

# Train ViT on Small Dataset With Translation Perceptibility



Huan Chen, Wentao Wei, Ping Yao

Institute of Computing Technology, Chinese Academy of Sciences;

University of Chinese Academy of Sciences; Chien-Shiung Wu College, Southeast University



## INTRODUCTION

- ◆ Vision Transformers (ViTs) exhibit deficiencies when trained on smaller datasets, specifically lacking locality, inductive biases, and hierarchical structure, which are inherent in convolutional approaches.
- ◆ Inspired by the **translation equivariance** of CNNs, we propose a novel self-supervised auxiliary task that enables ViTs to acquire **translation perceptibility**.
- ◆ Our method delivers competitive performance on small datasets across various resolutions without necessitating architectural modifications, and it can be seamlessly integrated with previous methods for enhanced utility.

## Translation Perceptibility

Consider "x" and "y" as the input and output respectively, let "TS" denote the translation-set, and let "F" and "trans" stand for the model and translation function respectively.

- ◆ **Translation Invariance** means that the system produces exactly the same response (output) regardless of how its input is translated.

$$y = F(\text{trans}(x, TS)) = F(x)$$

- ◆ **Translation Invariance** means that the system produces exactly the same response (output) regardless of how its input is translated.

$$y = F(\text{trans}(x, TS)) = \text{trans}(F(x), TS)$$

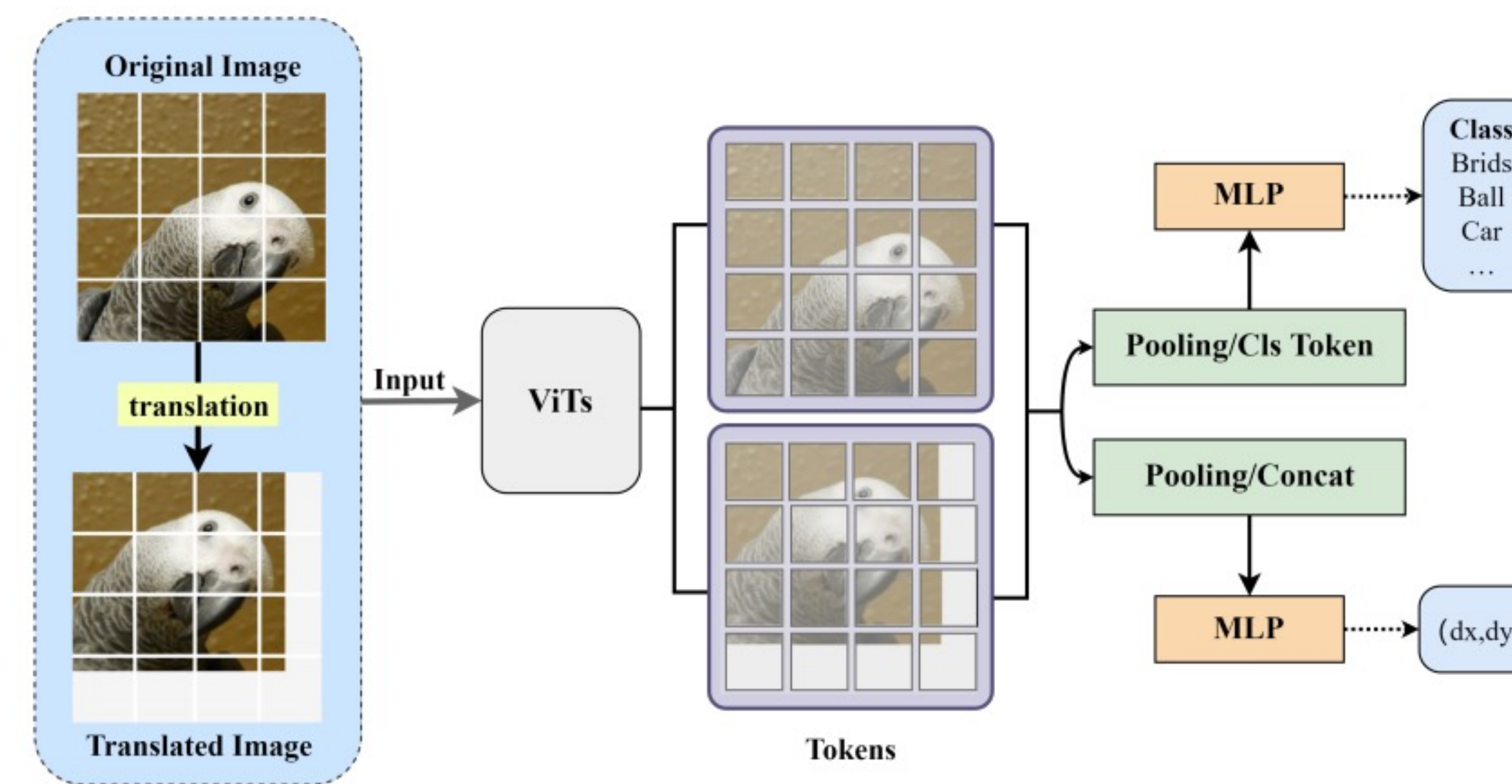
- ◆ **Translation Perceptibility** means that the system can work differently in different locations, but its output changes regularly with the input.

$$y = F(x)$$

$$y_{\text{trans}} = F(\text{trans}(x, TS))$$

$$TS = \text{MLP}(y, y_{\text{trans}})$$

## METHODS



**Pipe:** In order to guide the model in learning translation perceptibility, we first apply an arbitrary translation to the input image along any direction and generate the corresponding translation labels. Subsequently, both the original and translated images are fed into the network for processing. The output tokens are utilized for classification tasks as well as translation perception prediction tasks.

## RESULTS

Model	Imagenet-200	CIFAR10	CIFAR100	CINIC10	SVHN
ResNet18	53.32	90.44	64.49	77.79	96.78
ResNet56	56.51	94.65	74.44	85.34	97.61
ResNet101	59.77	95.27	76.18	86.81	97.82
EfficientNet B0	55.48	88.38	61.64	75.64	96.06
ViT(scratch)	54.07	93.58	73.81	83.73	97.82
SL-ViT	58.75	94.53	76.92	84.48	97.79
ViT-Drloc	54.44	81.00	58.29	71.50	94.02
ViT-vfsd(reproduce)	58.56	96.06	76.41	86.90	98.02
ViT-Trans(ours)	59.47	96.26	77.16	86.45	98.09
ViT-vfsd-Trans(ours)	<b>59.48</b>	<b>96.74</b>	<b>78.01</b>	<b>87.64</b>	<b>98.20</b>
Swin(scratch)	60.05	93.97	77.32	83.75	97.83
SL-Swin	64.95	94.93	79.99	87.22	97.92
Swin-Drloc	48.66	86.07	65.32	77.25	95.77
Swin-vfsd(reproduce)	64.28	96.52	80.67	87.96	98.02
Swin-Trans(ours)	62.27	96.87	80.28	88.26	98.15
Swin-vfsd-Trans(ours)	<b>65.05</b>	<b>97.08</b>	<b>81.25</b>	<b>88.63</b>	<b>98.17</b>

**Quantitative results:** Performance-wise, our method excels across datasets without extra inference parameters.

Model	Imagenet-200	Imagenet-100	Model	Imagenet-200	Imagenet-100
ViT(scratch)	54.07	62.56	Swin(scratch)	60.05	66.36
ViT-Trans(ours)	<b>59.47</b>	<b>65.50</b>	Swin-Trans(ours)	<b>62.27</b>	<b>69.00</b>
SL-ViT	58.75	66.96	SL-Swin	64.95	71.88
SL-ViT-Trans(ours)	<b>61.49</b>	<b>69.64</b>	SL-Swin-Trans(ours)	<b>66.80</b>	<b>74.81</b>
ViT-Drloc	54.44	64.52	Swin-Drloc	-	67.08
ViT-Drloc-Trans(ours)	<b>57.30</b>	<b>65.36</b>	Swin-Drloc-Trans(ours)	-	<b>69.96</b>
ViT-vfsd	58.56	65.38	Swin-vfsd	64.28	69.38
ViT-vfsd-Trans(ours)	<b>59.48</b>	<b>65.66</b>	Swin-vfsd-Trans(ours)	<b>65.05</b>	<b>71.30</b>

Model	WHU-RS19	UCMerced_LandUse	flowers102
ViT(scratch)	82.69	83.57	68.67
ViT-vfsd	89.76	91.66	69.01
ViT-Trans(ours)	91.83	94.52	73.65
ViT-vfsd-Trans(ours)	<b>93.27</b>	<b>95.24</b>	<b>74.72</b>
Swin(scratch)	85.10	88.81	79.13
Swin-vfsd	87.02	94.76	80.62
Swin-Trans(ours)	<b>94.71</b>	<b>97.62</b>	<b>85.37</b>
Swin-vfsd-Trans(ours)	<b>94.71</b>	96.43	84.66

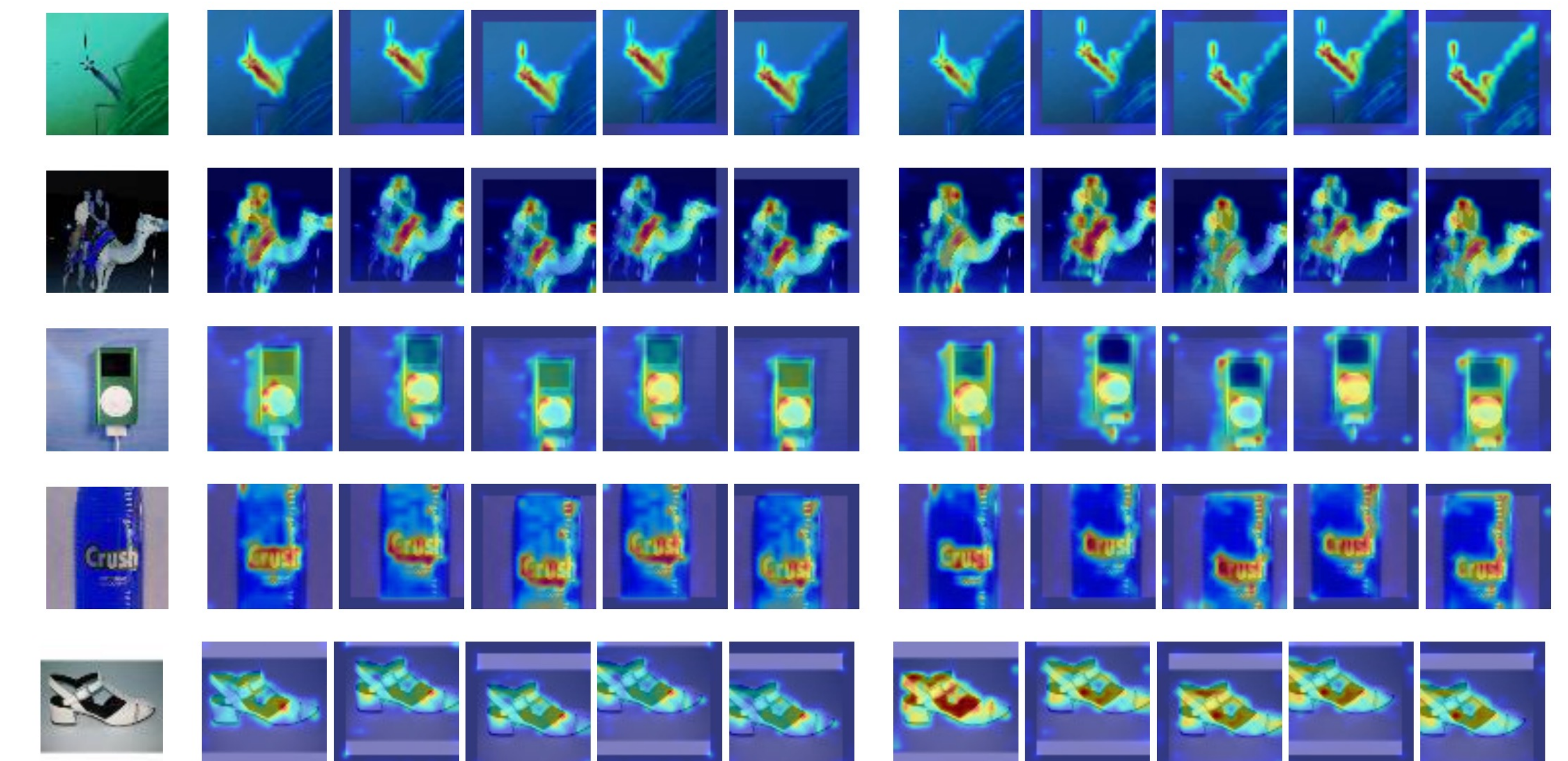
**Larger dimensions:** Our method excels with larger input dimensions.

Model	Imagenet-200	CIFAR10	CIFAR100
CaiT(scratch)	58.87	94.91	76.89
CaiT-vfsd	62.18	96.50	79.64
CaiT-Trans(ours)	62.00	96.73	80.66
CaiT-vfsd-Trans(ours)	<b>62.84</b>	<b>97.32</b>	<b>80.90</b>

**Models:** Aside from ViT/Swin, our method remains effective on CaiT as well.

**Extensibility:** Our method can be integrated with previous state-of-the-art methods to achieve even better performance.

## RESULTS



**Attention to salient regions.** Comparing our method (left) and vfsd[1] (right) using attention rollout on low-res Imagenet-100 samples. Our approach shows greater resilience to image translation.

## CONCLUSION

- ◆ We propose a self-supervised training method for Vision Transformers (ViTs) on small datasets, guiding ViTs to learn translation perceptibility
- ◆ Our approach outperforms state-of-the-art methods on small datasets with varying resolutions, and its benefits amplify as the input size increases.
- ◆ Our approach can integrate previously advanced methods, demonstrating its extensive extensibility.

## Reference

1. Gani H, Naseer M, Yaqub M. How to train vision transformer on small-scale datasets?[J]. arXiv preprint arXiv:2210.07240, 2022.