



**ICLR**  
International Conference On  
Learning Representations

# Ranking-Aware Adapter for Text-Driven Image Ordering with CLIP

Wei-Hsiang Yu<sup>1</sup> Yen-Yu Lin<sup>1</sup> Ming-Hsuan Yang<sup>2</sup> Yi-Hsuan Tsai<sup>3</sup>

<sup>1</sup>National Yang Ming Chiao Tung University   <sup>2</sup>UC Merced   <sup>3</sup>Atmanity Inc.

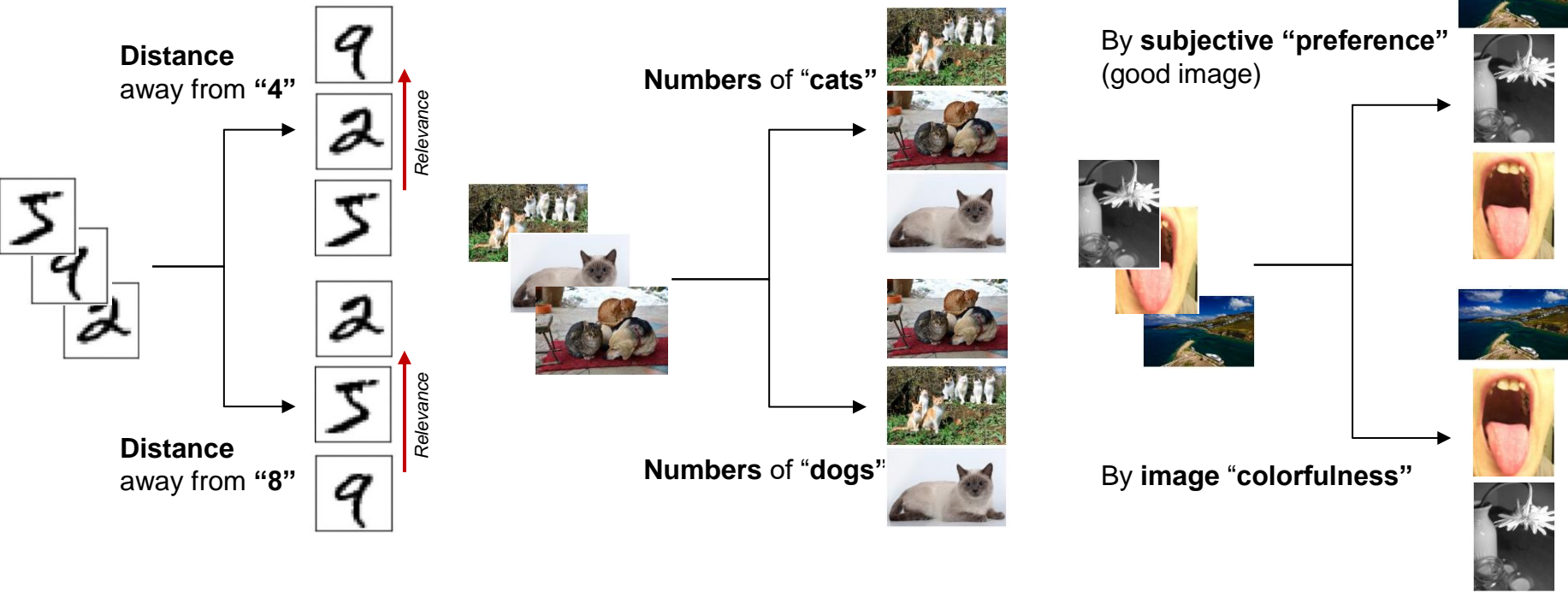
Website



# **Ranking-Aware Adapter** for **Text-Driven** **Image Ordering** with **CLIP**

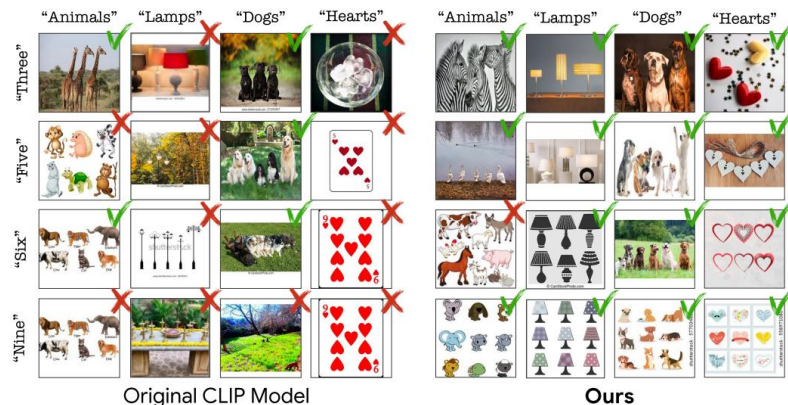
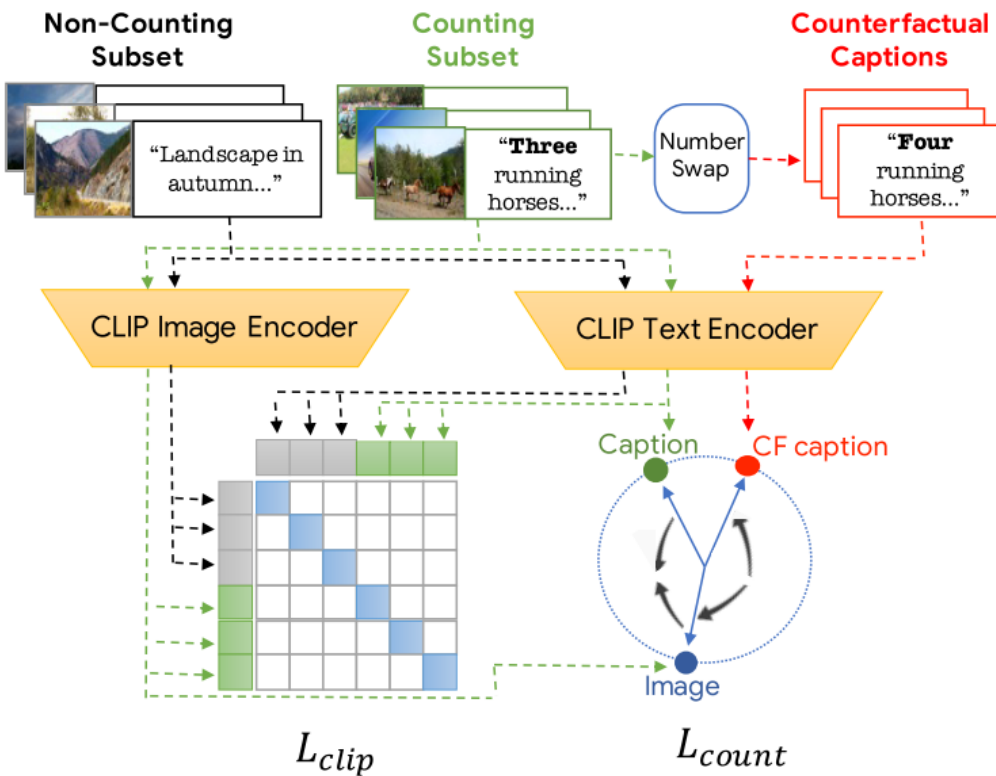
# Problem formulation

- Arrange a set of images according to user requirements.



# Previous efforts to text-driven image ordering

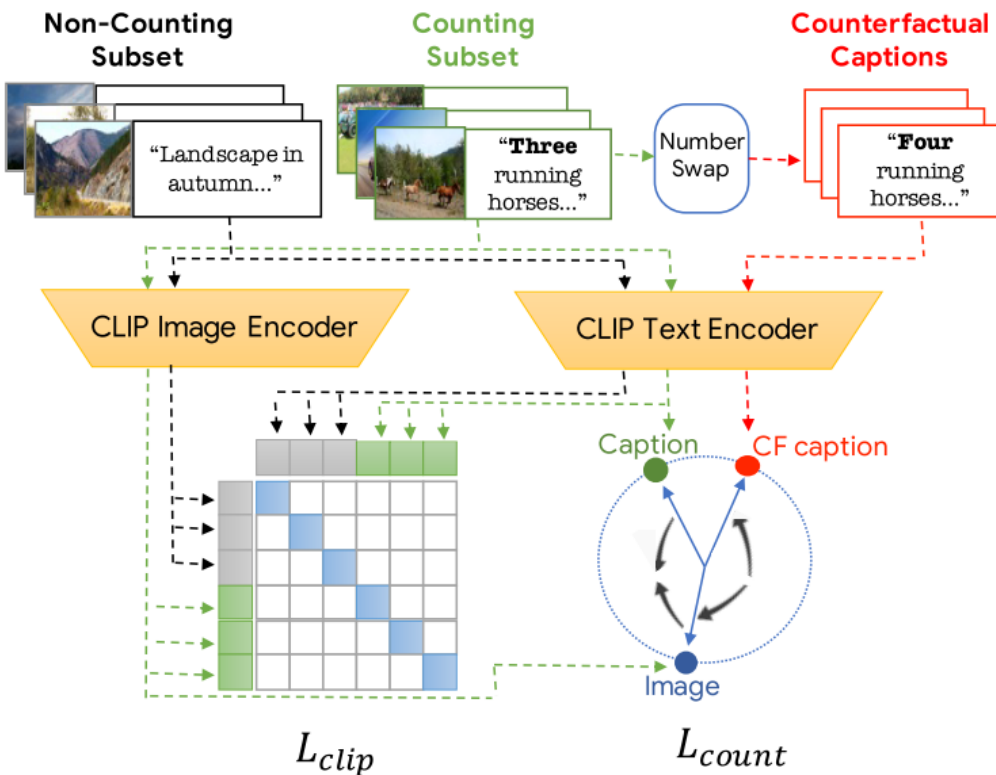
- Re-align “Number” to image using contrastive learning



*Teach CLIP to Count to Ten, ICCV'23*

# Previous efforts to text-driven image ordering

- Re-align “Number” to image using contrastive learning



Values span a wide range? (e.g., facial age)



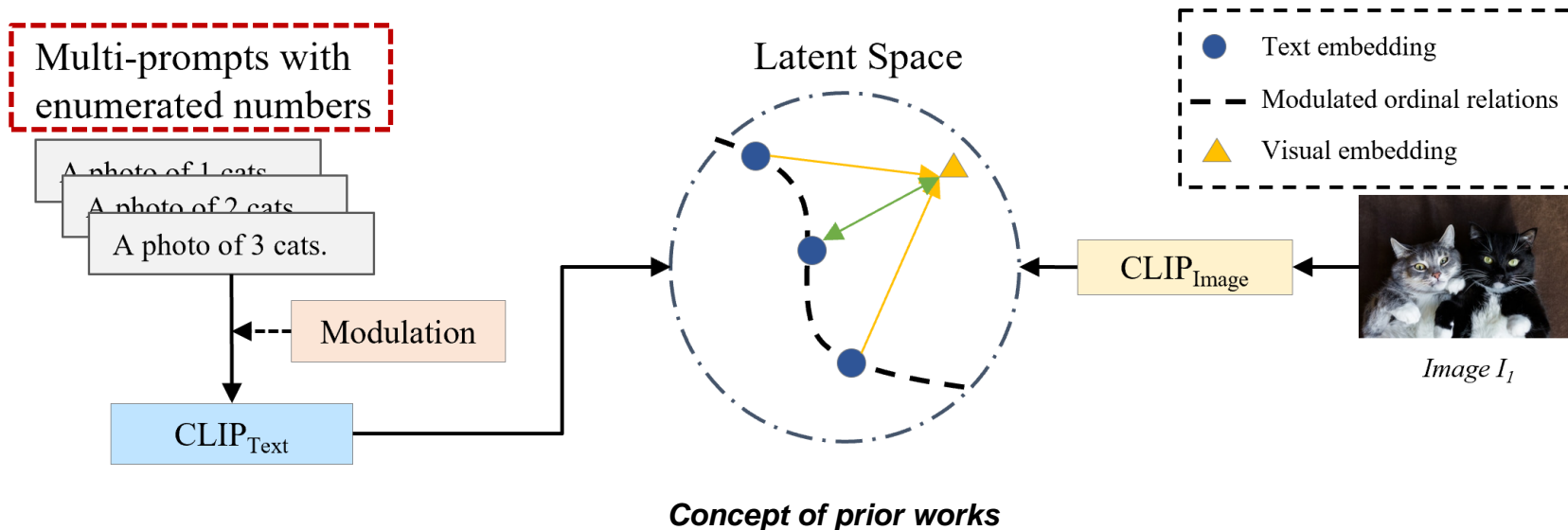
Continuous and subjective values?  
(e.g., image contrast, image aesthetics, etc.)



*Teach CLIP to Count to Ten, ICCV'23*

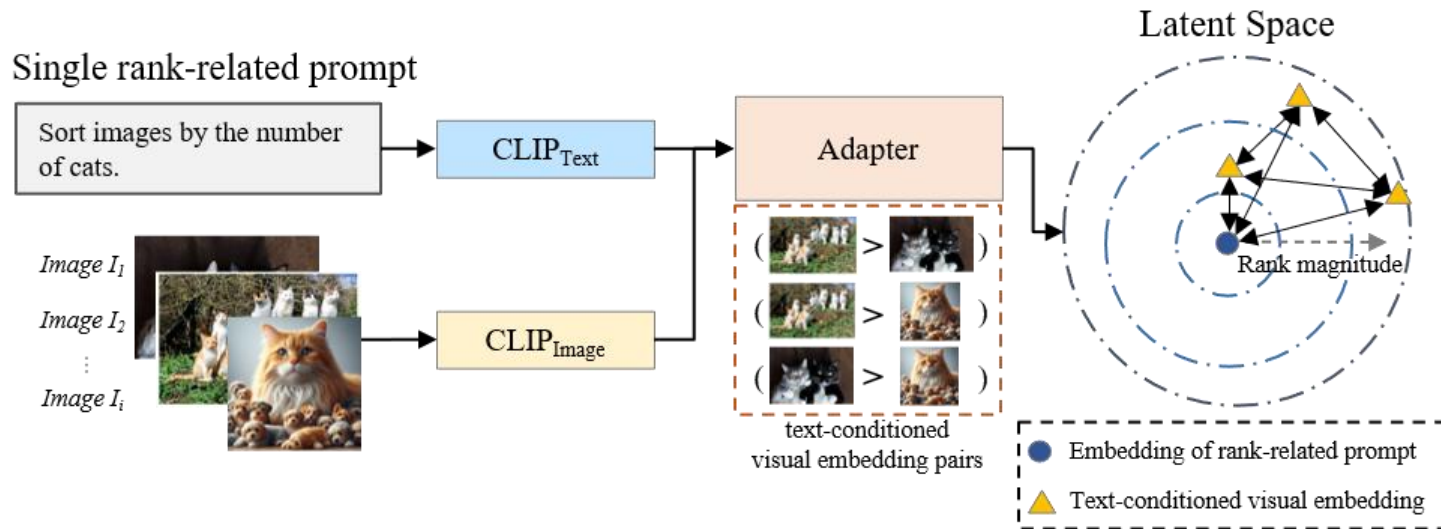
# Motivation 1 – Contrastive learning with enumerated numbers

- Existing methods focus on aligning numerical text prompts to the image
  - Difficult to handle continuous values spanning a wide range - binning
  - Requires to enumerate all attribute-value combinations – inefficient



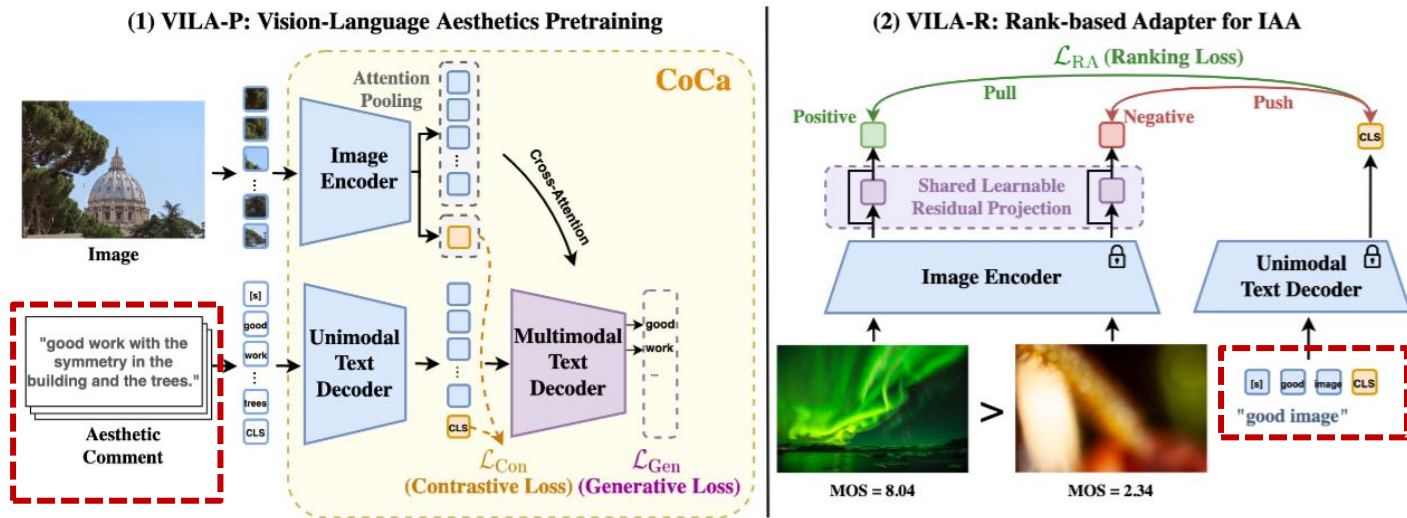
# Idea – Single rank-related prompt with image pairs

- Introduce a “**single rank-related prompt**” – The attribute used to sort images.
- Generate text-conditioned image embedding pairs and supervise using their numerical order with a learning-to-rank framework.



## Motivation 2 – Learning-to-Rank in Visual-Text Alignment

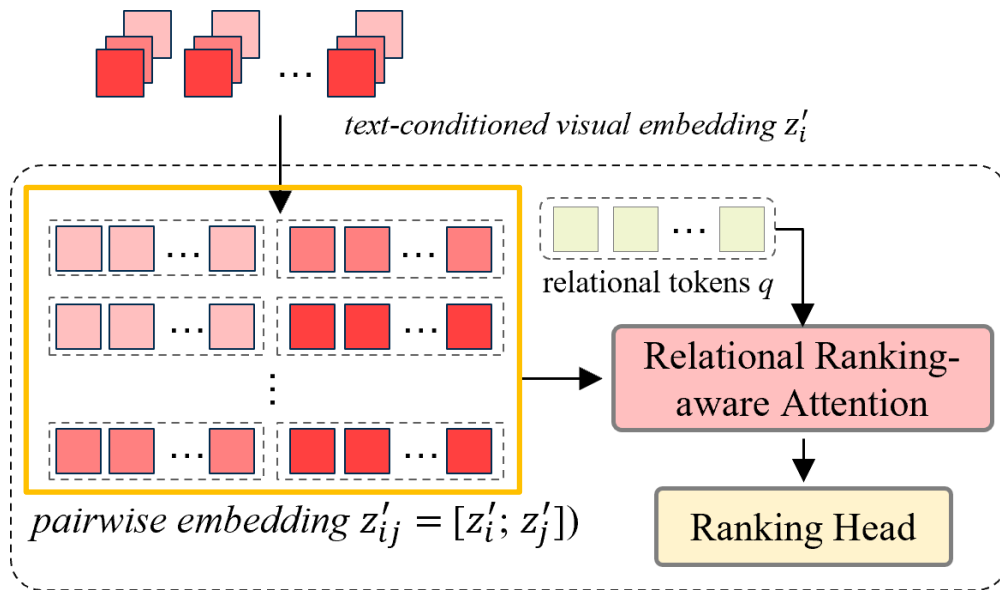
- Existing methods pretrain the text encoder on a **task-specific dataset** and then use a learning-to-rank framework to align images to specific attributes.
  - Requires fine-tuning the text encoder on a specific dataset (e.g., aesthetic comment).





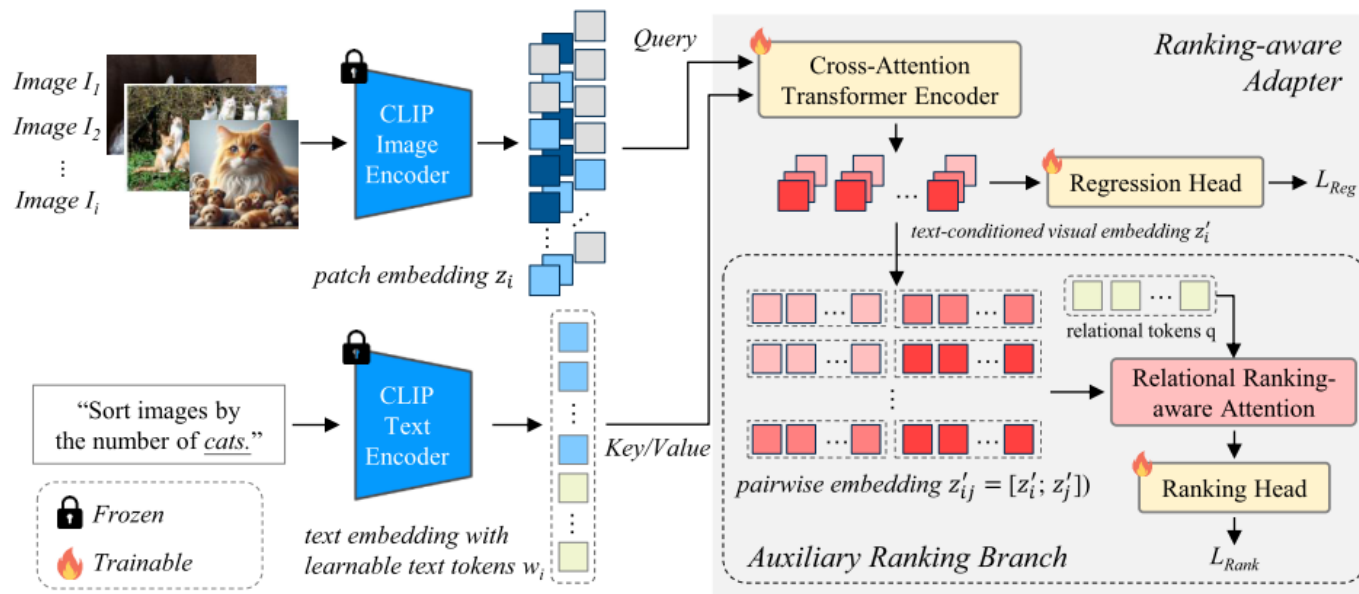
# Idea – Learning from text-conditioned visual distinctions

- Paired embedding difference  $\propto$  queried value difference between the image pair
  - Design a relational ranking-aware attention module to extract the text-conditioned visual distinctions



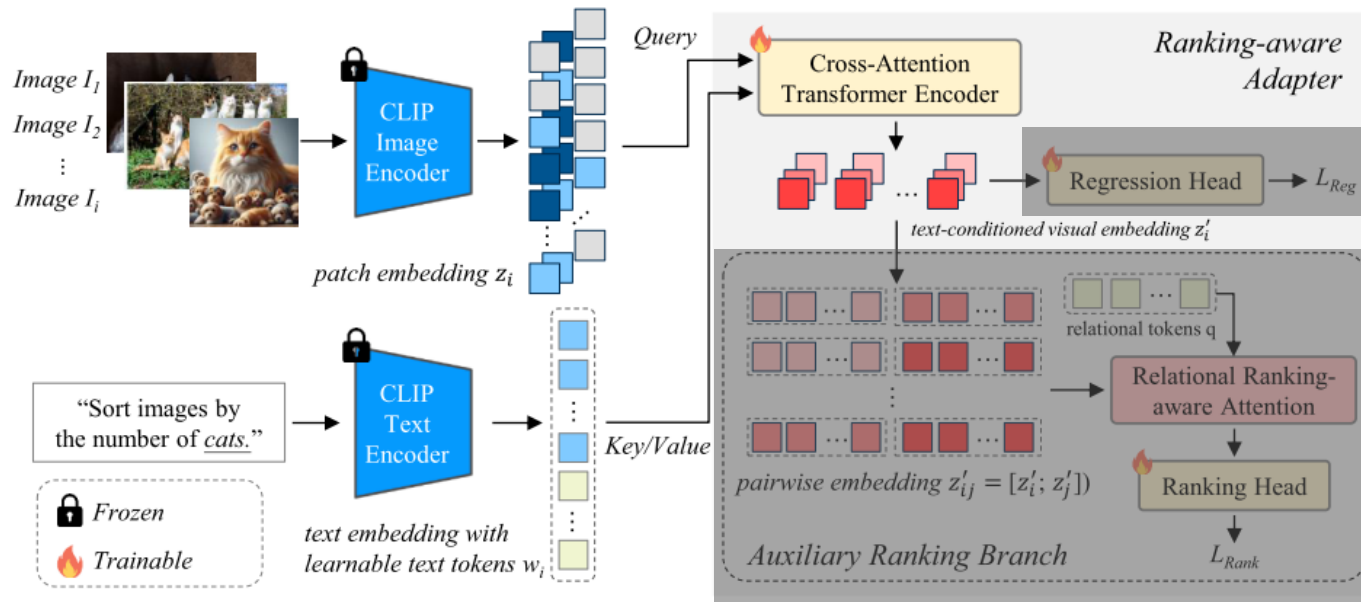
# Our Approach

## - Ranking-Aware Adapter for Text-Driven Image Ordering with CLIP

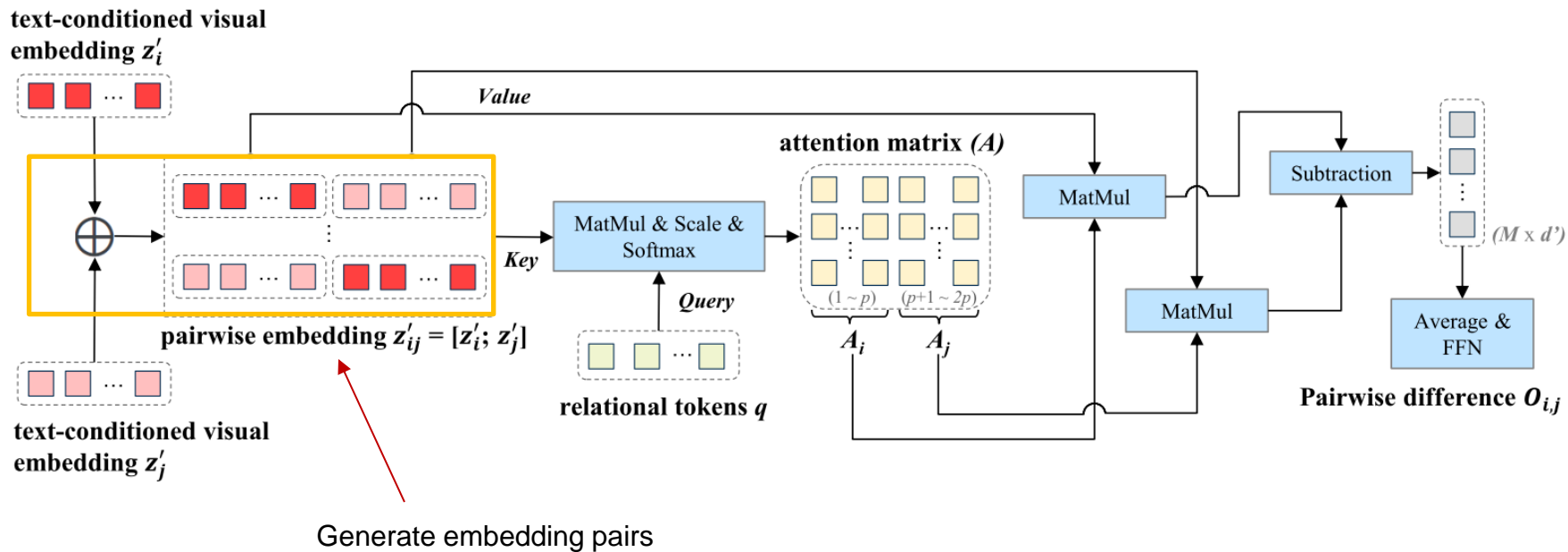


# Our Approach

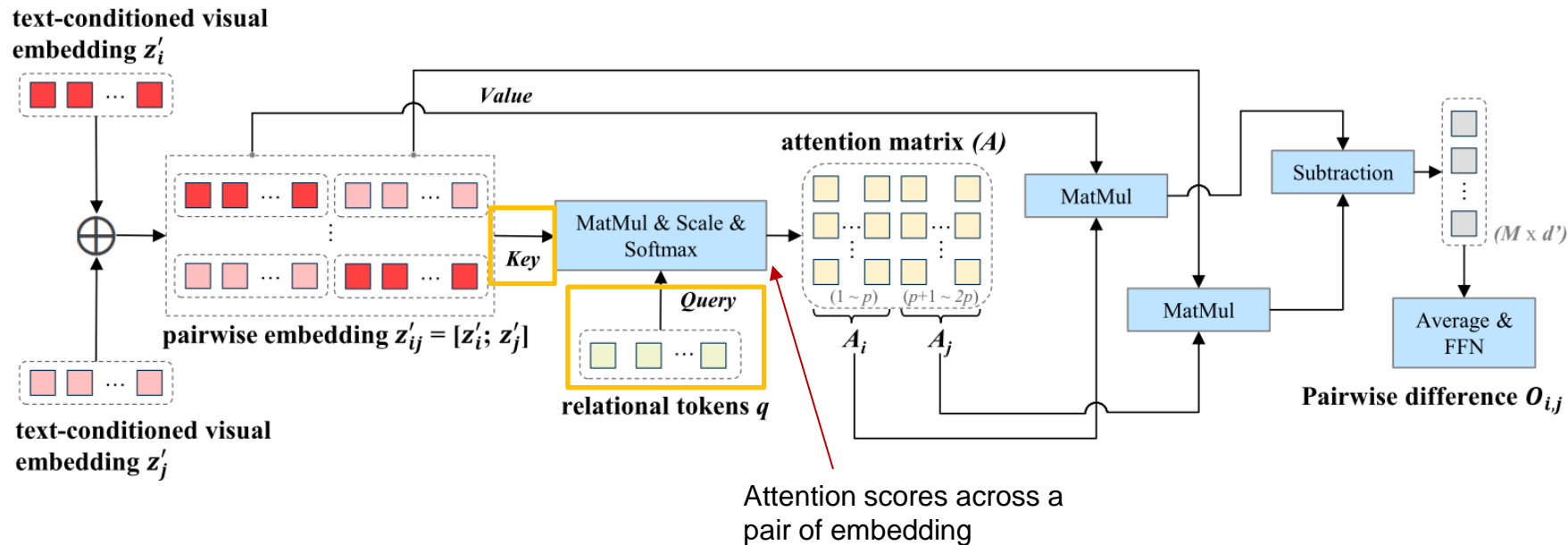
## - Ranking-Aware Adapter for Text-Driven Image Ordering with CLIP



# Relational Ranking-aware Attention

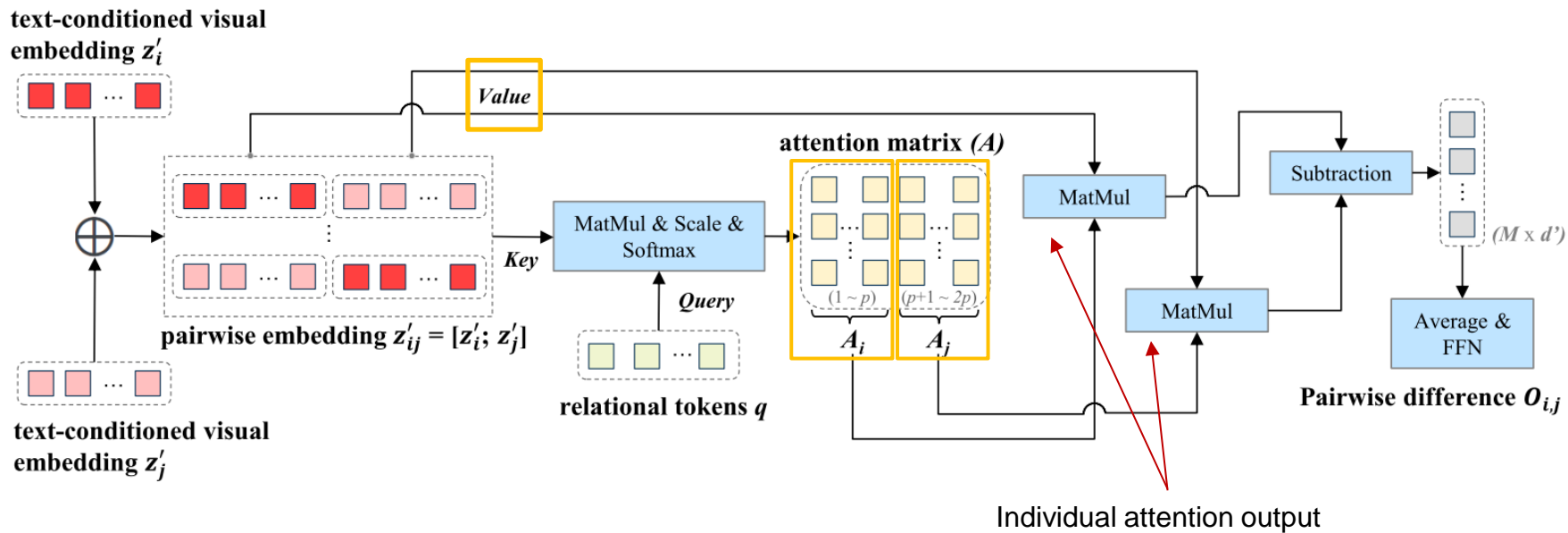


# Relational Ranking-aware Attention



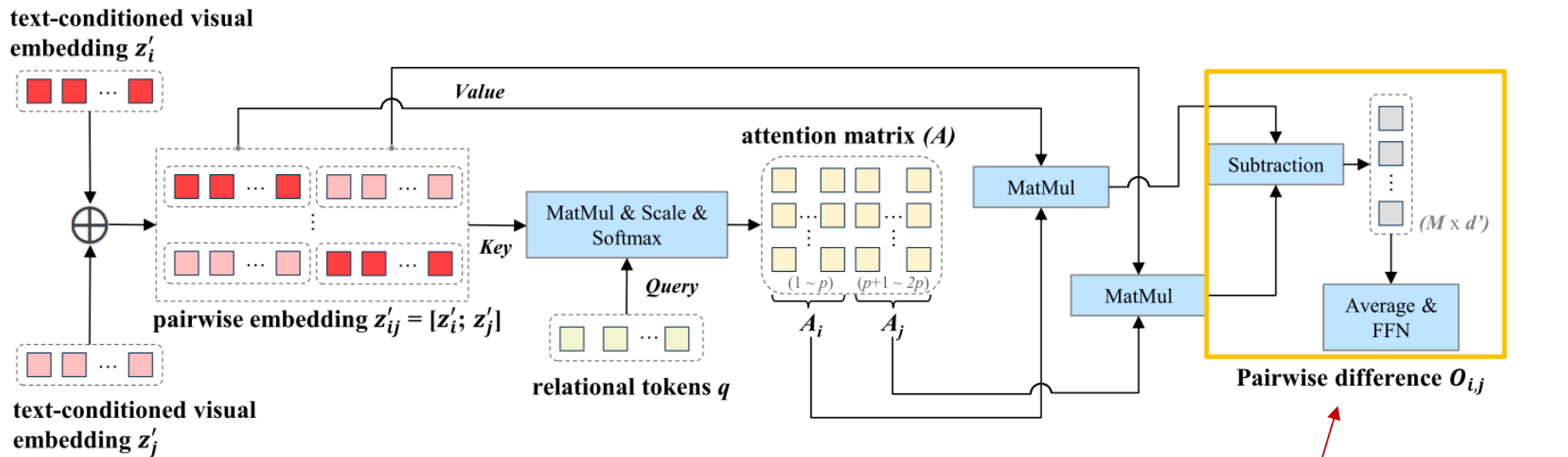
$$A = \text{Softmax}\left(\frac{q \cdot (k_i \oplus k_j)^T}{\sqrt{d'}}\right) = \text{Softmax}\left(\frac{q \cdot K^T}{\sqrt{d'}}\right)$$

# Relational Ranking-aware Attention



$$O_i = A_i \cdot V_i \quad \text{and} \quad O_j = A_j \cdot V_j.$$

# Relational Ranking-aware Attention

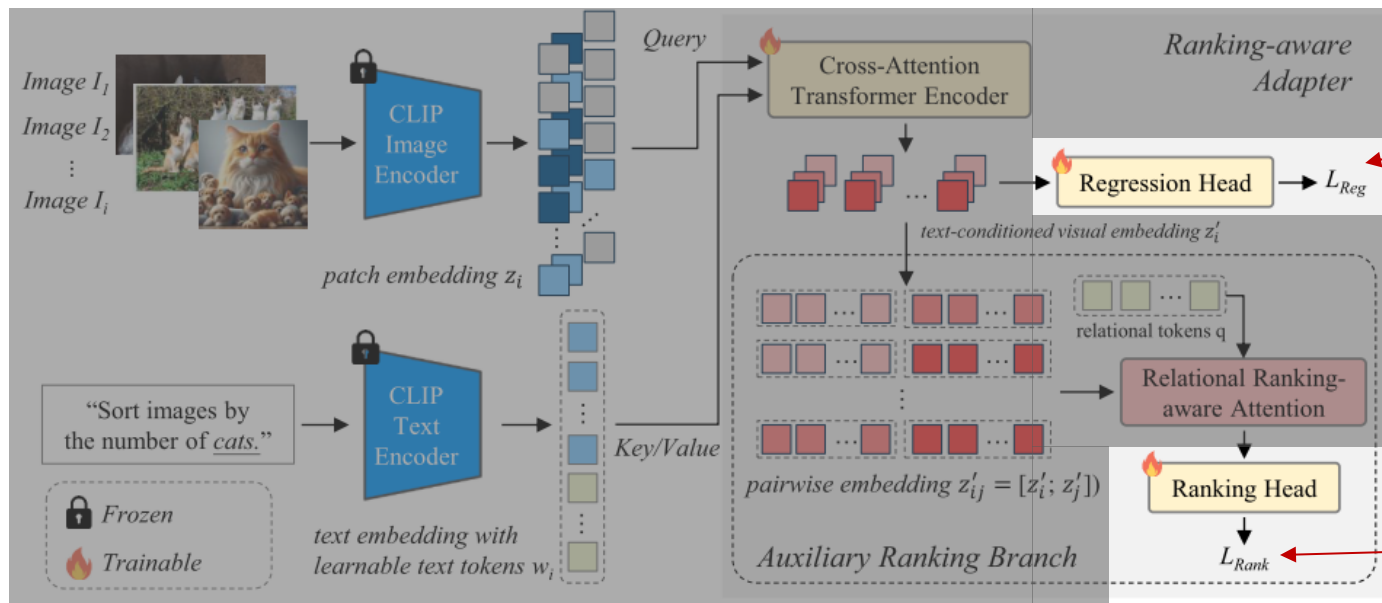


Enforce the relative difference

$$O_{i,j} = FFN\left(\frac{\sum_{m=1}^M (O_{i,m} - O_{j,m})}{M}\right)$$

# Our Approach

## - Ranking-Aware Adapter for Text-Driven Image Ordering with CLIP



Predict image ranking score

$$L_{reg} = \begin{cases} \frac{1}{2}(y_i - s_i)^2, & \text{if } |y_i - s_i| < 1, \\ |y_i - s_i| - 0.5, & \text{otherwise.} \end{cases}$$

Encourage relative ordering

$$L_{rank}(y, O_{i,j}) = \sum_{\{(i,j)|y_i > y_j\}} \max(0, 1 - O_{i,j}).$$



# Experiment Settings

- Datasets
  - **Visual Quantity**: COCO (object count sorting)
  - **Visual Quality**: *KonIQ-A-10k* (image quality) and *AVA* (image aesthetics)
  - **Perceptual Concepts**: *Adience* (facial age) and *historical colored image* (photo taken decade)
- Evaluation Metrics
  - Object count sorting and IQA/IAA – **Pearson's correlation** and **Spearman's correlation**
  - Facial aging and historical colored image aging – **Mean absolute error** (MAE) and **Accuracy**

# Quantitative Results

Method	Fine-tuning	PLCC ( $\uparrow$ )	SRCC ( $\uparrow$ )
BLIP-2		0.284	0.252
Flamingo (10-shot)		0.033	0.031
InstructBLIP		0.509	0.485
VLM-VILA		0.558	0.507
Zero-shot CLIP		0.026	0.001
CountingCLIP	✓	0.251	0.422
Paiss et al. (2023)			
Ours	✓	<b>0.624</b>	<b>0.557</b>

Method	Adience		HCI	
	Accuracy (%)	MAE	Accuracy (%)	MAE
Zero-shot CLIP	43.3 (3.6)	0.80 (0.02)	26.1 (0.6)	1.48 (0.03)
CoOp	60.6 (5.5)	0.50 (0.08)	51.9 (2.6)	0.76 (0.06)
OrdinalCLIP	61.2 (4.2)	0.47 (0.06)	56.4 (1.7)	0.67 (0.03)
L2RCLIP	<b>66.2 (4.4)</b>	<b>0.36 (0.05)</b>	67.2 (1.6)	0.43 (0.03)
NumCLIP	-	-	69.6 (2.0)	0.35 (0.03)
InstructBLIP	63.7	0.41	30.9	0.96
Ours	65.2 (2.9)	<b>0.36 (0.03)</b>	<b>72.8 (2.6)</b>	<b>0.32 (0.03)</b>

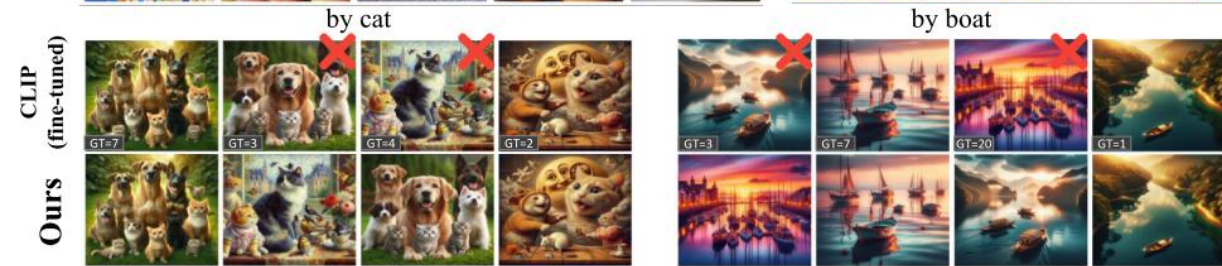
Method	Task-related pertaining	Fine-tuning	KonIQ-10k		AVA Dataset	
			PLCC ( $\uparrow$ )	SRCC ( $\uparrow$ )	PLCC ( $\uparrow$ )	SRCC ( $\uparrow$ )
<b>Purely vision-based (task-specific)</b>						
MUSIQ <a href="#">Ke et al. (2021)</a>		✓	0.924	0.937	0.726	0.738
<b>VLM-based (task-specific)</b>						
VILA-P <a href="#">Ke et al. (2023)</a>	✓		-	-	0.657	0.663
VILA-R <a href="#">Ke et al. (2023)</a>	✓	✓	<b>0.919</b>	<b>0.932</b>	<b>0.774</b>	<b>0.774</b>
CLIP (fine-tuned)		✓	0.245	0.216	0.162	0.160
InstructBLIP <a href="#">Dai et al. (2023)</a>			0.211	0.163	0.229	0.226
CLIP-IQA <a href="#">Wang et al. (2023c)</a>			0.695	0.727	0.420	0.415
CLIP-IQA+ <a href="#">Wang et al. (2023c)</a>		✓	0.895	0.909	0.677	0.587
<a href="#">Hentschel et al. (2022)</a>		✓	-	-	0.731	0.741
Ours		✓	<b>0.919</b>	<b>0.911</b>	<b>0.760</b>	<b>0.747</b>

# Qualitative Results

## Object count sorting task



Same set of images w/ different query



On the unseen domain  
(Generated Arts)

## IQA/IAA task

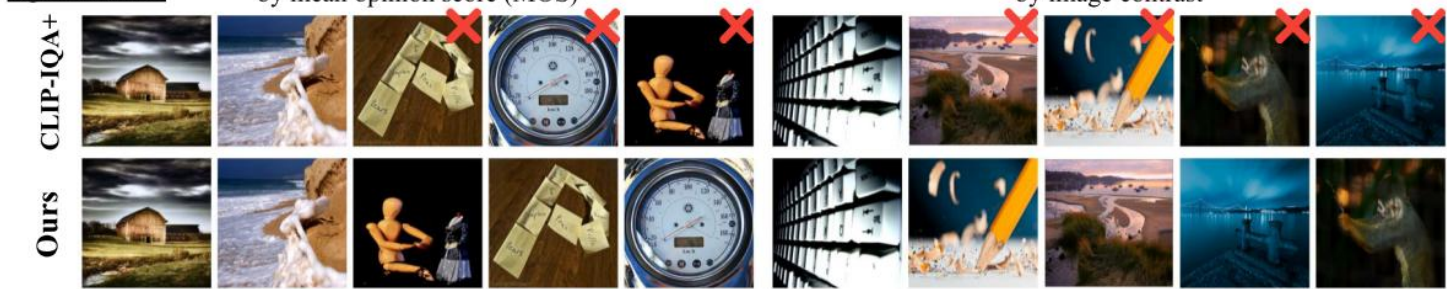


# Qualitative Results

## Object count sorting task



## IQA/IAA task

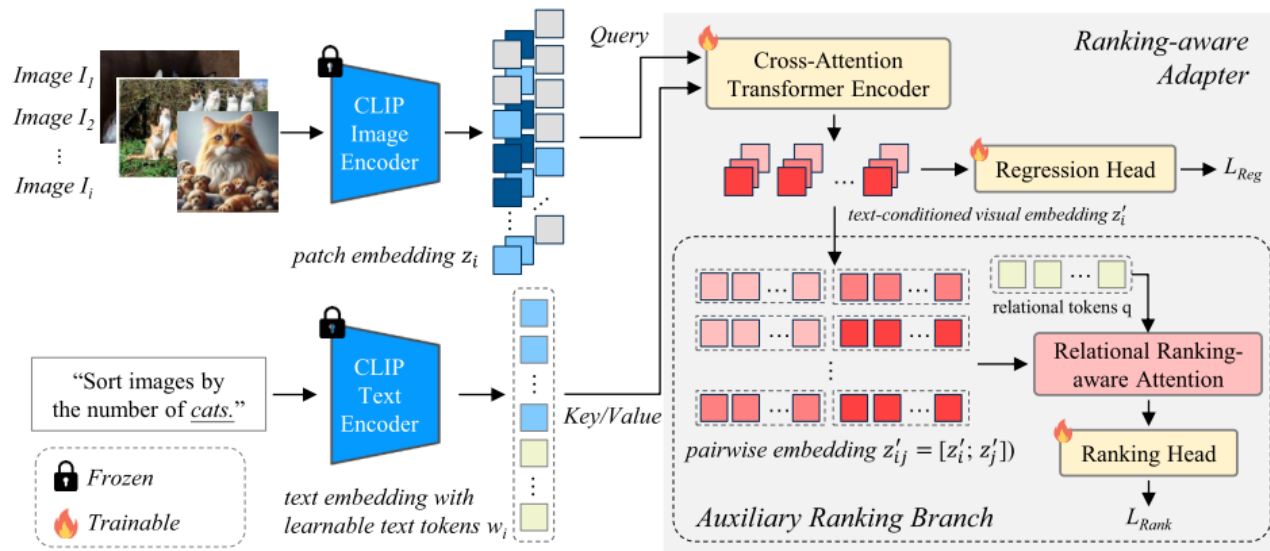


On IQA Tasks

# Ablation Study

LTR paradigm	Ranking Head	Ranking-aware Attention	HCI		Object count sorting		
			MAE ( $\downarrow$ )		PLCC ( $\uparrow$ )		SRCC ( $\uparrow$ )
-	-	-	1.113		0.251		0.422
✓	-	-	0.402		0.612		0.538
✓	✓	-	0.355	(+11.69%)	0.619	(+1.14%)	0.536 (-0.37%)
✓	✓	✓	0.317	(+21.14%)	0.624	(+1.96%)	0.557 (+3.53%)

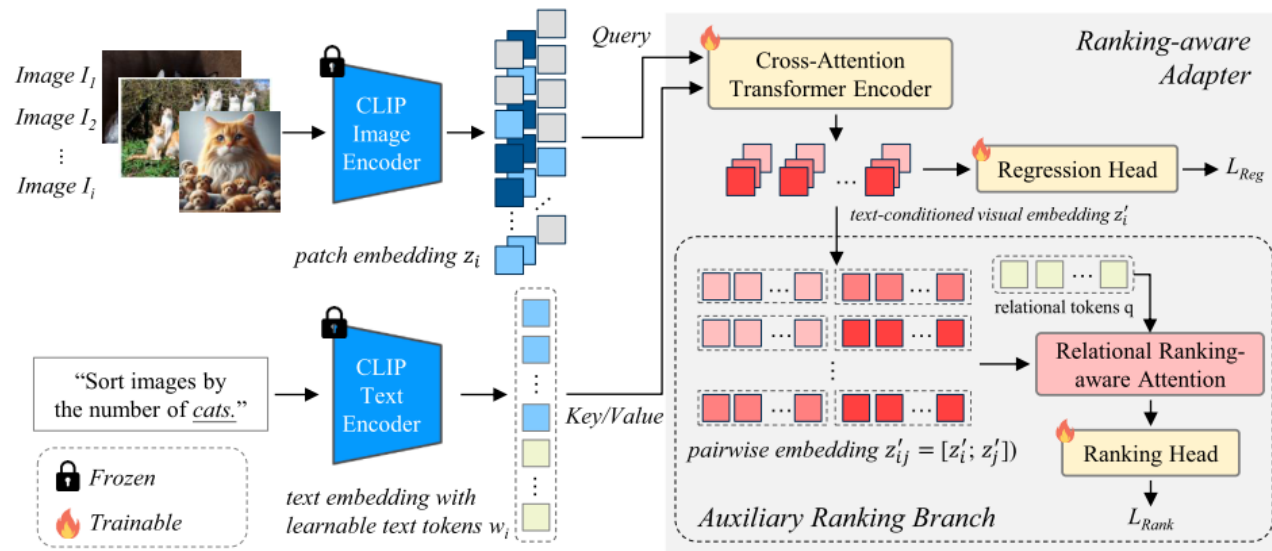
Baseline: CLIP w/ finetuning  
(contrastive Learning)





# Ablation Study

LTR paradigm	Ranking Head	Ranking-aware Attention	HCI		Object count sorting			
			MAE ( $\downarrow$ )		PLCC ( $\uparrow$ )		SRCC ( $\uparrow$ )	
-	-	-	1.113		0.251		0.422	
✓	-	-	0.402		0.612		0.538	
✓	✓	-	0.355	(+11.69%)	0.619	(+1.14%)	0.536	(-0.37%)
✓	✓	✓	0.317	(+21.14%)	0.624	(+1.96%)	0.557	(+3.53%)



# Conclusion

- We present an efficient and scalable framework **for text-driven image ranking** by reframing CLIP's image-text contrastive learning into an **LTR task**.
- By leveraging a lightweight adapter with our **ranking-aware attention** module, it can effectively capture **text-driven visual differences between image pairs**.
- Our work, **an all-in-one and end-to-end method**, surpass CLIP baselines and achieve results comparable to **SOTA methods tailored for specific task**, highlights the potential of leveraging VLMs with visual distinctions for developing sense of number sense.

**Thank You  
For Your Attention!**

Paper



Code

