



# Token Merging for Training-Free Semantic Binding in Text-to-Image Synthesis

**Taihang Hu<sup>1</sup>, Linxuan Li<sup>1</sup>, Joost van de Weijer<sup>3</sup>, Hongcheng Gao<sup>4</sup>  
Fahad Shahbaz Khan<sup>5,6</sup>, Jian Yang<sup>1</sup>, Ming-Ming Cheng<sup>1,2</sup>, Kai Wang<sup>3\*</sup>, Yaxing Wang<sup>1,2\*</sup>**

<sup>1</sup>VCIP, College of Computer Science, Nankai University, <sup>2</sup>NKIARI, Shenzhen Futian

<sup>3</sup>Computer Vision Center, Universitat Autònoma de Barcelona

<sup>4</sup>University of Chinese Academy of Sciences

<sup>5</sup>Mohamed bin Zayed University of AI, <sup>6</sup>Linkoping University

# Introduction: Problem Definition

- Semantic Binding: associating an object with its attribute (attribute binding) or linking it to related sub-objects (object binding).
- Existing Issues: Incorrect binding and missing attributes.

a dog wearing hat and a cat wearing sunglasses

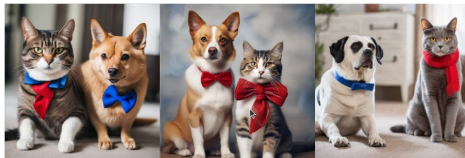


SDXL

SD3

*ToMe* (Ours)

a dog with blue bow tie and a cat with red scarf



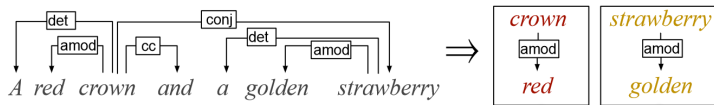
SDXL

PlayGround v2

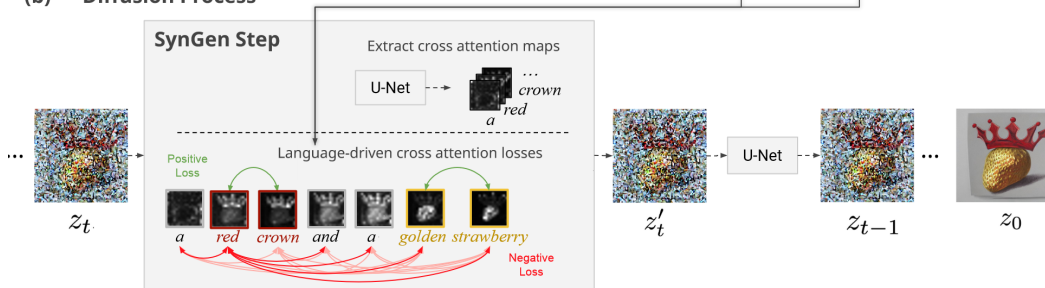
*ToMe* (Ours)

# Introduction: Related works

## (a) Entity-Modifier Identification

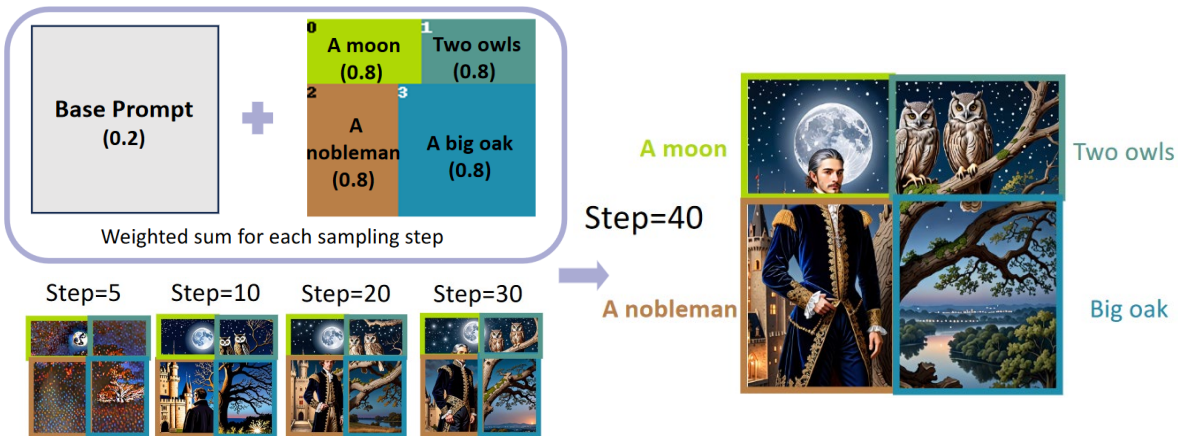


## (b) Diffusion Process



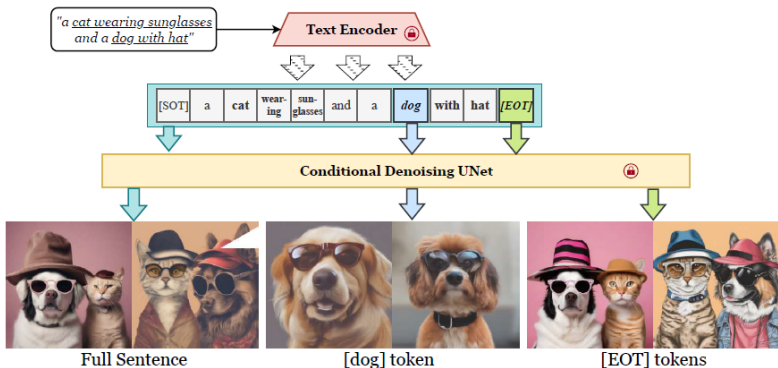
[1] Linguistic binding in diffusion models: Enhancing attribute correspondence through attention map, NeurIPS 2023 oral

# Introduction: Related works

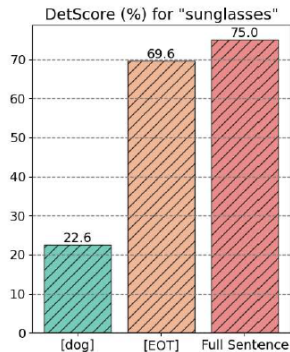


# Motivation

- Text tokens exhibit **information coupling**, even a single token can couple with preceding information.
- EOT(End of Text) has the ability to **contain all information**.



(a) T2I generation with various tokens



(b) DetScore probabilities

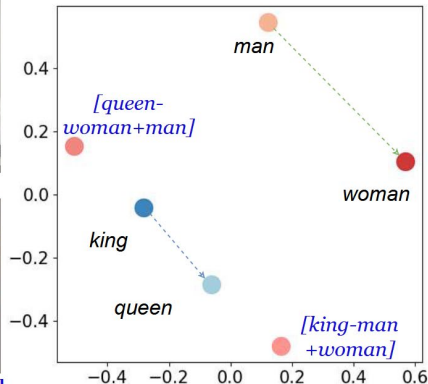
# Motivation

- Textual token embedding is **additive**, and the composite token obtained through element-wise addition has the ability to represent multiple objects.

(a) Additivity Text-to-Image Generations

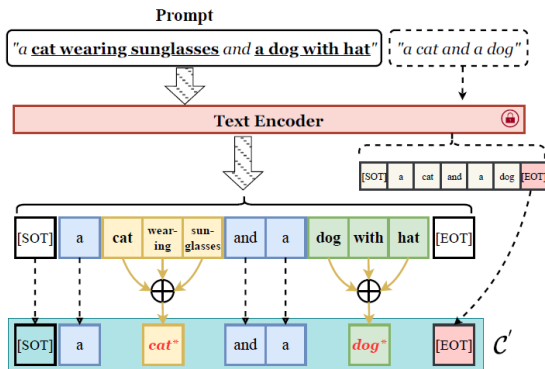


(b) PCA plot of text embeddings

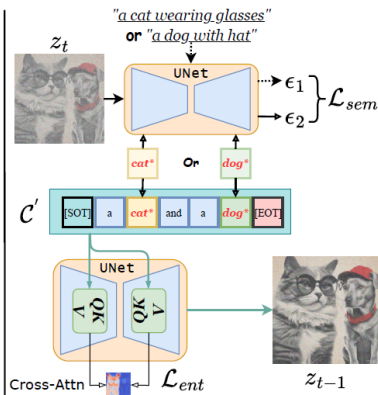


# Method

- Token merging: aggregating relevant tokens into a single composite token, aligning the object, its attributes, and sub-objects in the **same cross-attention map**



(a) Token Merging and End Token Substitution



(b) Iterative Composite Token Update

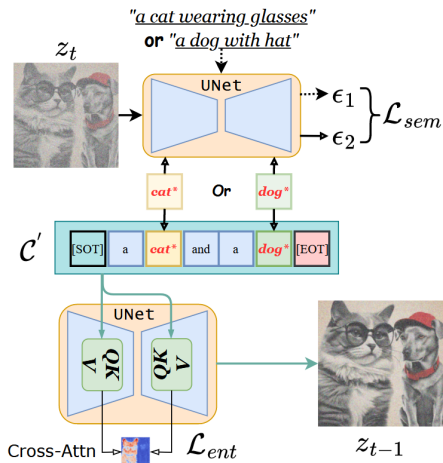
# Method

- Semantic binding loss:** Using a clean prompt as supervisory to eliminate any irrelevant semantic information within the composite token

$$\mathcal{L}_{sem} = \sum_{k \in [1, K]} \|\epsilon_{\theta}(z_t, \hat{c}_k, t) - \epsilon_{\theta}(z_t, \mathcal{C}, t)\|_2^2$$

- Entropy loss:** Ensure that tokens focus exclusively on their designated regions, preventing the cross-attention map from becoming overly divergent

$$\mathcal{L}_{ent} = \sum_{k \in [1, K]} \sum_{p_i \in A_k} -p_i \log(p_i)$$





# Experiments

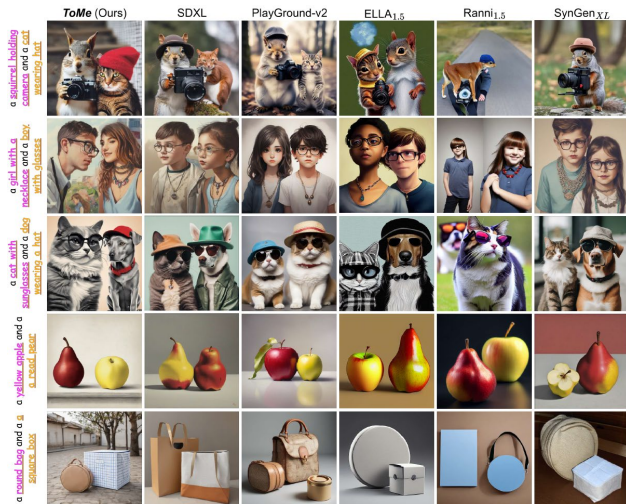


Figure 5: Qualitative comparison among various T2I generation methods with complex prompts.

# Experiments

Table 1: Quantitative results for semantic binding assessment on various benchmarking subsets. We denote the best score in **blue**, and the second-best score in **green**.

Method	Base Model	Train	BLIP-VQA $\uparrow$			Human-preference $\uparrow$			GPT-4o $\uparrow$
			Color	Texture	Shape	Color	Texture	Shape	
SDXL[53]	-	✓	0.6369	0.5637	0.5408	0.7798	0.5140	0.4029	0.4907
PlayG-v2[37]	-	✓	0.6208	0.6125	0.5087	-	-	-	0.5417
Ranni[21]	SD1.5	✓	0.2414	0.3029	0.2857	-0.8554	-0.6853	-0.8051	0.4166
ELLA[30]		✓	0.6911	0.6308	0.4938	0.6586	0.2963	0.0565	0.6481
SynGen[58]		✗	0.6619	0.6451	0.4661	0.4326	0.5072	0.0426	0.5545
CoMat[34]		✓	0.6561	0.6190	0.4975	-	-	-	-
Ranni[21]	SDXL	✓	0.6893	0.6325	0.4934	-	-	-	-
ELLA[30]		✓	0.7260	0.6686	0.5634	-	-	-	-
SynGen[58]		✗	0.7010	0.6044	0.5069	1.016	0.7867	0.4016	0.6458
CoMat[34]		✓	0.7774	0.6591	0.5262	-	-	-	-
<i>ToMe</i> (Ours)	SDXL	✗	0.7656	0.6894	0.6051	1.074	0.9281	0.5916	0.9549

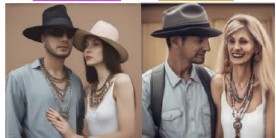
# Experiments



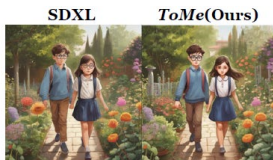
A dog with hat and a cat



A cat with scarf and a dog with tie



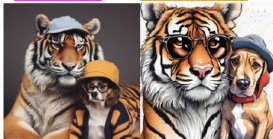
A man with hat and a girl with necklace



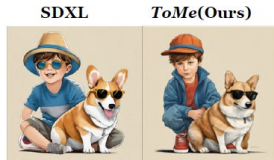
A boy with glasses and a girl



A bear with hat and a man with glasses



A tiger with glasses and a dog with hat



A boy with hat and a Corgi with sunglasses



A fox with sunglasses and a deer with crown



A squirrel holding guns and a bear with hat

# Experiments: Ablation Study

- Semantic binding loss filters out irrelevant information in tokens;
- Entropy loss makes the cross-attn map more focused.

Table 2: Ablation Study conducted on the T2I-CompBench benchmark.

Conf.	$ToMe$	$\mathcal{L}_{ent}$	$\mathcal{L}_{sem}$	BLIP-VQA		
				Color	Texture	Shape
A	×	×	×	0.6369	0.5637	0.5408
B	✓	×	×	0.6577	0.5828	0.5437
C	✓	✓	×	0.7525	0.6775	0.5797
D	×	✓	✓	0.5881	0.6194	0.5386
E	×	✓	×	0.5983	0.5798	0.5125
F	✓	×	✓	0.6804	0.6263	0.5645
<b>Ours</b>	✓	✓	✓	<b>0.7656</b>	<b>0.6894</b>	<b>0.6051</b>



Figure 6: Text-to-Image generation with various configurations.

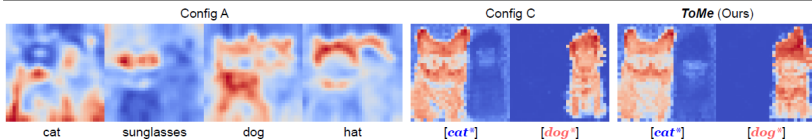


Figure 7: Cross-Attention maps visualization with various configurations.

# Experiments: Additional Applications

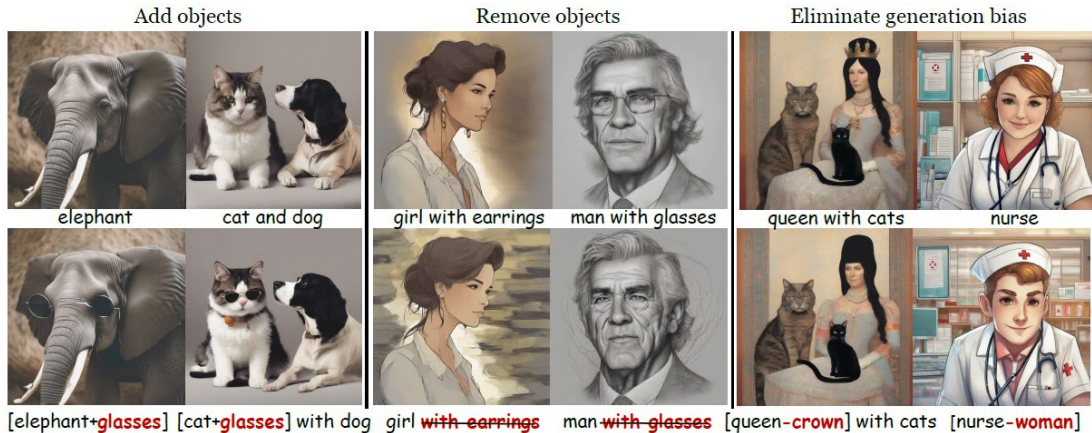


Figure 8: Additional applications of semantic additivity in text embedding.

# Conclusion

- The paper introduces Token Merging (ToMe), a training-free method that enhances semantic binding without the need for additional datasets, large language models, or extensive fine-tuning.
- By merging tokens for objects and their attributes into a **single composite token**, ToMe ensures that the generated image maintains **coherent cross-attention**, aligning the visual output closely with the intended semantics of the text prompt.

Code: <https://github.com/hutaihang/Tome>