

ResFormer: Scaling ViTs with Multi-Resolution Training

Rui Tian^{1,2} Zuxuan Wu^{1,2†} Qi Dai³ Han Hu³ Yu Qiao⁴ Yu-Gang Jiang^{1,2}

¹Shanghai Key Lab of Intell. Info. Processing, School of CS, Fudan University

²Shanghai Collaborative Innovation Center of Intelligent Visual Computing

³Microsoft Research Asia

⁴Shanghai AI Laboratory

Abstract

Vision Transformers (ViTs) have achieved overwhelming success, yet they suffer from vulnerable resolution scalability, i.e., the performance drops drastically when presented with input resolutions that are unseen during training. We introduce, ResFormer, a framework that is built upon the seminal idea of multi-resolution training for improved performance on a wide spectrum of, mostly unseen, testing resolutions. In particular, ResFormer operates on replicated images of different resolutions and enforces a scale consistency loss to engage interactive information across different scales. More importantly, to alternate among varying resolutions effectively, especially novel ones in testing, we propose a global-local positional embedding strategy that changes smoothly conditioned on input sizes. We conduct extensive experiments for image classification on ImageNet. The results provide strong quantitative evidence that ResFormer has promising scaling abilities towards a wide range of resolutions. For instance, ResFormer-B-MR achieves a Top-1 accuracy of 75.86% and 81.72% when evaluated on relatively low and high resolutions respectively (i.e., 96 and 640), which are 48% and 7.49% better than DeiT-B. We also demonstrate, moreover, ResFormer is flexible and can be easily extended to semantic segmentation, object detection and video action recognition. Code is available at <https://github.com/ruitian12/resformer>.

1. Introduction

The strong track record of Transformers in a multitude of Natural Language Processing [60] tasks has motivated an extensive exploration of Transformers in the computer vision community. At its core, Vision Transformers (ViTs) build upon the multi-head self-attention mechanisms for feature learning through partitioning input images into patches of identical sizes and processing them as sequences

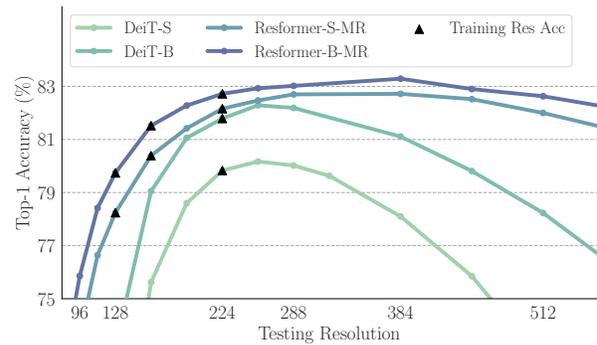


Figure 1. Comparisons between ResFormer and vanilla ViTs. ResFormer achieves promising results on a wide range of resolutions.

for dependency modeling. Owing to their strong capabilities in capturing relationships among patches, ViTs and their variants demonstrate prominent results in versatile visual tasks, e.g., image classification [41, 57, 73, 78], object detection [4, 35, 63], vision-language modeling [29, 46, 62] and video recognition [3, 34, 42, 72].

While ViTs have been shown effective, it remains unclear how to scale ViTs to deal with inputs with varying sizes for different applications. For instance, in image classification, the *de facto* training resolution of 224 is commonly adopted [41, 57, 58, 73]. However, among works in pursuit of reducing the computational cost of ViTs [44, 49], shrinking the spatial dimension of inputs is a popular strategy [7, 37, 64]. On the other hand, fine-tuning with higher resolutions (e.g., 384) is widely used [16, 41, 55, 58, 67, 70] to produce better results. Similarly, dense prediction tasks such as semantic segmentation and object detection also require relatively high resolution inputs [1, 35, 40, 63].

Despite of the necessity for both low and high resolutions, limited effort has been made to equip ViTs with the ability to handle different input resolutions. Given a novel resolution that is different from that used during training, a common practice adopted for inference is to keep the patch size fixed and then perform bicubic interpolation on positional embeddings directly to the corresponding scale. As

[†]Corresponding author.

Note that we use resolution, scale and size interchangeably.

shown in Sec. 3, while such a strategy is able to scale ViTs to relatively larger input sizes, the results on low resolutions plunge sharply. In addition, significant changes between training and testing scales also lead to limited results (*e.g.*, DeiT-S trained on a resolution of 224 degrades by 1.73% and 7.2% when tested on 384 and 512 respectively).

Multi-resolution training, which randomly resizes images to different resolutions, is a promising way to accommodate varying resolutions at test time. While it has been widely used by CNNs for segmentation [23], detection [25] and action recognition [66], generalizing such an idea to ViTs is challenging and less explored. For CNNs, thanks to the stacked convolution design, all input images, regardless of their resolutions, share the same set of parameters in multi-resolution training. For ViTs, although it is feasible to share parameters for all samples, bicubic interpolations of positional embeddings, which are not scale-friendly, are still needed when iterating over images of different sizes.

In this paper, we posit that positional embeddings of ViTs should be adjusted smoothly across different scales for multi-resolution training. The resulting model then has the potential to scale to different resolutions during inference. Furthermore, as images in different scales contain objects of different sizes, we propose to explore useful information across different resolutions for improved performance in a similar spirit to feature pyramids, which are widely used in hierarchical backbone designs for both image classification [25, 41] and dense prediction tasks [23, 24, 38].

To this end, we introduce ResFormer, which takes in inputs as multi-resolution images during training and explores multi-scale clues for better results. Trained in a single run, ResFormer is expected to generalize to a large span of testing resolutions. In particular, given an image during training, ResFormer resize it to different scales, and then use all scales in the same feed-forward process. To encourage information interaction among different resolutions, we introduce a scale consistency loss, which bridges the gap between low-resolution and high-resolution features by self-knowledge distillation. More importantly, to facilitate multi-resolution training, we propose a global-local positional embedding strategy, which enforces parameter sharing and changes smoothly across different resolutions with the help of convolutions. Given a novel resolution at testing, ResFormer dynamically generates a new set of positional embeddings and performs inference.

To validate the efficacy of ResFormer, we conduct comprehensive experiments on ImageNet-1K [14]. We observe that ResFormer makes remarkable gains compared with vanilla ViTs which are trained on single resolution. Given the testing resolution of 224, ResFormer-S-MR trained on resolutions of 128, 160 and 224 achieves a Top-1 accuracy of 82.16%, outperforming the 224-trained DeiT-S [57] by 2.24%. More importantly, as illustrated in Fig. 1, Res-

Former surpasses DeiT by a large margin on unseen resolutions, *e.g.*, ResFormer-S-MR outperforms DeiT-S by 6.67% and 56.04% when tested on 448 and 80 respectively. Furthermore, we also validate the scalability of ResFormer on dense prediction tasks, *e.g.*, ResFormer-B-MR achieves 48.30 mIoU on ADE20K [80] and 47.6 AP^{box} on COCO [39]. We also show that ResFormer can be readily adapted for video action recognition with different sizes of inputs via building upon TimeSFormer [3].

2. Related Work

Scaling Vision Models. Many studies in recent literature [40, 51, 56, 75] discuss how to scale vision models, with most of them focusing on the capacity of deep neural networks. For instance, EfficientNet [56] studies how model width, depth and input resolution affect convolutional neural networks. RegNet [47] designs manual designing space for CNNs and finds simple linear correlation between the search space (*e.g.* width) and performance. ResNet-RS [2] presents how different scaling strategies on depth and input resolution can affect the model capacity.

Recent approaches have investigated scaling of transformers [40, 75]. For example, V-Moe [51] scales vision transformers to large model sizes with sparse mixture-of-experts. Several studies [18, 20, 32, 71] explore the aspect of data scaling under self-supervised framework, *i.e.*, how data sizes affect the model performance. In contrast, limited effort has been made towards the scaling abilities of models towards input resolutions. Attempts have been made by Liu *et al.* [40] to scale up to larger resolutions, while neglecting lower resolutions. Instead, our work takes the initiative to scale ViTs to various resolutions, both lower and higher, satisfying the practical needs from varied visual tasks.

Positional embedding. The self-attention architecture is clueless about spatial relationships among patches. Therefore, to overcome permutation-invariance, various positional embedding strategies have been proposed to enable Transformers to perceive the sequence order of input tokens. Absolute positional embeddings (APE) infuse global spatial information into Transformers, *e.g.*, sine-cosine APE [60] proposed for NLP tasks and learned APE adopted in the vanilla Vision Transformer [16]. Meanwhile, efficacy of relative position embeddings (RPB) is widely validated in both language [13, 52] and vision tasks [5, 40, 68]. For instance, Wu *et al.* replaces APE with the relative strategy of iRPE [68] for performance gains on classification and detection. ConViT [17] also suggests that adding gated relative positional embeddings to self-attention blocks brings about soft convolutional inductive biases. Moreover, dynamic positional embeddings are introduced to model local information from input tokens, *e.g.*, Twins [9] adopt conditional positional embeddings (CPE) [10] implemented by

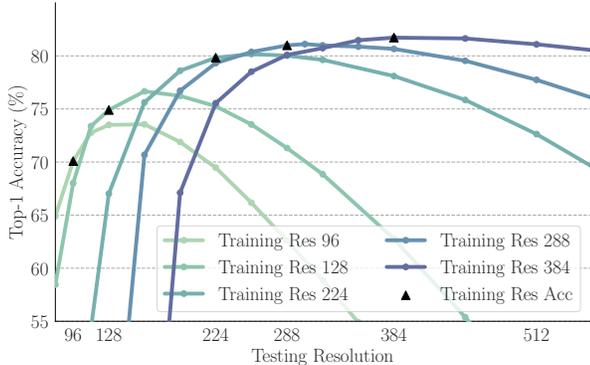


Figure 2. Top-1 accuracy of DeiT-S trained with 5 different resolutions and tested on resolutions varying from 80 to 576. During testing, we follow the common pre-processing steps (*i.e.*, `Resize` and `CenterCrop` in Pytorch implementation) and set the cropping rate to 0.875.

convolutions. In order to improve performance and scalability simultaneously, we propose to inject spatial embeddings in ResFormer from both from global and local perspectives.

Multi-scale training. In early CNNs, multi-scale data augmentations [53] are employed for image classification by randomly sampling training images from a certain range of scales. Later in dense prediction tasks, multi-scale training and testing become a widely-adopted paradigm [4, 33, 54]. In addition, the idea has also been explored in action recognition. Wu *et al.* introduce a multigrid strategy [66], which enables efficient training by sampling data with different grids of temporal span, spatial span and temporal stride. Video ResKD [43] achieves excellent efficiency by employing high-resolution features from large models as teachers to improve low-resolution performance.

Most approaches in multi-scale training rely on CNNs, as convolutions can be readily applied to varying sizes of inputs. In contrast, vanilla ViTs are equipped with tokens of fixed-dimension, other related attempts lay emphasis on multi-scale spatial dimension of features instead of input [19, 21] or perform in an unsupervised way [48], yet limited effort has been made to explore multi-resolution supervised training for ViTs. In this paper, we make the first step to investigate such a strategy for ViTs which not only leads to good performance on training resolutions but can also generalize towards novel resolutions.

3. Resolution Generalization

In this section, we conduct a set of pilot experiments to show the scalability of ViTs towards different resolutions.

Generalizing to different resolutions. As revealed by previous work [59], CNNs suffer from distribution shifts between training and testing due to different pre-processing methods, *i.e.*, the *de facto* “random resizing and cropping”

strategy for training and “center cropping” for testing result in different distributions of cropped regions in images. For ViTs, theoretically, the discrepancy persists since the same pre-processing strategies of training and testing are employed. However, there lacks a comprehensive study on how ViTs behave towards input scales varied from the training process. To this end, we feed pre-trained ViTs with testing samples of varying sizes. In particular, we instantiate ViT models with DeiT-S [57] and initialize the model with weights pre-trained on ImageNet-1K. We then fine-tune the model on a resolution of 96, 128, 224, 288 and 384 respectively.¹ These derived models are then tested on a broad spectrum of resolutions. Following [41, 57], we simply resize the position embeddings with bicubic interpolation on different testing resolutions. The results are shown in Fig. 2. We observe the following trends for scaling up or down:

- **Scaling down:** All models undergo severe performance drop when directly adapted to small-scale inputs, especially for ones pre-trained on larger resolutions. For example, the Top-1 accuracy of DeiT-S with a training resolution of 384 decreases by 6.18% when tested on 224 and even plummets below 30% when the testing resolution is further reduced to 160.
- **Scaling up:** Ideally, increasing testing resolutions results in improved accuracy, which is also suggested as byproduct of train-test distribution discrepancy in [59]. However, models yield unsatisfactory performance when gap enlarges. *e.g.*, DeiT-S with a low training resolution of 128 stops growing in accuracy when the testing resolution reaches 256. It achieves a Top-1 accuracy of 75.26% with testing resolution set to 224, which is 4.57% lower than model trained with a resolution of 224, directly.

Above all, ViTs are vulnerable to resolution discrepancies between training and testing, particularly when evaluated on low-resolution inputs. This motivates us to equip ViTs with scalability towards a wide range of test resolutions so as to meet the need of versatile applications.

4. Method

Our goal is to train a vision transformer that not only performs well on resolutions the network has seen during training, but more importantly it is able to adapt to a wide range of unseen resolutions without significant performance drop during testing. To this end, we first introduce a resolution scaling transformer in Sec. 4.1, ResFormer, which operates on input samples of multiple resolutions in the training stage. Since the size of objects varies in different scales, we also introduce a scale consistency loss to fully explore information from all resolutions for improved accuracy. Furthermore, as mentioned in Sec. 3, directly interpolating positional embeddings to unseen resolutions during inference

¹Please refer to Appendix B.1 for detailed fine-tuning setup.

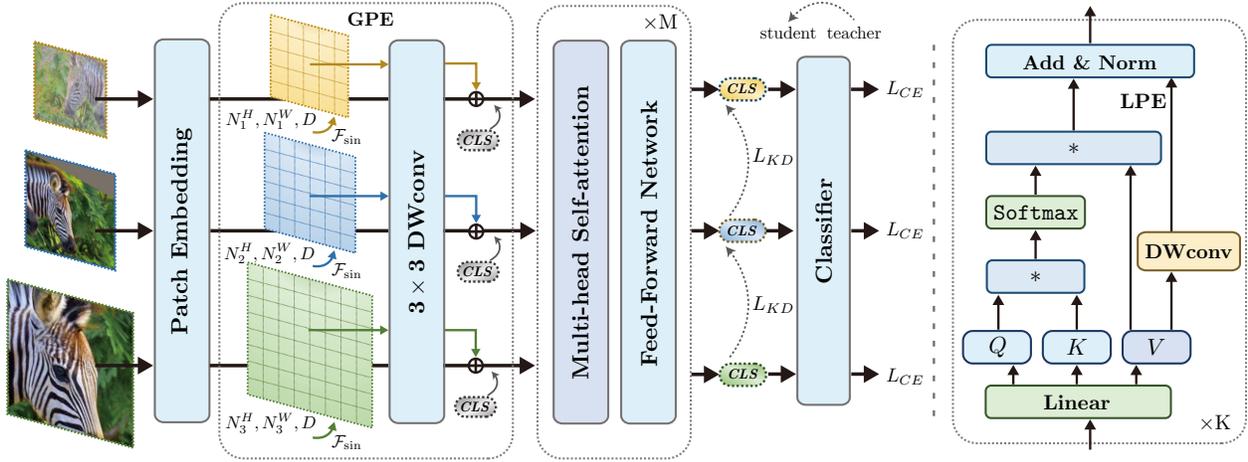


Figure 3. **Left:** The overview of ResFormer framework. **Right:** The pipeline of generating local positional embedding.

produces unsatisfactory results. To mitigate this issue, ResFormer builds upon carefully designed global-local positional embeddings, which are generated conditioned on input resolutions, as will be described in Sec. 4.2.

4.1. ResFormer

Following the vanilla ViT [16], given an input image X whose height and width are H and W , respectively, we first split it into $N^H \times N^W$ patches, where the patch size is set to t and $N_H = H/t, N_W = W/t$. Each image patch is projected into a D -dimension feature by patch embedding and is denoted as a ‘‘token’’. Subsequently, a global class token cls is concatenated with image tokens before they are fed into Transformer blocks.

Unlike standard ViTs operating on single-scale images, ResFormer takes inputs of different resolutions during training so as to better model objects of varying sizes in different scales and generalize better during inference. More specifically, as shown in Fig. 3, we replicate a given training image for r times, where r denotes the number of resolutions used. For the i -th data replica, resizing and cropping operations² are applied to obtain a training sample X^i with a spatial size of $3 \times H_i \times W_i$. Afterwards, we apply random pre-processing strategies involved in ViTs training paradigm³ on each scale of inputs separately. As a result, one mini-batch is composed of groups of multi-resolution inputs sharing identical labels, which is roughly equivalent to extending the base batch size by r times. In addition, for the input sample X^i , we feed the global class token output of the last transformer block as inputs into classification head to compute final predictions Y^i . Naturally, the classification losses

²In practice, we realize it with `RandomResizedCrop` in PyTorch.

³Random pre-processing includes Auto-Augment [11], RandAugment [12], random erasing [79], MixUp [76] and CutMix [74].

can be written as:

$$L(\Theta) = \mathbb{E}_{(X, T) \sim \mathcal{D}} \sum_{i=1}^r L_{CE}(Y^i; X^i, T, \Theta), \quad (1)$$

where T denotes the ground-truth label and L_{CE} represents the cross-entropy loss. In addition, \mathcal{D} and Θ denote the training set and the parameters of the network, respectively.

Scale consistency loss. Given that larger inputs generally produce better recognition results compared to their smaller counterparts, we take advantage of knowledge distillation through enforcing consistencies among different resolutions. In particular, we use a smooth l_1 loss with feature whitening [65], denoted as L_{KD} , to transfer knowledge from the class token of a higher resolution to that of a lower resolution. This is achieved by serving cls_i as the teacher of cls_{i+1} with $H_i = W_i, H_i > H_{i+1}$. Combining with Eq. (1), the loss can be written as:

$$L(\Theta) = \mathbb{E}_{(X, T) \sim \mathcal{D}} \frac{1}{r} \left[\sum_{i=1}^r L_{CE}(Y^i; X^i, T, \Theta) + \sum_{i=1}^{r-1} L_{KD}(\text{cls}_{i+1}, \text{cls}_i) \right]. \quad (2)$$

Especially, teacher class tokens are detached from the gradient computational graph. At last, the loss is divided by r to ensure stability of training.

4.2. Global-Local Positional Embedding

The commonly-used positional embeddings highly depend on the size of input samples. As a result, when multiple resolutions are involved in the training process, positional embeddings need to be carefully adjusted when iterating images of different scales, as simple interpolations

incur performance drops. Therefore, we propose to use conditional positional embeddings both globally and locally to bridge the resolution gap among a broad range of resolutions. Below, we first introduce the global positional embedding and then describe its local counterpart.

Global position embedding. To incorporate location information in patch embeddings, a typical way is to add absolute position embedding (APE). Given the input sample X , x^{img} refers to the output image tokens of the patch embedding, whose spatial dimension equals $N_H \times N_W$ and feature dimension is D . For simplicity, we denote x as concatenation of class token `cls` and image tokens x^{img} , the absolute positional embedding p can be expressed as,

$$x = x + p, \quad p \in \mathbb{R}^{1 \times (N^H \times N^W + 1) \times D}. \quad (3)$$

The most straightforward way is implementing p with learned parameters, as widely-adopted in [16, 57, 63]. Another common tactic is fulfilled with sinusoidal mapping $\mathcal{F}_{\text{sine}}$ [22, 60], through which p is generated on-the-fly by a fixed function dependent on N_H, N_W and D . Due to space limitation, we provide explicit expression in Appendix B.1. Furthermore, compared with learned APE, the sine-cosine APE changes more weakly between different input scales, as displayed in Sec. 5.2. Therefore, we build our method upon sine-cosine APE with the assumption that smoother positional embedding would contribute to better resolution scalability.

We further improve the sine-cosine positional embedding with conditional computation such that the embeddings are tailored to the model during training. As illustrated in Fig. 3, a simple yet effective depth-wise convolution is applied so as to generate the final positional embedding conditioned on sinusoidal encoding. Since convolutions should be performed in 2-D dimension, we leave out the class token by concatenating a zero padding shaped of $\mathbb{R}^{1 \times 1 \times D}$ with output embeddings of DWconv. In general, the strategy introduced above aims at injecting smooth spatial information of global context into ViT, thereby we denote it as global positional embedding (GPE).

Local positional embedding. Positional embeddings introduced in [15, 34, 77] share the same design philosophy since they are both dynamically generated by input tokens and carry spatial information of local neighbourhood. It has been unveiled that such strategies can effectively introduce translation invariance into ViTs and hence facilitate generalizing to various resolutions. We refer to the positional embedding conditioned on local input feature as local positional embedding (LPE) and hypothesize that LPE is orthogonal with GPE in modelling spatial information of image tokens. Consequently, the combination of LPE and GPE may result in best resolution scalability.

To this end, we incorporate local positional embeddings into attention blocks in a similar fashion to [15]. Given a

multi-head self-attention block, a query Q , a key K and a value V are obtained through a linear projection, and the output z can be derived as:

$$z = \text{Softmax}(QK^T/\sqrt{D})V. \quad (4)$$

In particular, local spatial information of V is utilized. We first set the class token aside and reshape the value matrix to get $V' \in \mathbb{R}^{M \times D' \times N_H \times N_W}$, where M denotes the number of attention heads and D' satisfies $D = D' \cdot M$. Inspired by [40, 41], we generate dynamic positional embeddings conditioned on V' separately for each head. Therefore, a 3×3 depth-wise convolution is implemented to obtain the LPE for each head. The above operations can be denoted as mapping \mathcal{H} conditioned on V . Therefore, Eq. (4) can be re-written as:

$$z = \text{Softmax}(QK^T/\sqrt{D})V + \mathcal{H}(V). \quad (5)$$

By virtue of convolutions, LPE can be dynamically generated regardless of input scales. Eventually, in ResFormer, global and local positional embeddings are combined to ensure better generalization to novel resolutions.

5. Experiments

Implementation details. We instantiate ResFormer with DeiT [57] due to its simplicity. Given an input image, we resize it to 128, 160 and 224, respectively, for multi-resolution training throughout the experiments, unless specified otherwise. The resulting images are then used as inputs to ResFormer. For image classification, we use AdamW [31] as our optimizer and apply a cosine decay learning rate scheduler. Small and tiny models are trained with a batch size of 1024 and a learning rate of $5e^{-4}$, yet a learning rate of $8e^{-4}$ is used for the base model. We keep all augmentation and regularization settings in [57] for fair comparisons. For all experiments, we follow the official training and testing split as well as the evaluation metrics. More details can be found in Appendix B.1. For testing, we report results on a wide range of resolutions. Note that ResFormer only uses a single scale during testing.

5.1. Main Results

Effectiveness of ResFormer in image classification.

Tab. 1 presents the results of ResFormer and comparisons with DeiT [57] using various settings. In particular, we use ResFormer-M-R to denote a variant of ResFormer, where M represents the model size (*i.e.*, T, S, B for tiny, small and base models respectively) and R indicates the resolution used for training (*i.e.*, MR denotes multiple resolution; if R is a number, it represents the resolution itself).

When evaluated with a testing resolution of 224, ResFormers achieve highly competitive results—ResFormer-S-MR and ResFormer-B-MR offers an accuracy of 82.16%

Table 1. Top-1 Accuracy of DeiT and ResFormer on ImageNet-1K. Columns highlighted with grey background refer to the training resolutions of given models. Specifically, ResFormer adopts training resolutions of 128, 160 and 224 for multi-resolution training.

Model	Testing resolution										
	96	112	128	160	192	224	288	384	448	512	640
DeiT-T [57]	8.06	34.22	52.16	65.68	70.18	72.14	73.1	71.29	67.43	66.07	59.31
ResFormer-T-MR	61.40	64.93	67.78	71.09	72.97	73.85	74.85	75.04	74.39	73.77	71.65
DeiT-S [57]	17.55	54.34	67.02	75.62	78.60	79.83	80.02	78.10	75.85	72.63	63.86
ResFormer-S-128	70.25	73.91	75.47	77.06	77.48	76.89	74.78	69.55	64.54	58.34	45.25
ResFormer-S-160	67.34	72.26	75.05	78.06	78.94	79.19	78.25	74.86	71.38	66.65	54.77
ResFormer-S-224	57.80	66.36	71.35	76.99	79.63	80.83	81.42	80.65	79.28	77.73	73.26
ResFormer-S-MR	73.59	76.64	78.24	80.39	81.42	82.16	82.70	82.72	82.52	82.00	80.72
DeiT-B [57]	27.86	64.46	73.18	79.05	81.06	81.79	82.19	81.11	79.81	78.23	74.23
ResFormer-B-MR	75.86	78.42	79.74	81.52	82.28	82.72	83.02	83.29	82.9	82.63	81.72

Table 2. Results and comparisons of different backbones on ADE20K. All backbones are pre-trained on ImageNet-1k, among which MAE [22] uses unsupervised pre-training.

Backbone	#Param	Lr sched	mIoU	ms + flip
DeiT-S [57]	52.1M	80k	42.96	43.79
XCiT-S12/16 [1]	52.4M	160k	45.90	46.72
ResFormer-S-224	51.7M	80k	45.47	46.61
ResFormer-S-MR	51.7M	80k	46.31	47.45
DeiT-B [57]	120.6M	160k	45.36	47.16
XCiT-S24/16 [1]	109.0M	160k	47.69	48.57
ViT-B + MAE [22]	176.5M	160k	48.13	48.70
ResFormer-B-MR	119.8M	160k	48.30	49.28

and 82.72%, respectively, outperforming their DeiT counterparts by 2.33% and 0.93%. We also see from Tab. 1 that ResFormer trained with multi-resolution images outperforms models trained with single scale inputs with clear margins on all “seen” resolutions. For instance, ResFormer-S-MR outperforms ResFormer-S-128, ResFormer-S-160, ResFormer-S-224 by 2.77%, 2.33% and 1.33% respectively. Similar trends can also be found for ResFormer-B and ResFormer-T. This highlights the effectiveness of multi-resolution training.

Furthermore, for “unseen” resolutions, ResFormer demonstrates clear scaling capabilities. In particular, given a test resolution of 384, ResFormer-S-224 achieves a Top-1 accuracy of 80.65%, which is 2.55% higher than its DeiT-S counterpart (78.10%). This suggests that global-local positional embeddings can indeed improve generalization of different resolutions. ResFormer-S-MR further boosts the accuracy to 82.72%, demonstrating the benefit of multi-resolution training. Besides, ResFormer consistently generalize well to lower resolutions. Compared with DeiT, ResFormer-S-MR and ResFormer-B-MR increase by 56.24% and 48.00% when evaluated on a resolution of 80,

Table 3. Results and comparisons of different backbones on the mini-val set of COCO2017 using Mask R-CNN [23] and 3× training schedule. All backbones are pre-trained on ImageNet-1k in the supervised setting. Part of results are credited to [1, 8].

Backbone	#Param	AP ^b	AP ^b ₅₀	AP ^b ₇₅	AP ^m	AP ^m ₅₀	AP ^m ₇₅
PVT-Small [63]	44.1M	43.0	65.3	46.9	39.9	62.5	42.8
XCiT-S12/16 [1]	44.3M	45.3	67.0	49.5	40.8	64.0	43.8
ViT-S [36]	43.8M	44.0	66.9	47.8	39.9	63.4	42.2
ViTDet-S [35]	45.7M	44.5	66.9	48.4	40.1	63.6	42.5
ResFormer-S-MR	45.6M	46.4	68.5	50.4	40.7	64.7	43.4
PVT-Large [63]	81.0M	44.5	66.0	48.3	40.7	63.4	43.7
XCiT-M24/16 [1]	101.1M	46.7	68.2	51.1	42.0	65.6	44.9
ViT-B [36]	113.6M	45.8	68.2	50.1	41.3	65.1	44.4
ViTDet-B [35]	121.3M	46.3	68.6	50.5	41.6	65.3	44.5
ResFormer-B-MR	115.3M	47.6	69.0	52.0	41.9	65.9	44.4

highlighting the effectiveness of ResFormer when dealing with significant resolution shifts during inference.

Semantic segmentation. To show flexibility of ResFormer, We evaluate for semantic segmentation on ADE20K [80] with UperNet [69]. As shown in Tab. 2, ResFormer-S-224 improves DeiT-S by 2.51 measured by mIoU. Both ResFormer-S-MR and ResFormer-B-MR, which are pre-trained with the multi-resolution strategy, achieve better results. In particular, ResFormer-S-MR reaches up to 47.45 and ResFormer-B-MR hits the peak of 49.28 mIoU. This suggests that ResFormer effectively models multi-scale and high-resolution features for pixel-level dense predictions. Note that ResFormer-B-MR and ResFormer-S-MR are directly used as pre-trained backbones and we do not perform multi-resolution fine-tuning on ADE20K, since segmentation tasks already require images with a size of 512×512 as inputs, and multi-resolution training would be computationally expensive. Nonetheless, results in Tab. 2 demon-

Table 4. Top-1 Accuracy of TimeSformer on Kinetics-400. MR stands for multi-resolution training.

Model	Testing resolution				
	96	128	160	224	288
TimeSformer [3]	26.28	61.94	70.60	75.54	75.45
ResFormer-B-224	58.61	68.50	73.09	76.32	76.78
ResFormer-B-160	67.28	71.64	74.56	75.98	75.18
ResFormer-B-128	64.66	72.32	74.13	74.19	72.51
ResFormer-B-MR	70.56	74.33	76.38	77.32	77.56

strate the great potential of transferring models that are pre-trained with multiple resolutions for dense prediction tasks.

Object detection. We further explore performance of ResFormer on COCO2017 [39] for object detection and instance segmentation, following the designs of ViTDet [35] by appending simple feature pyramids on the feature maps of last-layer outputs and using both non-shifted window attention and global self-attention blocks. In addition, To adapt from ResFormer pre-trained on ImageNet-1K, we also adopt global positional embedding and inject local positional embeddings into all attention blocks. According to results reported in Tab. 3, ResFormer achieves promising results, *e.g.* ResFormer-S-MR outperforms ViTDet-S by 2.0 box AP and 0.6 mask AP and ResFormer-B-MR surpasses ViTDet-B by 1.3 box AP and 0.3 mask AP. We believe that the improved resolution scalability of ResFormer contributes to better performance on object detection.

Video action recognition. We also evaluate ResFormer for video action recognition on Kinetics400. For an easy adaptation from our pre-trained image models to the video domain, we choose the TimeSformer [3] framework with a divided spatial and temporal attention design. In particular, we initialize the backbone with weights of model pre-trained on ImageNet-1K and conduct multi-resolution training on Kinetics400 with clip sizes set to $8 \times 224 \times 224$, $8 \times 160 \times 160$ and $8 \times 128 \times 128$ respectively. As Tab. 4 demonstrates, ResFormer-B fine-tuned with single resolution outweighs vanilla TimeSformer by 0.78% on a testing resolution of 224 and generalizes better to clips of both higher and lower resolutions. On top of that, by implementing multi-resolution training on video samples, ResFormer improves performance on each training resolution by a large margin, *e.g.*, the Top-1 accuracy on testing resolution of 160 grows from 73.09% to 76.08%.

5.2. Discussion

Training resolutions. We experiment with 3 different settings using a small model, *i.e.*, (128, 160, 224), (160, 224, 288) and (128, 224, 384), which we denoted as (a), (b) and (c) respectively. The results are summarized in Fig. 4. We see that ResFormer achieves outstanding performance on a



Figure 4. Top-1 Accuracy of ResFormer-S-MR with different training resolutions on ImageNet-1K.

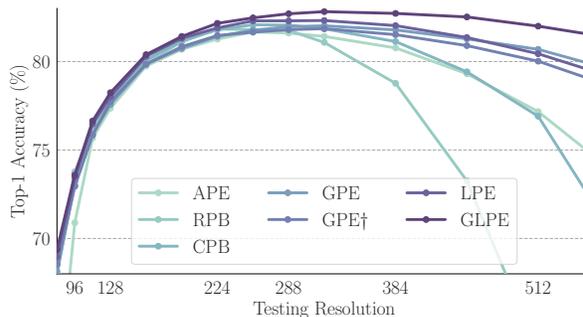


Figure 5. Results of different positional embedding strategies on a broader range of resolutions.

wide range of resolutions. More specifically, compared with (a), (b) adopts higher training resolutions, consequently reflecting on performance rise in high-resolution inputs and incurring a drop on low resolution. Regarding differences between setting of (b) and (c), the spectrum of training resolutions expands in both directions. Despite the wide range between 128 and 384, we witness an all-around improvement of (c) over (b), highlighting that ResFormer is able to deal with significant resolution variations.

Positional embedding. We evaluate the performance of ResFormer with different positional embedding strategies using a small model. In particular, we compare with (1) APE, which stands for vanilla absolute positional embedding in DeiT [57]. In practice, we set the spatial dimension of APE according to the highest training resolution and downsample it for lower resolutions. For inference, the position embedding can be re-scaled to any test resolution with bicubic interpolation; (2) APE* which uses an individual APE for each training resolution; (3) RPB, which is introduced in [41] and we use the RPB of highest resolution for interpolation during inference; (4) CPB, which is a resolution-agnostic strategy [40] and images of arbitrary scales can be input into ViTs with CPB directly; (5) GPE, which is our global positional embedding; (6) GPE†, which represents the plain sine-cosine absolute positional embedding without convolutional enhancement; (7) LPE, which

Table 5. Results of ResFormer-S-MR with different positional embedding (PE) strategies on ImageNet-1K. The performance gain compared to single-resolution training is indicated in the bracket.

PE	Testing resolution		
	128	160	224
APE	77.36 ($\uparrow 3.99$)	79.74 ($\uparrow 2.46$)	81.27 ($\uparrow 1.44$)
APE*	77.31 ($\uparrow 3.93$)	79.58 ($\uparrow 2.21$)	81.42 ($\uparrow 1.59$)
RPB	77.90 ($\uparrow 2.74$)	79.92 ($\uparrow 2.04$)	81.84 ($\uparrow 1.27$)
CPB	77.64 ($\uparrow 1.33$)	79.77 ($\uparrow 1.77$)	81.74 ($\uparrow 1.13$)
GPE \dagger	77.73 ($\uparrow 2.73$)	79.83 ($\uparrow 2.27$)	81.47 ($\uparrow 1.29$)
GPE	77.57 ($\uparrow 2.76$)	79.63 ($\uparrow 2.05$)	81.42 ($\uparrow 1.40$)
LPE	78.02 ($\uparrow 2.62$)	80.29 ($\uparrow 2.29$)	81.90 ($\uparrow 1.28$)
GLPE	78.24 ($\uparrow 2.77$)	80.39 ($\uparrow 2.33$)	82.16 ($\uparrow 1.33$)

is our local positional embedding; (8) GLPE, which is the combination our GPE and LPE.

Tab. 5 shows the results of different positional embeddings on training resolutions. We see that interpolating APE makes no differences compared with maintaining multiple APEs (*i.e.*, APE*), which suggests that ViTs can be trained to deal with different scales of inputs with shared spatial information. Furthermore, all positional embeddings coupled with multi-resolution training demonstrate better results compared to their counterparts trained with single resolutions, *i.e.*, steady gains are made by all positional embeddings when testing on 128, 160 and 224, (gains are shown in the bracket in Tab. 5). Fig. 5 further presents the results of generalizing to more resolutions. Clear performance drops can also be observed in Fig. 5 when APE, RPB and CPB are scaled up to unseen large resolutions, especially RPB. In contrast, LPE and GPE decreases slowly towards extremely large resolutions. GLPE, the combination of LPE and GPE, offers the best results.

Knowledge distillation. To strengthen the interaction between different resolutions, we use a smooth-L1 loss to distill information from class tokens. We also experiment with a L2 loss (*i.e.*, Mean squared error). Further, as inputs of different resolutions output features with different scales, we additionally follow the practice in DeiT [57] by distilling logits with a Kullback-Leibler divergence loss. The experiments are conducted on ResFormer-S-MR for 100 epochs for efficiency purposes. Tab. 6 shows the ablation results. We observe the efficacy of distilling with class tokens compared to logits. In addition, the smooth L_1 loss have similar performance with L_2 loss with slightly better results on high resolutions (*i.e.*, 224).

Training strategies. We also explore a widely-used multi-resolution training strategy [23,25] without cross-scale consistency loss, where one iteration consists of randomly sampled images of one certain resolution. In particular, we feed samples of different scales (*i.e.*, 128, 160, 224) iter-

Table 6. Results of ResFormer-S-MR with different distillation strategies on ImageNet-1K for 100ep. Performance gains over result of training **without distillation** are shown in the bracket.

Distillation		Testing resolution		
Target	Loss	128	160	224
logit	KL	73.50 ($\leftrightarrow 0.0$)	76.45 ($\uparrow 0.07$)	78.82 ($\uparrow 0.26$)
cls	L_2	74.71 ($\uparrow 1.21$)	77.27 ($\uparrow 0.89$)	79.33 ($\uparrow 0.77$)
cls	smooth L_1	74.71 ($\uparrow 1.21$)	77.39 ($\uparrow 1.01$)	79.68 ($\uparrow 1.12$)

Table 7. Results of ResFormer-S-MR with different training strategies on ImageNet-1K. We append performance gain/drop compared with single-resolution training in the bracket.

Training Strategy	Testing resolution		
	128	160	224
MR (iter)	75.70 ($\uparrow 0.23$)	78.31 ($\uparrow 0.25$)	80.32 ($\downarrow 0.51$)
MR (epoch)	75.18 ($\downarrow 0.29$)	78.05 ($\downarrow 0.01$)	80.26 ($\uparrow 0.16$)
MR w/o KD	77.72 ($\uparrow 2.25$)	79.66 ($\uparrow 1.60$)	81.77 ($\uparrow 0.94$)
MR	78.24 ($\uparrow 2.77$)	80.39 ($\uparrow 2.33$)	82.16 ($\uparrow 1.33$)

atively based on two settings: (1) iteration-based, where each mini-batch uses one resolution and resolutions vary for different training iterations; (2) epoch-based, where a fixed resolution is used for each epoch and the change of resolutions only occur at the epoch-level. As Tab. 7 shows, both iteration-based and epoch-based multi-resolution training generate worse results compared to single resolution training. In contrast, our strategy demonstrates strong advantages on all training resolutions by clear margins, even without the scale-consistency loss, highlighting the importance of enforcing consistencies of all resolutions in a mini-batch.

Qualitative visualizations. We visualize two positional embeddings on resolutions of 128, 160, 224 and 384, respectively. As shown in Fig. 6, compared with APEs shifted by interpolation, our GPEs that are generated with convolutions demonstrate a smoother variations among input scales. In addition, Fig. 5 suggests that GPE generalizes better to higher resolutions unseen in training.

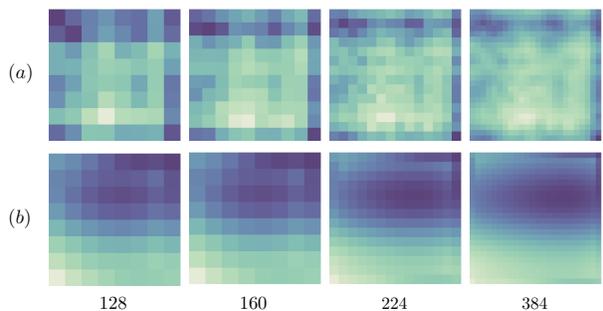


Figure 6. Heatmaps of different PE averaged on each token. (a): Absolute Positional Embeddings (APE), (b): Global Positional Embeddings (GPE).

6. Conclusion

We introduced ResFormer, a ViT framework to encourage excellent all-round performance on a wide range of resolutions. In particular, ResFormer was motivated by training on sample of different scales and aided by a scale-consistency loss. A global-local positional embedding strategy was also introduced to facilitate better generalization on unseen resolutions. Extensive experiments demonstrated promising scalabilities of ResFormer in a broad range of resolutions. We also observe that ResFormer can be readily adapted to downstream tasks, *e.g.*, semantic segmentation, object detection and video action recognition.

Acknowledgement This project was supported by NSFC under Grant No. 62102092 and No. 62032006.

References

- [1] Alaaeldin Ali, Hugo Touvron, Mathilde Caron, Piotr Bojanowski, Matthijs Douze, Armand Joulin, Ivan Laptev, Natalia Neverova, Gabriel Synnaeve, Jakob Verbeek, et al. Xcit: Cross-covariance image transformers. In *NeurIPS*, 2021. 1, 6
- [2] Irwan Bello, William Fedus, Xianzhi Du, Ekin Dogus Cubuk, Aravind Srinivas, Tsung-Yi Lin, Jonathon Shlens, and Barret Zoph. Revisiting resnets: Improved training and scaling strategies. In *NeurIPS*, 2021. 2
- [3] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *ICML*, 2021. 1, 2, 7, 13, 14
- [4] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020. 1, 3
- [5] Chun-Fu Chen, Rameswar Panda, and Quanfu Fan. Regionvit: Regional-to-local attention for vision transformers. In *ICLR*, 2022. 2
- [6] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, et al. Mmdetection: Open mmlab detection toolbox and benchmark. *arXiv preprint arXiv:1906.07155*, 2019. 13
- [7] Mengzhao Chen, Mingbao Lin, Ke Li, Yunhang Shen, Yongjian Wu, Fei Chao, and Rongrong Ji. Cf-vit: A general coarse-to-fine method for vision transformer. *arXiv preprint arXiv:2203.03821*, 2022. 1
- [8] Zhe Chen, Yuchen Duan, Wenhai Wang, Junjun He, Tong Lu, Jifeng Dai, and Yu Qiao. Vision transformer adapter for dense predictions. In *ICLR*, 2023. 6
- [9] Xiangxiang Chu, Zhi Tian, Yuqing Wang, Bo Zhang, Haibing Ren, Xiaolin Wei, Huaxia Xia, and Chunhua Shen. Twins: Revisiting the design of spatial attention in vision transformers. In *NeurIPS*, 2021. 2
- [10] Xiangxiang Chu, Zhi Tian, Bo Zhang, Xinlong Wang, Xiaolin Wei, Huaxia Xia, and Chunhua Shen. Conditional positional encodings for vision transformers. *arXiv preprint arXiv:2102.10882*, 2021. 2
- [11] Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: Learning augmentation policies from data. In *CVPR*, 2019. 4
- [12] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *CVPR Workshops*, 2020. 4
- [13] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G Carbonell, Quoc Viet Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. In *ACL*, 2019. 2
- [14] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 2
- [15] Xiaoyi Dong, Jianmin Bao, Dongdong Chen, Weiming Zhang, Nenghai Yu, Lu Yuan, Dong Chen, and Baining Guo. Cswin transformer: A general vision transformer backbone with cross-shaped windows. In *CVPR*, 2022. 5
- [16] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2020. 1, 2, 4, 5
- [17] Stéphane d’Ascoli, Hugo Touvron, Matthew L Leavitt, Ari S Morcos, Giulio Biroli, and Levent Sagun. Convit: Improving vision transformers with soft convolutional inductive biases. In *ICML*, 2021. 2
- [18] Alaaeldin El-Nouby, Gautier Izacard, Hugo Touvron, Ivan Laptev, Hervé Jegou, and Edouard Grave. Are large-scale datasets necessary for self-supervised pre-training? *arXiv preprint arXiv:2112.10740*, 2021. 2
- [19] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. Multiscale vision transformers. In *ICCV*, 2021. 3
- [20] Priya Goyal, Mathilde Caron, Benjamin Lefaudeux, Min Xu, Pengchao Wang, Vivek Pai, Mannat Singh, Vitaliy Liptchinsky, Ishan Misra, Armand Joulin, et al. Self-supervised pretraining of visual features in the wild. *arXiv preprint arXiv:2103.01988*, 2021. 2
- [21] Jiaqi Gu, Hyoukjun Kwon, Dilin Wang, Wei Ye, Meng Li, Yu-Hsin Chen, Liangzhen Lai, Vikas Chandra, and David Z Pan. Multi-scale high-resolution vision transformer for semantic segmentation. In *CVPR*, 2022. 3
- [22] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *CVPR*, 2022. 5, 6
- [23] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *ICCV*, 2017. 2, 6, 8, 13
- [24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. *TPAMI*, 2015. 2
- [25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 2, 8
- [26] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, Dawn Song, Jacob Steinhardt,

- and Justin Gilmer. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *ICCV*, 2021. 12
- [27] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *ICLR*, 2019. 12
- [28] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. In *CVPR*, 2021. 12
- [29] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *ICML*, 2021. 1
- [30] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017. 13
- [31] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 5
- [32] Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and Neil Houlsby. Big transfer (bit): General visual representation learning. In *ECCV*, 2020. 2
- [33] Hei Law and Jia Deng. Cornernet: Detecting objects as paired keypoints. In *ECCV*, 2018. 3
- [34] Kunchang Li, Yali Wang, Peng Gao, Guanglu Song, Yu Liu, Hongsheng Li, and Yu Qiao. Uniformer: Unified transformer for efficient spatiotemporal representation learning. In *ICLR*, 2022. 1, 5
- [35] Yanghao Li, Hanzi Mao, Ross Girshick, and Kaiming He. Exploring plain vision transformer backbones for object detection. *arXiv preprint arXiv:2203.16527*, 2022. 1, 6, 7
- [36] Yanghao Li, Saining Xie, Xinlei Chen, Piotr Dollár, Kaiming He, and Ross Girshick. Benchmarking detection transfer learning with vision transformers. *arXiv preprint arXiv:2111.11429*, 2021. 6
- [37] Mingbao Lin, Mengzhao Chen, Yuxin Zhang, Ke Li, Yunhang Shen, Chunhua Shen, and Rongrong Ji. Super vision transformer. *arXiv preprint arXiv:2205.11397*, 2022. 1
- [38] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *CVPR*, 2017. 2
- [39] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 2, 7, 13
- [40] Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang, Li Dong, et al. Swin transformer v2: Scaling up capacity and resolution. In *CVPR*, 2022. 1, 2, 5, 7
- [41] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, 2021. 1, 2, 3, 5, 7
- [42] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin transformer. In *CVPR*, 2022. 1
- [43] Chuofan Ma, Qiushan Guo, Yi Jiang, Zehuan Yuan, Ping Luo, and Xiaojuan Qi. Rethinking resolution in the context of efficient video recognition. In *NeurIPS*, 2022. 3
- [44] Lingchen Meng, Hengduo Li, Bor-Chun Chen, Shiyi Lan, Zuxuan Wu, Yu-Gang Jiang, and Ser-Nam Lim. Adavit: Adaptive vision transformers for efficient image recognition. In *CVPR*, 2022. 1
- [45] MMSegmentation Contributors. OpenMMLab Semantic Segmentation Toolbox and Benchmark, 2020. 13
- [46] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 1
- [47] Ilija Radosavovic, Raj Prateek Kosaraju, Ross Girshick, Kaiming He, and Piotr Dollár. Designing network design spaces. In *CVPR*, 2020. 2
- [48] Kanchana Ranasinghe, Muzammal Naseer, Salman Khan, Fahad Shahbaz Khan, and Michael S Ryoo. Self-supervised video transformer. In *CVPR*, 2022. 3
- [49] Yongming Rao, Wenliang Zhao, Benlin Liu, Jiwen Lu, Jie Zhou, and Cho-Jui Hsieh. Dynamicvit: Efficient vision transformers with dynamic token sparsification. In *NeurIPS*, 2021. 1
- [50] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishal Shankar. Do imagenet classifiers generalize to imagenet? In *ICML*, 2019. 12
- [51] Carlos Riquelme, Joan Puigcerver, Basil Mustafa, Maxim Neumann, Rodolphe Jenatton, André Susano Pinto, Daniel Keysers, and Neil Houlsby. Scaling vision with sparse mixture of experts. In *NeurIPS*, 2021. 2
- [52] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations. In *NAACL*, 2018. 2
- [53] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 3
- [54] Bharat Singh and Larry S Davis. An analysis of scale invariance in object detection snip. In *CVPR*, 2018. 3
- [55] Andreas Peter Steiner, Alexander Kolesnikov, Xiaohua Zhai, Ross Wightman, Jakob Uszkoreit, and Lucas Beyer. How to train your vit? data, augmentation, and regularization in vision transformers. *TMLR*, 2022. 1
- [56] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *ICML*, 2019. 2
- [57] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *ICML*, 2021. 1, 2, 3, 5, 6, 7, 8, 11, 12
- [58] Hugo Touvron, Matthieu Cord, and Hervé Jégou. Deit iii: Revenge of the vit. *arXiv preprint arXiv:2204.07118*, 2022. 1
- [59] Hugo Touvron, Andrea Vedaldi, Matthijs Douze, and Hervé Jégou. Fixing the train-test resolution discrepancy. In *NeurIPS*, 2019. 3
- [60] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia

Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 1, 2, 5

[61] Haoan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global representations by penalizing local predictive power. In *NeurIPS*, 2019. 12

[62] Junke Wang, Dongdong Chen, Zuxuan Wu, Chong Luo, Luwei Zhou, Yucheng Zhao, Yujia Xie, Ce Liu, Yu-Gang Jiang, and Lu Yuan. Omnivl: One foundation model for image-language and video-language tasks. In *NeurIPS*, 2022. 1

[63] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *ICCV*, 2021. 1, 5, 6

[64] Yulin Wang, Rui Huang, Shiji Song, Zeyi Huang, and Gao Huang. Not all images are worth 16x16 words: Dynamic transformers for efficient image recognition. In *NeurIPS*, 2021. 1

[65] Yixuan Wei, Han Hu, Zhenda Xie, Zheng Zhang, Yue Cao, Jianmin Bao, Dong Chen, and Baining Guo. Contrastive learning rivals masked image modeling in fine-tuning via feature distillation. *arXiv preprint arXiv:2205.14141*, 2022. 4

[66] Chao-Yuan Wu, Ross Girshick, Kaiming He, Christoph Feichtenhofer, and Philipp Krahenbuhl. A multigrid method for efficiently training video models. In *NeurIPS*, 2020. 2, 3

[67] Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang. Cvt: Introducing convolutions to vision transformers. In *ICCV*, 2021. 1

[68] Kan Wu, Houwen Peng, Minghao Chen, Jianlong Fu, and Hongyang Chao. Rethinking and improving relative position encoding for vision transformer. In *ICCV*, 2021. 2

[69] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *ECCV*, 2018. 6, 13

[70] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han Hu. Simmim: A simple framework for masked image modeling. In *CVPR*, 2022. 1

[71] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Yixuan Wei, Qi Dai, and Han Hu. On data scaling in masked image modeling. In *CVPR*, 2023. 2

[72] Zhen Xing, Qi Dai, Han Hu, Jingjing Chen, Zuxuan Wu, and Yu-Gang Jiang. Svformer: Semi-supervised video transformer for action recognition. In *CVPR*, 2023. 1

[73] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zi-Hang Jiang, Francis EH Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. In *ICCV*, 2021. 1

[74] Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *ICCV*, 2019. 4, 13

[75] Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers. In *CVPR*, 2022. 2

[76] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *ICLR*, 2018. 4, 13

[77] Linfeng Zhang, Chenglong Bao, and Kaisheng Ma. Self-distillation: Towards efficient and compact neural networks. In *NeurIPS*, 2021. 5

[78] Xiaosong Zhang, Yunjie Tian, Lingxi Xie, Wei Huang, Qi Dai, Qixiang Ye, and Qi Tian. Hivit: A simpler and more efficient design of hierarchical vision transformer. In *ICLR*, 2023. 1

[79] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data augmentation. In *AAAI*, 2020. 4

[80] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ade20k dataset. *IJCV*, 2019. 2, 6, 13

A. More Experiments

A.1. More about Resolution Scalability

Scalability of vanilla ViTs. As displayed in Fig. 7 and Fig. 8, in order to provide more comprehensive insights into resolution scalability, we further test tiny and base models of DeiT [57] which are pre-trained on training resolutions of 196, 128, 224, 288 and 384, respectively. The evaluation is conducted by generalizing models to different testing resolutions ranging from 80 to 576. We can observe that the

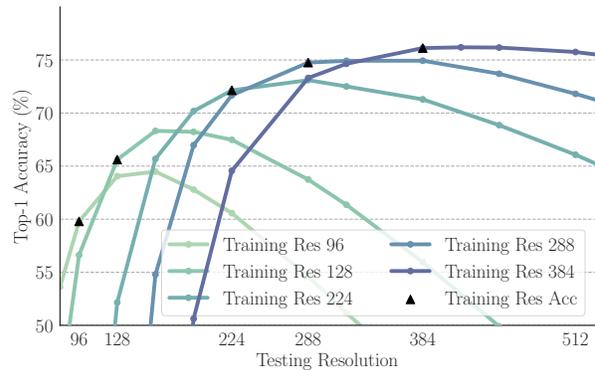


Figure 7. Top-1 accuracy of DeiT-T trained with 5 different resolutions and tested on resolutions varying from 80 to 576.

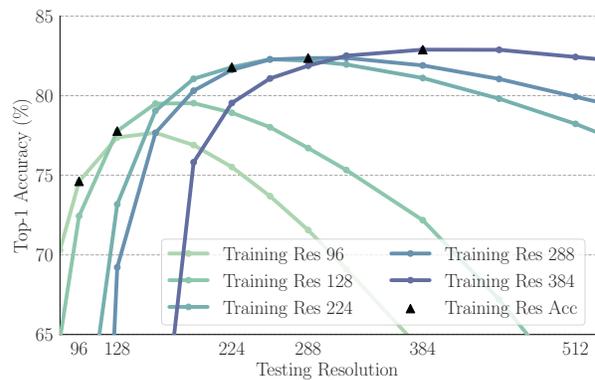


Figure 8. Top-1 accuracy of DeiT-B trained with 5 different resolutions and tested on resolutions varying from 80 to 576.

Table 8. Comparison of Top-1 accuracy between DeiT and ResFormer on ImageNet-1K with high testing resolutions.

Model	Testing resolution			
	512	640	800	1024
DeiT-S-224	72.63	63.86	49.31	31.45
ResFormer-S-MR (224)	82.00	80.72	78.12	72.49
DeiT-S-384	81.09	79.35	75.67	67.61
ResFormer-S-MR (384)	83.86	83.71	83.37	82.58

trends towards scaling up and scaling down testing resolutions are consistent with ones on DeiT-S.

Extending the range of testing resolutions. To further explore the potential for ResFormer, we extend the range of testing resolutions to 1024. As shown in Tab. 8, compared with DeiT, ResFormer achieves much more decent performance on fairly testing resolutions.

A.2. Evaluation on Robustness Datasets

We also evaluate our models on ImageNet-related robustness datasets, *i.e.*, ImageNet-Rendition (IN-R) [26], ImageNet-A (IN-A) [28], ImageNet-Sketch (IN-SK) [61], ImageNet-C (IN-C) [27] and ImageNetv2 (IN-v2) [50]. As reported in Tab. 9, we observe that ResFormer achieves promising performance on robustness as well. For example, ResFormer-S-224 is superior to DeiT-S on each dataset while ResFormer-S-MR makes further improvements. In particular, on IN-A, ResFormer-S-MR surpasses ResFormer-S-224 by 7.88 % and DeiT-S by 10.07%. This suggests that training with multi-scale inputs facilitates ViTs to cope with hard as well as out-of-distribution inputs.

Table 9. Performance on ImageNet-based robustness benchmarks. mCE [27] is employed for IN-C while Top-1 accuracy is used for IN-R, IN-A and IN-SK.

Model	IN-R \uparrow	IN-A \uparrow	IN-SK \uparrow	IN-C \downarrow	INv2 \uparrow
DeiT-S [57]	41.93	19.84	29.09	54.60	68.47
ResFormer-S-224	43.95	22.03	30.91	52.31	69.81
ResFormer-S-MR	45.08	29.91	31.47	51.03	71.68
DeiT-B [57]	44.66	28.15	31.96	48.52	70.91
ResFormer-B-MR	45.38	33.89	33.06	48.83	71.88

A.3. Training Efficiency

During the training of ResFormer, each input sample is replicated by r times, and thus this increases the training time. For efficiency, we reduce the total number of training epochs to 200, 150 and 100 respectively while keeping other hyperparameters unchanged. As shown in Fig. 9, ResFormer demonstrates competitive performance on training efficiency. For instance, ResFormer-S-MR with 200-epoch

training surpasses the 300-epoch counterparts ResFormer-S-224 and DeiT-S by 0.83% and 1.83% in Top-1 accuracy despite that they share similar training time.

As depicted in Fig. 10, training with a single lower resolution (*i.e.*, 160) significantly saves time. Nevertheless, ResFormer-S-MR still has an edge on time-performance trade-off, *e.g.*, ResFormer-S-160 with 450-epoch training is more time-consuming than ResFormer-S-MR with 200-epoch training while the accuracy is 0.61% lower.

B. Implementation Details

B.1. Image Classification

Sine-Cosine positional embedding. We demonstrate the explicit mapping function \mathcal{F}_{sine} for sine-cosine positional embedding p as follows. Firstly, image tokens are placed in a 2D spatial dimension as $x^{img} \in \mathbb{R}^{N_H \times N_W \times D}$. We denote the positional embedding for the token coordinated at (m, n) as $p_{m,n} \in \mathbb{R}^{1 \times D}$. Particularly, d -th dimension of $p_{m,n}$ can be mapped with $\mathcal{F}_{sine}(m, n, d)$ as below,

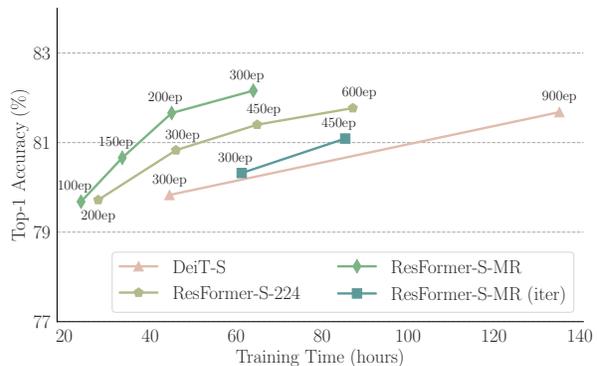


Figure 9. Trade-off between training time and Top-1 Accuracy on ImageNet-1K with a testing resolution of 224. Same hardware and software settings are adopted for all experiments, *i.e.*, we utilize $8 \times$ V100-32GB GPUs and set the per-GPU batch size to 128.

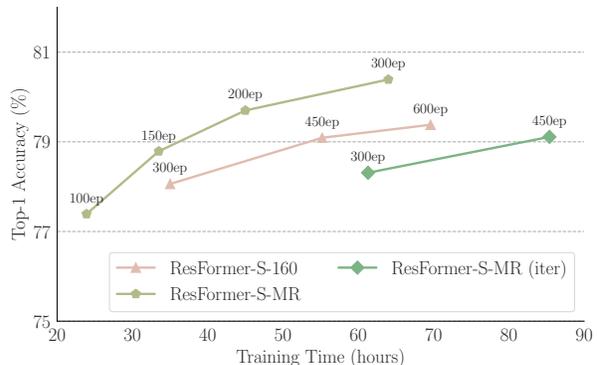


Figure 10. Trade-off between training time and Top-1 Accuracy on ImageNet-1K with a testing resolution of 160. Same hardware and software settings are adopted for all experiments.

$$\mathcal{F}_{sin}(m, n, d) = \begin{cases} f_{sin}(m, d, N_H, D) & \text{if } d < D/2 \\ f_{sin}(n, d, N_W, D) & \text{otherwise} \end{cases},$$

$$f_{sin}(pos, d, N, D) = \begin{cases} \sin(\frac{pos}{N+\epsilon}/T^{2d/D}) & \text{if } d\%2 = 0 \\ \cos(\frac{pos}{N+\epsilon}/T^{2(d-1)/D}) & \text{otherwise} \end{cases},$$

where the temperature T and ϵ is set to 10000 and $1e^{-6}$ respectively, and a normalization is also used to ensure better continuity among varying resolutions. For simplicity, N_i^H, N_i^W, D are omitted from function parameters.

Detailed hyperparameters. For experiments of image classification on ImageNet-1K, we set the hyperparameters for training ResFormer-T, ResFormer-S, ResFormer-B from scratch and fine-tuning on DeiT according to Tab. 11.

Augmentation strategy. Motivated by unsupervised learning, we apply separate random augmentation on different scales of inputs. In particular, to ensure the consistency of class tokens between different scales, as an exception, we apply MixUp [76] and CutMix [74] across different scales with same variables. As shown in Tab. 10, separate augmentation slightly outperform its counterpart, especially on the lowest testing resolution.

Table 10. Ablation study of augmentation strategies.

Model	Sep Aug	Testing resolution		
		128	160	224
ResFormer-S-MR		77.50	80.14	81.93
ResFormer-S-MR	✓	78.24	80.39	82.16

B.2. Semantic Segmentation

We follow the common practice on ADE20K [80] by training on 512×512 inputs for 80k iterations for ResFormer-S and for 160k iterations for ResFormer-B, respectively. In addition, we employ the AdamW optimizer with a learning rate of $6e^{-5}$, a weight decay of 0.01 and a batch-size of 16. We base our implementation on MMSEgmentation [45] and adopt the corresponding augmentations, *i.e.*, random resizing with the ratio range set to (0.5, 2.0), random horizontal flipping with probability of 0.5 and random photometric distortion. Despite that ResFormer employs a columnar structure, we simply extract features from different layers (*i.e.* the 2nd, 5th, 8th and 11th layers) as inputs of UperNet [69] without FPN-like necks. We report results in two different testing settings. For the first one, inputs are scaled to having a shorter side of 512. In addition, we apply flipping on inputs of multiple scales that are varied in $(0.5, 0.75, 1.0, 1.25, 1.5, 1.75) \times$ of training resolutions.

Table 11. Hyperparameters for training on ImageNet-1K.

Hyperparameters	Tiny / Small	Base	Fine-tune
Epochs	300	200	30
Base learning rate	5e-4	8e-4	5e-5
Warmup epochs	5	20	5
Stoch. depth	0.1	0.2	0.1
Gradient clipping	✗	5.0	✗
Batch size		1024	
Weight decay		0.05	
Optimizer		AdamW	
Learning rate schedule		Cosine	
Repeated augmentation		✓	
Random erasing		0.25	
Random augmentation		9/0.5	
Mixup		0.8	
Cutmix		1.0	
Color jitter		0.4	

B.3. Object Detection

To further validate the efficacy of ResFormer on dense prediction tasks, we evaluate ResFormer on COCO 2017 [39] for object detection and instance segmentation. In particular, we adopt Mask R-CNN [23] as our framework based on MMDetection [6] and train with the $3 \times$ schedule. Furthermore, we utilize AdamW optimizer with a learning rate of $1e^{-4}$, weight decay of 0.05 and a batch size of 16. It is worth noting that we follow the common multi-scale training for object detection instead of fine-tuning with the multi-resolution strategy. Therefore, training samples are resized randomly so that the shorter sizes vary from 480 to 800 with step of 32 and the longer sides are within 1333.

B.4. Video Action Recognition

Similar to the implementation for images, we train ResFormer on videos by replicating video clips to get multi-scale copies. Specifically, given a certain sampling rate s of $1/32$, a clip X of $F = 8$ frames is sampled and replicated into r copies. Different cropping sizes are applied on each sequence of frames. Consequently the i -th training copy X_i is sized in $\mathbb{R}^{F \times H_i \times W_i}, i \in \{1, \dots, r\}$. We also keep the augmentation strategy used for images by applying separate random augmentations [3] on each clip.

We follow the divided attention design adopted in TimeSFormer [3], in which attention computation is conducted along spatial dimension and temporal dimension separately. In order to align with image models, we only incorporate global and local positional embeddings into into spatial dimensions. For training on Kinetics-400 [30], we

adopt the same strategy with TimeSFormer [3]. In particular, the training epoch is set to 15 and the initial learning rate is set to $5e^{-3}$. In addition, we employ a SGD optimizer and a multi-step scheduler which divides the learning rate by 10 times at the 11th and the 14th epoch respectively.

In particular, we observe that ResFormer achieves better performance on videos with L_2 scale consistency loss. In order to improve performance by ensuring coherence in pre-training and fine-tuning. We adapt ResFormer-B-MR for L_2 loss for an extended fine-tuning of 100 epochs which matches the common 300-epoch pre-training. For fair comparison, we initiate all ResFormers in Kinetics-400 downstream tasks with same pre-trained weights.