We study different dataset sampling strategies by comparing the training behavior and fine-tuned performance of models pre-trained on IN1K (10%) and IN100. In IN1K (10%), the images are

⁴The least squares method is used for linear fit.

	w/ train loss	w/ val loss		w/ train loss	w/ val loss
overfit	+0.26	-0.79	overfit	+0.17	-0.54
non-overfit	-0.64	-0.90	non-overfit	-0.46	-0.78
(a) ImageNet-1K			(b) iNaturalist 2018		
	w/ train loss	w/ val loss		w/ train loss	w/ val loss
overfit	+0.54	-0.81	overfit	+0.62	-0.86
non-overfit	-0.35	-0.83	non-overfit	-0.31	-0.85
(c) COCO Object Detection			(d) COCO Instance Segmentation		
	w/ train loss	w/ val loss		w/ train loss	w/ val loss
overfit	+0.75	-0.91	overfit	+0.75	-0.91
non-overfit	-0.14	-0.90	non-overfit	-0.14	-0.90
(e) ADE-20K Semantic Segmentation					

Table 7: Pearson correlation coefficients between pre-training losses (training and validation losses) and fine-tuning performances on five downstream tasks.

Encoder	Decoder	# Params	Training loss	Validation loss	Top-1 accuracy
SwinV2-B	linear	90.0M	0.46	0.47	84.4
SwinV2-B	4-blocks	140.4M	0.44	0.48	84.4
SwinV2-B	8-blocks	190.8M	0.41	0.50	84.5

Table 8: Results of different decoders, including converged training and validation losses of MIM pre-training, and fine-tuning performance (top-1 accuracy) on ImageNet-1K image classification.

uniformly sampled from each category, and we randomly sample 100 categories from ImageNet-1K as IN100. Experiments are conducted on SwinV2-L with 500K training iterations. Table 9 shows the training loss, validation loss and fine-tuning top-1 accuracy of ImageNet-1K. For the two models pre-trained on IN1K (10%) and IN100, all three metrics are very similar. Figure 6 further illustrates the training dynamics of the two models, and we find both their training loss curves and validation loss curves are almost overlapping. These results show the disparity caused by different dataset sampling strategies is minor.

4 Related Work

Masked Image Modeling Masked Image Modeling learns representations by reconstructing the masked content of images, and its early exploration can be traced back to context encoder [25] and denoising autoencoder [34]. Recently, iGPT [4], BEiT [1], MAE [15] and SimMIM [36] recall this approach on training vision transformer. iGPT [4] sequentially predicted the pixels by auto-regressive manner. BEiT [1] proposed to predict the discrete visual tokens. MAE [15] and SimMIM [36] concurrently find predicting the raw pixels with a high masking ratio can work well. In this work, we use SimMIM as the default masked image modeling approach, because of its simplicity and no restrictions on the architecture of vision encoder like MAE.

Vision Transformer Transformer [33] was first applied to natural language processing and became the dominant architecture, and has recently attracted a lot of attention in computer vision. The pioneering work ViT [10] first shows that the transformer architecture works well in image classification when trained on large amounts of data. DeiT [31] proposed a better training recipe based on ViT and demonstrated that vision Transformer has promising performance when only using ImageNet-1K dataset. Swin Transformer [21] improves plain ViT by inducing the hierarchical architecture and non-overlapping local attention and successfully demonstrates the effectiveness of vision transformer on a wide range of vision tasks. Swin Transformer V2 [20] further addresses the training stability

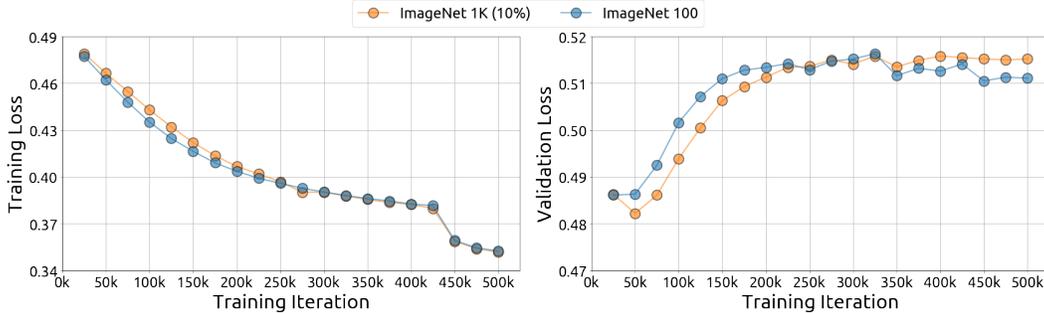


Figure 6: The training loss and validation loss of MIM pre-training with different dataset sampling strategies, ImageNet-1K (10%) and ImageNet-100. *Best viewed in color.*

Dataset	# Images	# Classes	Training loss	Validation loss	Top-1 Accuracy
IN1K (10%)	1.28×10^5	1000	0.351	0.515	83.5
IN100	1.27×10^5	100	0.352	0.511	83.4

Table 9: Results on different dataset sampling strategies (ImageNet-1K (10%) and ImageNet-100), include converged training and validation losses of MIM pre-training, and fine-tuning performance (top-1 accuracy) on ImageNet-1K image classification.

issue of [21] in model scaling and illustrates better performance than the original Swin Transformer, and thus we use it as the default vision encoder in this work.

Scaling Vision Models Many works [29, 28, 38, 20] examine how to scale vision models, but most are more concerned with exploring the perspective of model architecture designs. For example, EfficientNet [29] extensively studied how model width, model depth and input resolution affect the convolutional neural networks; [28] proposed to scale vision model with sparse mixture-of-expert; [38] and [20] studied how to scale ViT and Swin Transformer, respectively.

Only a few works explored the perspective of data scaling under the pre-training fine-tuning paradigm. BiT [18] revisited the supervised pre-training on a wide range of data scales up to 1M images. SEER [14] studied the effectiveness of data scaling in the contrastive learning framework with up to one billion images. Recently, SplitMask [11] find that masked image modeling is robust to the size of pre-training data and challenges the data scaling capability of masked image modeling, which is most relevant to our work.

5 Conclusion

In our work, we systematically study the data scaling capability of masked image modeling at different model sizes and training lengths. Based on the extensive experiments, we demonstrate that masked image modeling is not only a model scalable learner but also a data scalable learner, which challenges the conclusion of previous literature that a large dataset may not be necessary in masked image modeling. The reason behind this is that they overlooked a key factor, namely training length. In addition, a strong correlation between the validation loss of masked image modeling and the fine-tuning performance is observed. This observation suggests that validation loss can be considered as a good proxy metric for evaluating pre-trained models, and makes it possible to reduce the experimental overhead of measuring models by fine-tuning.

While these findings deepen our understanding of masked image modeling in data scaling angles and can facilitate future research, our study still has limitations. First, the maximum model size used in our study reaches only one billion parameters, which we speculate leaves the overfitting phenomenon on the ImageNet-1K dataset unobserved; Second, there is a lack of research on the effect of encoder specifications (*e.g.*, depth and width) on data scaling. Third, our study does not involve the study angle of data augmentation which is a common technique to alleviate data scarcity and overfitting.

References

- [1] Bao, H., Dong, L., and Wei, F. (2021). Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*.
- [2] Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. (2020). Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.
- [3] Chen, K., Wang, J., Pang, J., Cao, Y., Xiong, Y., Li, X., Sun, S., Feng, W., Liu, Z., Xu, J., et al. (2019). Mmdetection: Open mmlab detection toolbox and benchmark. *arXiv preprint arXiv:1906.07155*.
- [4] Chen, M., Radford, A., Child, R., Wu, J., and Jun, H. (2020a). Generative pretraining from pixels. *Advances in Neural Information Processing Systems*.
- [5] Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. (2020b). A simple framework for contrastive learning of visual representations. *ICML*.
- [6] Cubuk, E. D., Zoph, B., Shlens, J., and Le, Q. V. (2020). 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.
- [7] Dai, Z., Liu, H., Le, Q. V., and Tan, M. (2021). Coatnet: Marrying convolution and attention for all data sizes.
- [8] 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.
- [9] 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*.
- [10] 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*.
- [11] El-Nouby, A., Izacard, G., Touvron, H., Laptév, I., Jegou, H., and Grave, E. (2021). Are large-scale datasets necessary for self-supervised pre-training? *arXiv preprint arXiv:2112.10740*.
- [12] Fang, Y., Dong, L., Bao, H., Wang, X., and Wei, F. (2022). Corrupted image modeling for self-supervised visual pre-training. *arXiv preprint arXiv:2202.03382*.
- [13] Fedus, W., Zoph, B., and Shazeer, N. (2021). Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity.
- [14] Goyal, P., Caron, M., Lefaudeux, B., Xu, M., Wang, P., Pai, V., Singh, M., Liptchinsky, V., Misra, I., Joulin, A., et al. (2021). Self-supervised pretraining of visual features in the wild. *arXiv preprint arXiv:2103.01988*.
- [15] 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*.
- [16] He, K., Gkioxari, G., Dollár, P., and Girshick, R. (2017). Mask r-cnn. In *ICCV*, pages 2961–2969.
- [17] Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., Gray, S., Radford, A., Wu, J., and Amodei, D. (2020). Scaling laws for neural language models.
- [18] Kolesnikov, A., Beyer, L., Zhai, X., Puigcerver, J., Yung, J., Gelly, S., and Houlsby, N. (2020). Big transfer (bit): General visual representation learning. In *European conference on computer vision*, pages 491–507. Springer.
- [19] 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 *European conference on computer vision*, pages 740–755. Springer.

- [20] 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*.
- [21] 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*.
- [22] Liu, Z., Ning, J., Cao, Y., Wei, Y., Zhang, Z., Lin, S., and Hu, H. (2021c). Video swin transformer.
- [23] Microsoft (2020). Turing-nlg: A 17-billion-parameter language model by microsoft.
- [24] Microsoft (2021). Using deepspeed and megatron to train megatron-turing nlg 530b, the world’s largest and most powerful generative language model.
- [25] Pathak, D., Krahenbuhl, P., Donahue, J., Darrell, T., and Efros, A. A. (2016). Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2536–2544.
- [26] Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., and Sutskever, I. (2019). Language models are unsupervised multitask learners.
- [27] Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu, P. J. (2020). Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67.
- [28] Riquelme, C., Puigcerver, J., Mustafa, B., Neumann, M., Jenatton, R., Pinto, A. S., Keysers, D., and Houlsby, N. (2021). Scaling vision with sparse mixture of experts. *arXiv preprint arXiv:2106.05974*.
- [29] Tan, M. and Le, Q. (2019). Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR.
- [30] Tong, Z., Song, Y., Wang, J., and Wang, L. (2022). Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training. *arXiv preprint arXiv:2203.12602*.
- [31] Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., and Jégou, H. (2021). Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning*, pages 10347–10357. PMLR.
- [32] 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.
- [33] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.
- [34] Vincent, P., Larochelle, H., Bengio, Y., and Manzagol, P.-A. (2008). Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine learning*, pages 1096–1103.
- [35] Xiao, T., Liu, Y., Zhou, B., Jiang, Y., and Sun, J. (2018). Unified perceptual parsing for scene understanding. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 418–434.
- [36] 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*.
- [37] 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.
- [38] Zhai, X., Kolesnikov, A., Houlsby, N., and Beyer, L. (2021). Scaling vision transformers.

- [39] Zhang, H., Cisse, M., Dauphin, Y. N., and Lopez-Paz, D. (2017). mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*.
- [40] 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.
- [41] Zhou, B., Zhao, H., Puig, X., Xiao, T., Fidler, S., Barriuso, A., and Torralba, A. (2018). Semantic understanding of scenes through the ade20k dataset. *International Journal on Computer Vision*.

A Hyper-parameters and training details

We illustrate the training details of pre-training and fine-tuning for different tasks and different models. Table 10 presents pre-training details. Table 11 presents the fine-tuning details on ImageNet-1K image classification. Table 12 presents the fine-tuning details on iNaturalist 2018. Table 13 presents the fine-tuning details on COCO dataset. Table 14 presents the fine-tuning details on ADE20K dataset.

B Training dynamics of masked image modeling

We show the training curves and validation curves of different models trained by masked image modeling to better illustrate the training dynamics. In Figure 7, each row presents the training and validation loss curves for training with the same model but different dataset. The training loss is computed on its corresponding training dataset and the validation loss is computed on the ImageNet-1K validation set. We make the following observations: First, all models have the overfitting issues when using small datasets. Second, for the non-overfitting cases, the training and validation losses are similar using different sizes of datasets for training. In Figure 8, the training/validation loss curves of different models but using the same training dataset are presented at each row. We make the following observations: First, larger models have lower training losses than smaller models for all datasets. Second, the validation loss of the larger model is lower than the smaller model in the non-overfitting cases but higher than the smaller model in the over-fitting cases.

Pre-training setting of all models	
Input size	192 ²
Window size	12
Patch size	4
Mask patch size	32
Mask ratio	0.6
Training iterations	125,000 / 250,000 / 500,000
Batch size	2048
Optimizer	AdamW
Init. learning rate	4e-4
Weight decay	0.05
Adam ϵ	1e-8
Adam β	(0.9, 0.999)
Learning rate scheduler	Step
Step learning rate ratio	0.1
Step iterations	109,375 / 218,750 / 437,500
Warm-up iterations	6250
Gradient clipping	5.0
Stochastic depth	0.1
Rand crop scale	[0.67, 1]
Rand resize ratio	[3/4, 4/3]
Rand horizontal flip	0.5
Reconstruction target	Norm. with sliding window [12]
Norm. patch size	47

Table 10: Details and hyper-parameters for SimMIM pre-training.

Hyperparameters	SwinV2				
	Small(S)	Base(B)	Large(L)	Huge(H)	giant(g)
Input size			224 ²		
Window size			14		
Patch size			4		
Training epochs	100	100	100	50	50
Warm-up epochs	20	20	20	10	10
Layer decay	0.8	0.75	0.7	0.65	0.65
Batch size			2048		
Optimizer			AdamW		
Base learning rate			5e-3		
Weight decay			0.05		
Adam ϵ			1e-8		
Adam β			(0.9, 0.999)		
Learning rate scheduler			cosine		
Gradient clipping			5.0		
Stochastic depth			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 11: Details and hyper-parameters for ImageNet-1K fine-tuning.

Hyperparameters	SwinV2		
	Small(S)	Base(B)	Large(L)
Input size			224 ²
Window size			14
Patch size			4
Training epochs			100
Warm-up epochs			20
Layer decay	0.8	0.75	0.7
Batch size			2048
Optimizer			AdamW
Base learning rate			1.6e-2
Weight decay			0.1
Adam ϵ			1e-8
Adam β			(0.9, 0.999)
Learning rate scheduler			cosine
Gradient clipping			5.0
Stochastic depth			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 12: Details and hyper-parameters for iNaturalist 2018 fine-tuning.

Hyperparameters	SwinV2		
	Small(S)	Base(B)	Large(L)
Detector	Mask R-CNN		
Window size	14		
Patch size	4		
Training input size	(1024, 1024)		
Testing input size	(800, 1333)		
Training epochs	36		
Warm-up iterations	500		
Batch size	32		
Optimizer	AdamW		
Base learning rate	8e-5		
Weight decay	0.05		
Adam ϵ	1e-8		
Adam β	(0.9, 0.999)		
Learning rate scheduler	Step		
Step learning rate ratio	0.1		
Step epochs	(27, 33)		
Stochastic depth	0.1	0.1	0.2
Rand horizontal flip	0.5		
Scale Jittering	[0.1, 2.0]		

Table 13: Details and hyper-parameters for fine-tuning on the COCO dataset.

Hyperparameters	SwinV2		
	Small(S)	Base(B)	Large(L)
Architecture	UPerNet		
Window size	20		
Patch size	4		
Training input size	(640, 640)		
Test input size	(640, 2560)		
Slide test stride	(426, 426)		
Training iterations	80,000		
Warm-up iterations	750		
Layer decay	0.95	0.95	0.9
Batch size	32		
Optimizer	AdamW		
Base learning rate	[1e-4, 3e-4]		
Weight decay	0.05		
Adam ϵ	1e-8		
Adam β	(0.9, 0.999)		
Learning rate scheduler	Linear		
Stochastic depth	0.1		
Rand horizontal flip	0.5		
Scaling Jittering	[0.5, 2.0]		
Photo Metric Distortion	✓		

Table 14: Details and hyper-parameters for fine-tuning on the ADE20K dataset.

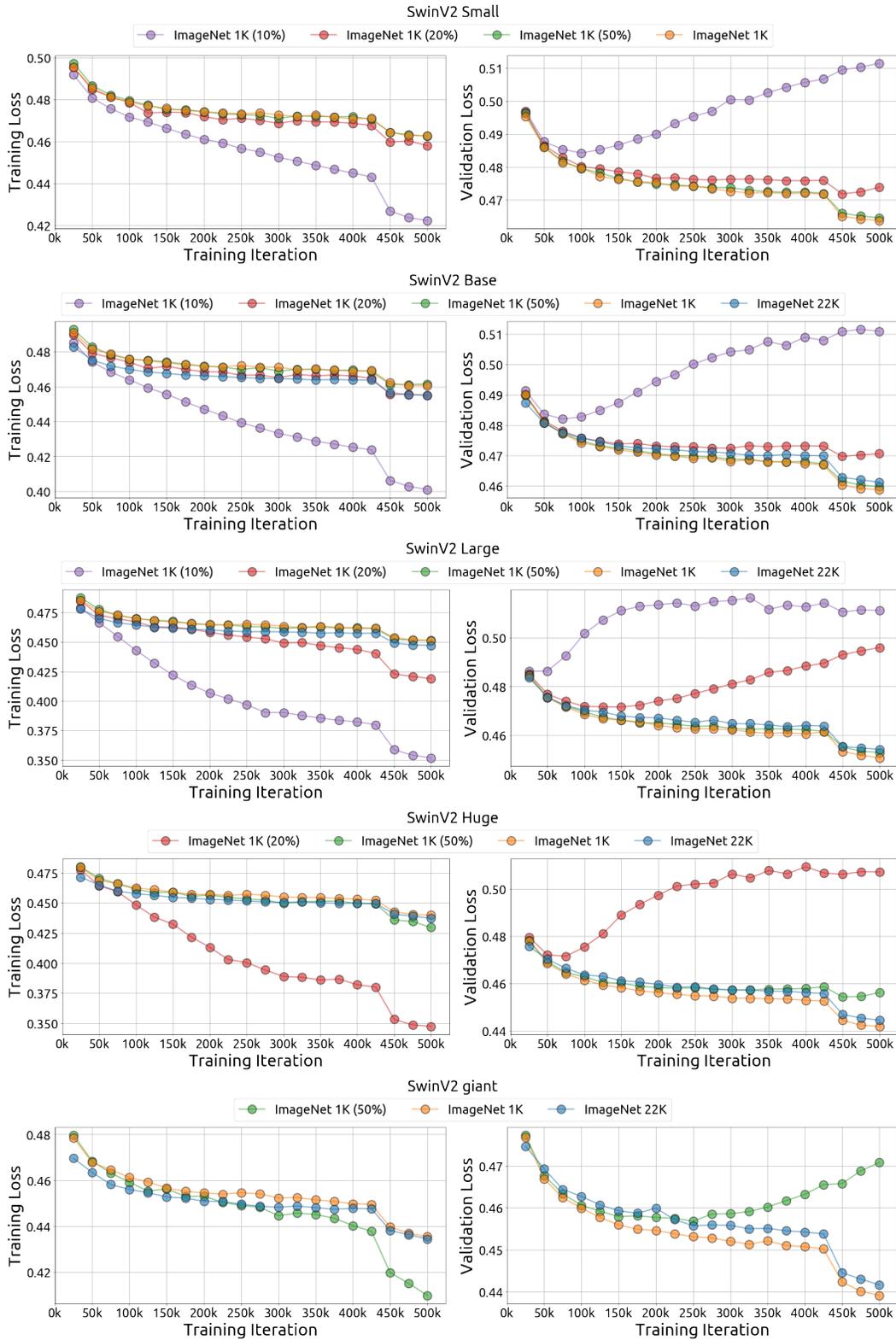


Figure 7: Each row presents the training and the validation loss curves for training with the same model (e.g., SwinV2 giant at the last row) but different datasets. The training loss is computed on its corresponding training dataset, and the validation loss is computed on the ImageNet-1K validation set. *Best viewed in color.*

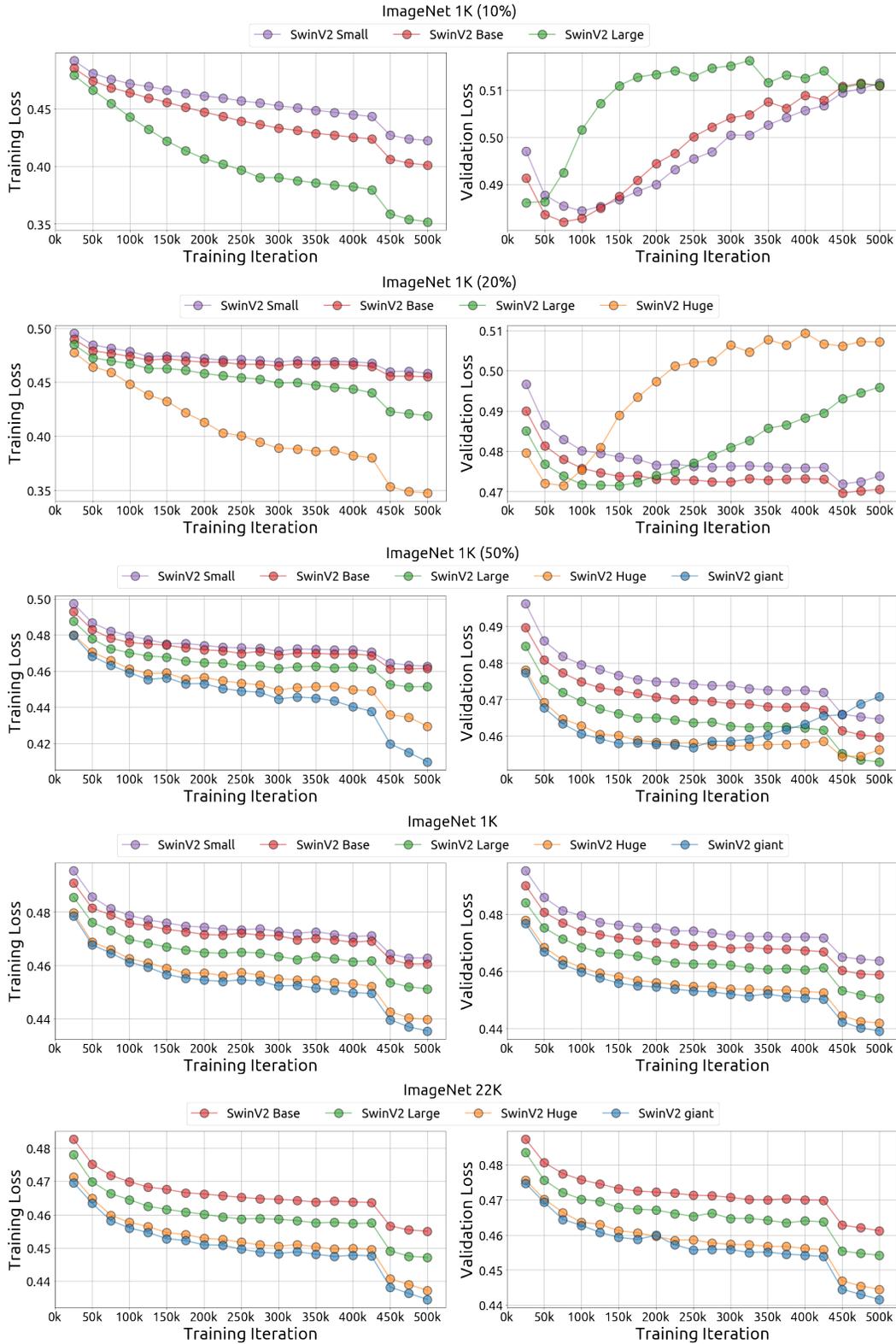


Figure 8: Each row presents the training and the validation loss curves for training with the same loss curves for training with the same dataset (e.g., ImageNet22K at the last row) but different models. The training loss is computed on its corresponding training dataset, and the validation loss is computed on the ImageNet-1K validation set. *Best viewed in color.*