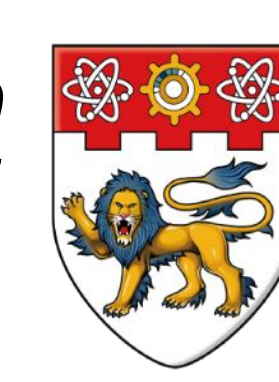




# TACA: Rethinking Cross-Modal Interaction in Multimodal Diffusion Transformer

Zhengyao Lv\*1, Tianlin Pan\*2,3, Chenyang Si\*2, Zhaoxi Chen4, Wangmeng Zuo5, Ziwei Liu4†, Kwan-Yee K. Wong1†  
HKU1 NJU2 UCAS3 NTU4 HIT5



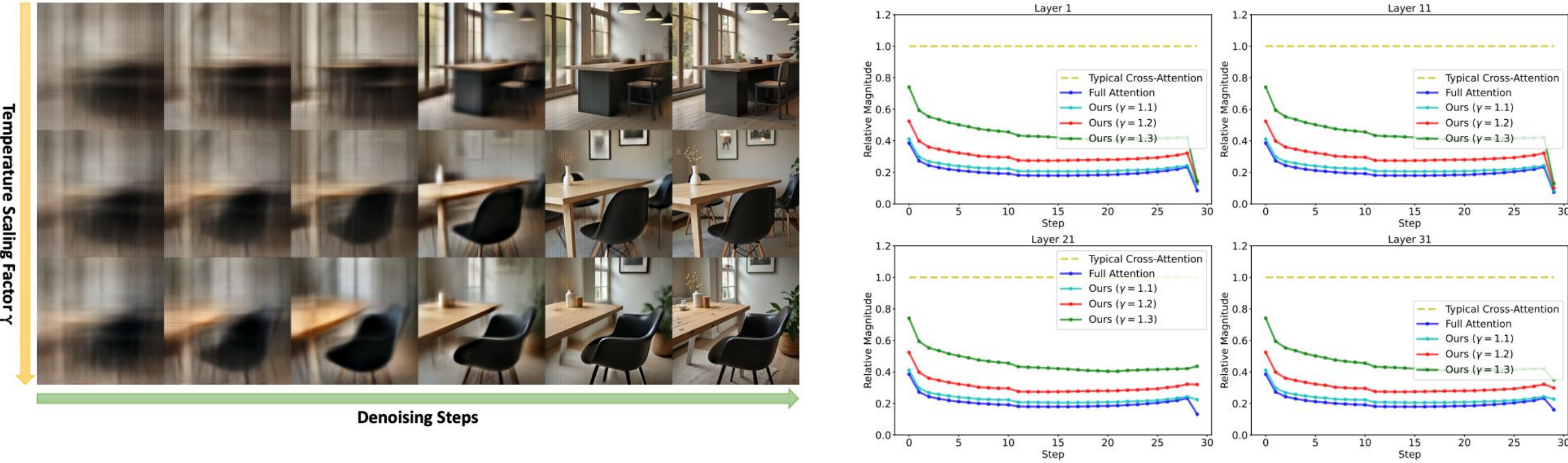
NANYANG  
TECHNOLOGICAL  
UNIVERSITY  
SINGAPORE



ICCV HONOLULU  
OCT 19-23, 2025 HAWAII

## Motivation

- Task:** We propose **TACA**, a parameter-efficient method that dynamically rebalances cross-modal attention in multimodal diffusion transformers to **improve text-image alignment**.
- Motivation:**
  - In **MM-DiT**, the cross-modal attention between visual and text tokens is suppressed **due to the significant imbalance in their numbers**
  - The attention weighting in DiT **does not adapt to the varying needs** of the denoising process across different timesteps.



## Method

### Temperature-Adjusted Cross-modal Attention (TACA)

$$\text{In MM-DiT: } \text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{D}}\right)V,$$

where the  $QK^T$  term can be expanded to

$$QK^T = \begin{pmatrix} W_c^Q c (W_c^K c)^T & W_c^Q c (W_x^K x)^T \\ W_x^Q x (W_c^K c)^T & W_x^Q x (W_x^K x)^T \end{pmatrix} = \begin{pmatrix} Q_{\text{txt}} K_{\text{txt}}^T & Q_{\text{txt}} K_{\text{vis}}^T \\ Q_{\text{vis}} K_{\text{txt}}^T & Q_{\text{vis}} K_{\text{vis}}^T \end{pmatrix}.$$

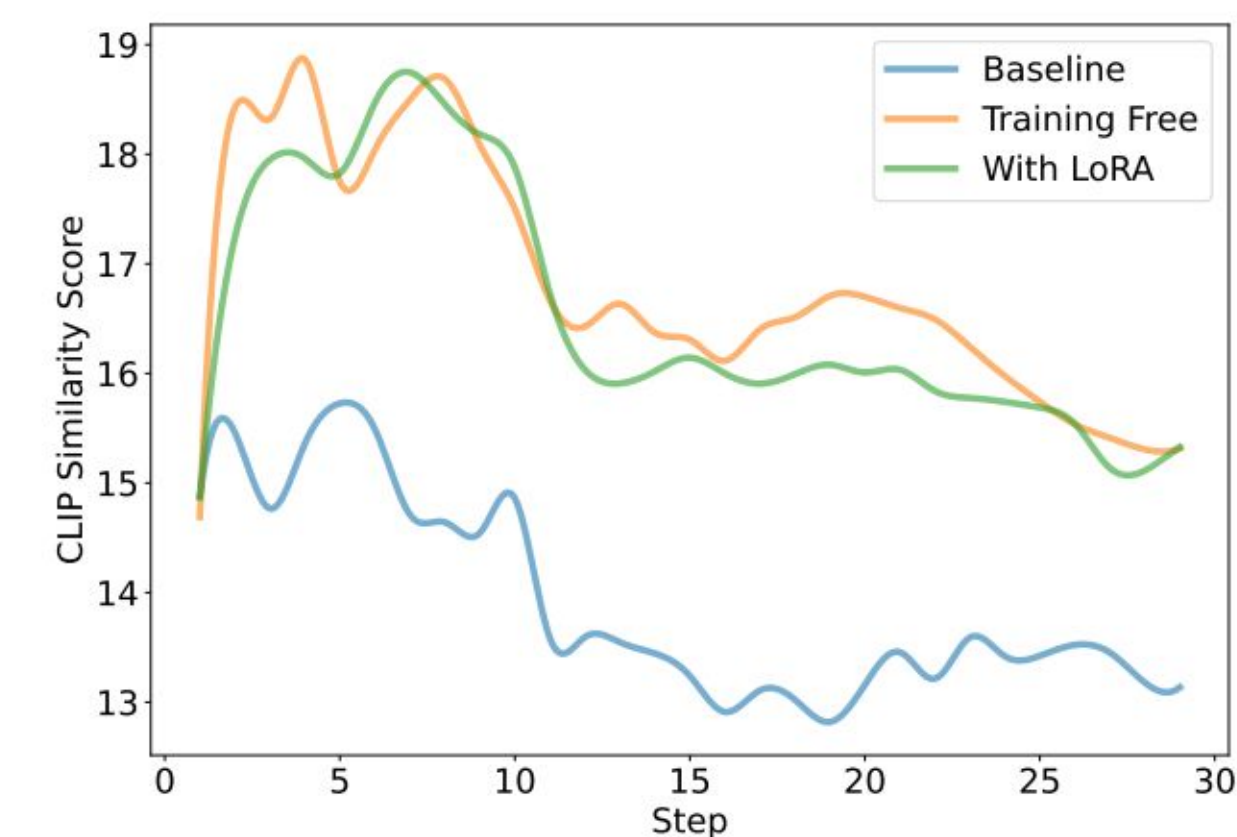
$$\text{After Softmax: } P_{\text{vis-txt}}^{(i,j)} = \frac{e^{\gamma s_{ij}^{\text{vt}}/\tau}}{\sum_{k=1}^{N_{\text{txt}}} e^{\gamma s_{ik}^{\text{vt}}/\tau} + \sum_{k=1}^{N_{\text{vis}}} e^{\gamma s_{ik}^{\text{vv}}/\tau}},$$

where  $s_{ik}^{\text{vt}} = Q_{\text{vis}}^{(i)} K_{\text{txt}}^{(k)} / \sqrt{D}$  and  $s_{ik}^{\text{vv}} = Q_{\text{vis}}^{(i)} K_{\text{vis}}^{(k)} / \sqrt{D}$

- Only enhance the early timesteps:

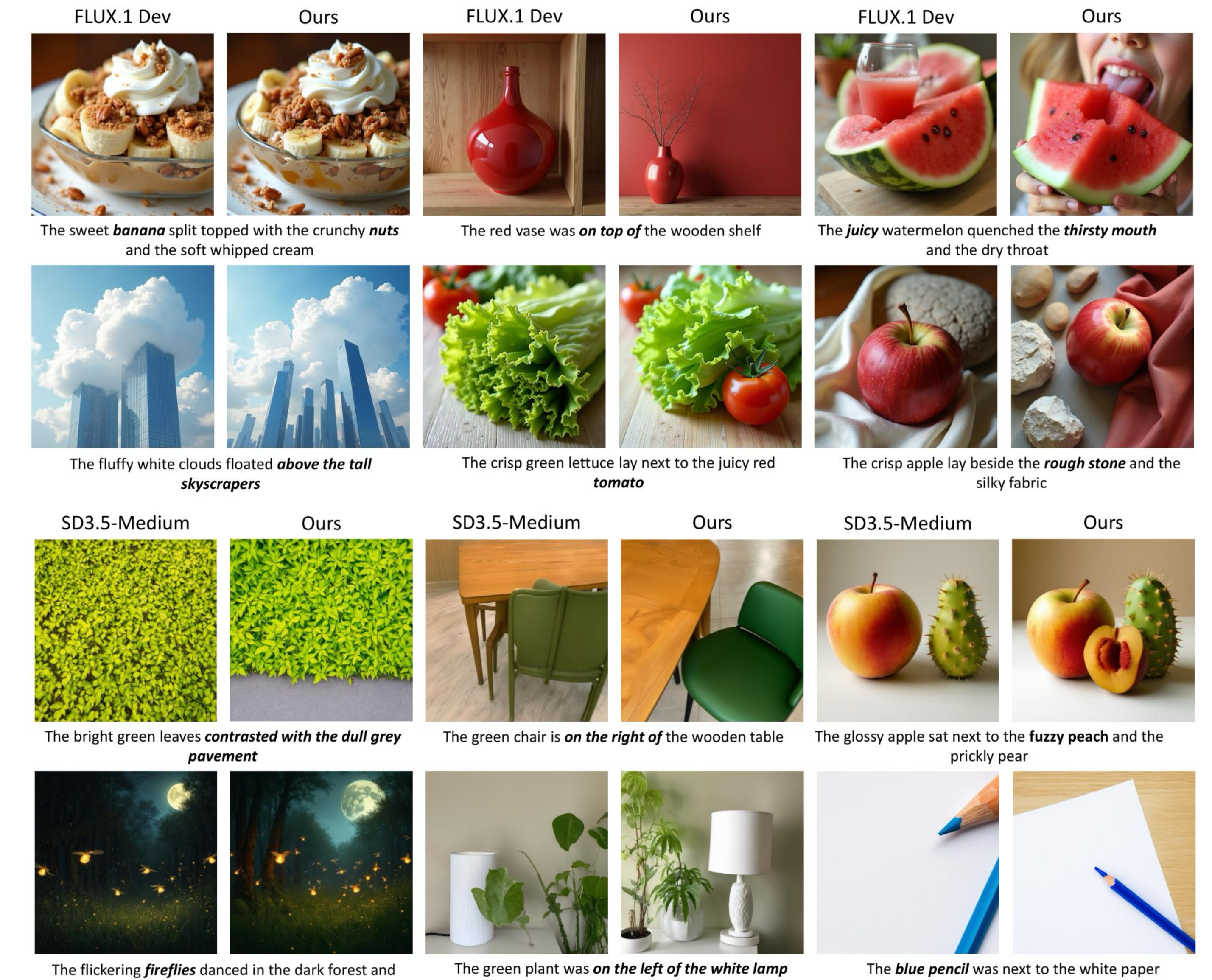
$$\gamma(t) = \begin{cases} \gamma_0 & t \geq t_{\text{thresh}} \\ 1 & t < t_{\text{thresh}} \end{cases}$$

- We trained a LoRA to mitigate artifacts caused by the shift in the output distribution of each attention layer



## Results

Model	Attribute Binding			Object Relationship		Complex ↑
	Color ↑	Shape ↑	Texture ↑	Spatial ↑	Non-Spatial ↑	
FLUX.1-Dev	0.7678	0.5064	0.6756	0.2066	0.3035	0.4359
FLUX.1-Dev + TACA (r = 64)	<b>0.7843</b>	<b>0.5362</b>	<b>0.6872</b>	<b>0.2405</b>	0.3041	<b>0.4494</b>
FLUX.1-Dev + TACA (r = 16)	0.7842	0.5347	0.6814	0.2321	<b>0.3046</b>	0.4479
SD3.5-Medium	0.7890	0.5770	0.7328	0.2087	0.3104	0.4441
SD3.5-Medium + TACA (r = 64)	<b>0.8074</b>	<b>0.5938</b>	<b>0.7522</b>	<b>0.2678</b>	0.3106	0.4470
SD3.5-Medium + TACA (r = 16)	0.7984	0.5834	0.7467	0.2374	<b>0.3111</b>	<b>0.4505</b>



Project  
Page

GitHub  
Repo

Paper

## Investigation

- Temperature scaling helps visual-text alignment.**
- To mitigate the suppression of cross-attention caused by the dominance of visual tokens ( $N_{\text{vis}} \gg N_{\text{txt}}$ ), we amplify the logits of visual-text interactions through a temperature coefficient gamma

