

# Learning and Transferring Visual Relation with Diffusion Transformers

Yan Gong<sup>1</sup>, Yiren Song<sup>2</sup>, Yicheng Li<sup>1</sup>, Chenglin Li<sup>1</sup>, Yin Zhang<sup>1\*</sup>

<sup>1</sup>Zhejiang University, <sup>2</sup>National University of Singapore



Presenter: Yan Gong

# Learning and Transferring Visual Relation with Diffusion Transformers

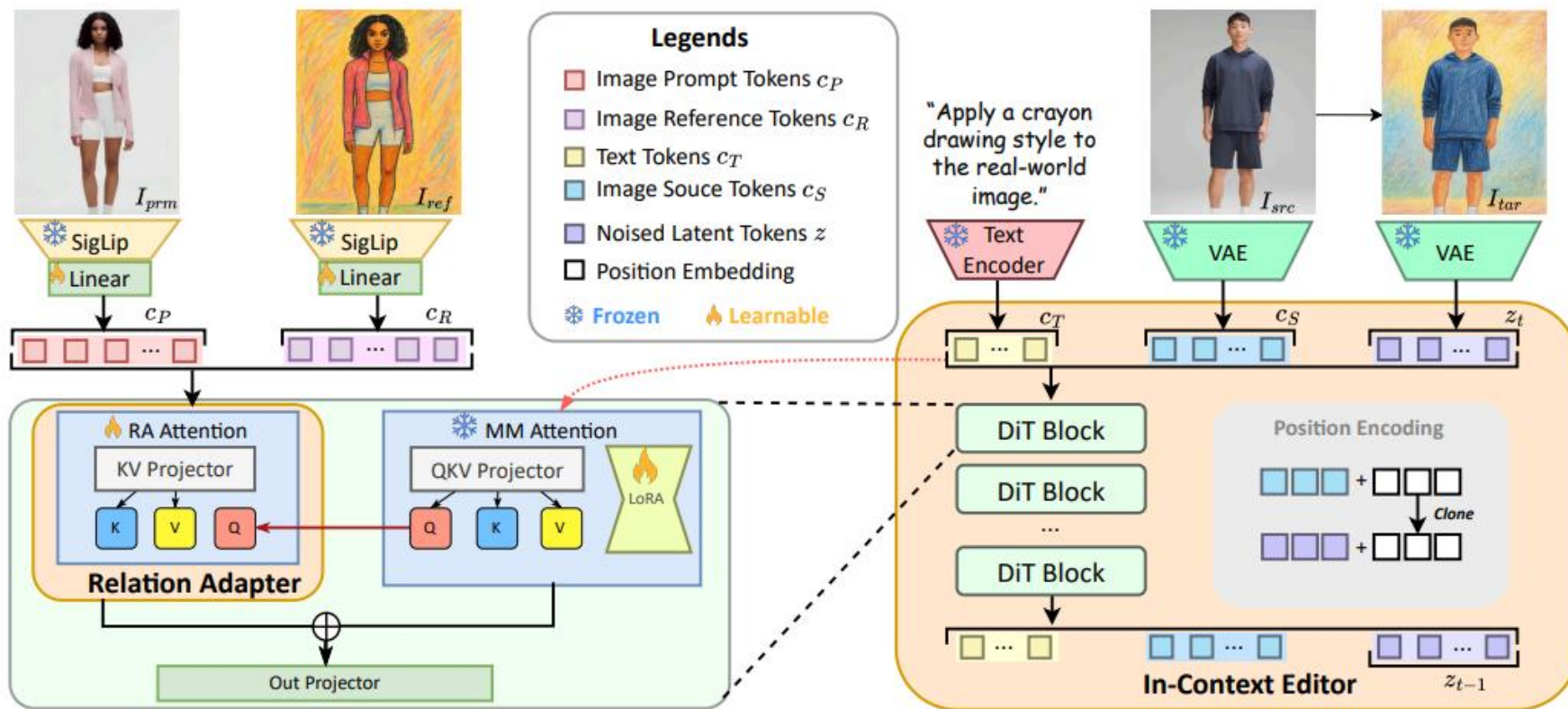
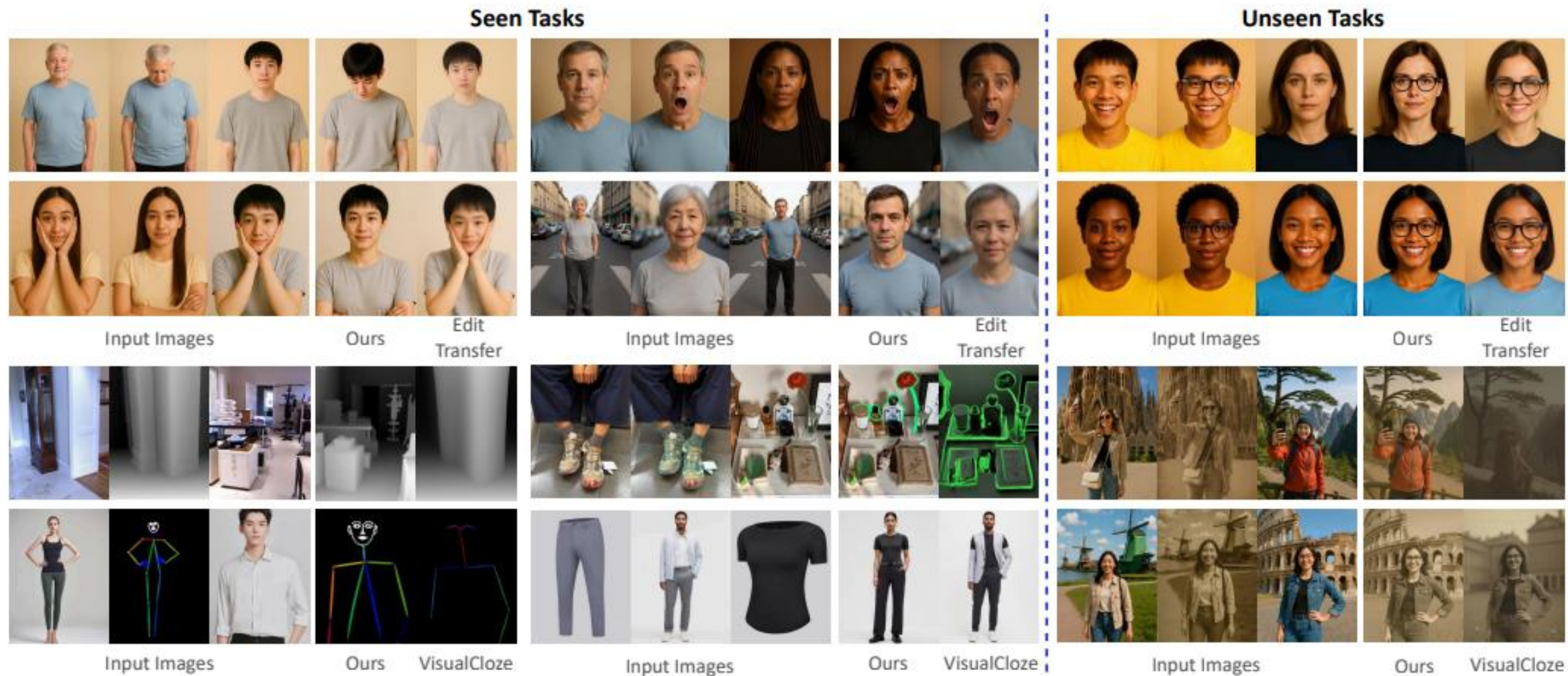


Figure 2: **The overall architecture and training paradigm of RelationAdapter.** We employ the RelationAdapter to decouple inputs by injecting visual prompt features into the MMAttention module to control the generation process. Meanwhile, a high-rank LoRA is used to train the In-Context Editor on a large-scale dataset. During inference, the In-Context Editor encodes the source image into conditional tokens, concatenates them with noise-added latent tokens, and directs the generation via the MMAttention module.



# Learning and Transferring Visual Relation with Diffusion Transformers



Compared with in-context based methods, our architecture demonstrates **image consistency, and editing effectiveness on both seen and unseen tasks.**

# Learning and Transferring Visual Relation with Diffusion Transformers

## ❖ Evaluation

- ◆ Outperforms baseline methods on common tasks, achieving lower MSE and higher CLIP-I, GPT-C, and GPT-A scores.
- ◆ The results are improved on both seen and unseen tasks, demonstrating its effectiveness and generalization ability in different editing scenarios.

Table 1: Quantitative Comparison of Baseline Methods Trained on a Common Task (ET: Edit Transfer, VC: VisualCloze). The best results are denoted as Bold.

Method	<i>MSE</i> ↓	<i>CLIP-I</i> ↑	<i>GPT-C</i> ↑	<i>GPT-A</i> ↑
EditTransfer	0.043	0.827	4.234	3.508
<b>Ours</b> ∩ <b>ET</b>	<b>0.020</b>	<b>0.905</b>	<b>4.437</b>	<b>4.429</b>
VisualCloze	0.049	0.802	3.594	3.411
<b>Ours</b> ∩ <b>VC</b>	<b>0.025</b>	<b>0.894</b>	<b>4.084</b>	<b>3.918</b>

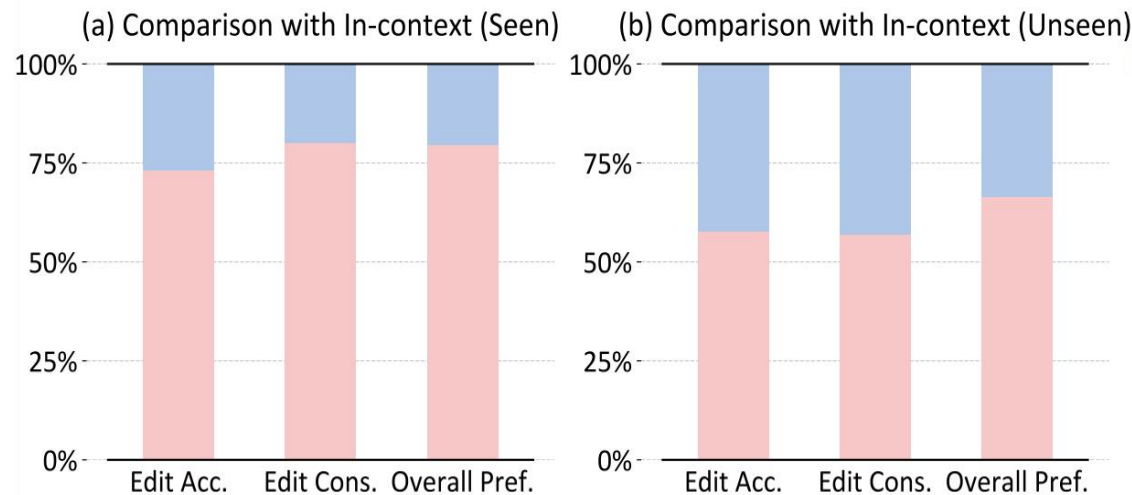
Table 2: Ablation Study on the Effectiveness of the RelationAdapter(RA) in Seen and Unseen Tasks (-S for Seen, -U for Unseen). The best results are denoted as Bold.

Method	<i>MSE</i> ↓	<i>CLIP-I</i> ↑	<i>GPT-C</i> ↑	<i>GPT-A</i> ↑
w/o RA -S	0.055	0.787	3.909	3.597
<b>Ours -S</b>	<b>0.044</b>	<b>0.852</b>	<b>4.079</b>	<b>4.106</b>
w/o RA -U	0.061	0.778	3.840	3.566
<b>Ours -U</b>	<b>0.053</b>	<b>0.812</b>	<b>4.187</b>	<b>4.173</b>

## ❖ Overall performance

- ◆ Efficiency improvement of our method over the In-Context baseline.
- ◆ The inference time is shortened (<9.0 seconds vs. 13.0+ seconds), the inference speed is increased by 30.8%, and the training speed is increased by 6.8%, with significant improvements in both training and inference efficiency.

## ❖ User Study



Thanks for listening!



# Learning and Transferring Visual Relation with Diffusion Transformers

Yan Gong<sup>1</sup>, Yiren Song<sup>2</sup>, Yicheng Li<sup>1</sup>, Chenglin Li<sup>1</sup>, Yin Zhang<sup>1\*</sup>

<sup>1</sup>Zhejiang University, <sup>2</sup>National University of Singapore



Resources



Arxiv

Feel free to connect!



Wechat



Snapchat