

# VerbDiff

## Text-only Diffusion Models with Enhanced Interaction Awareness

25' CVPR

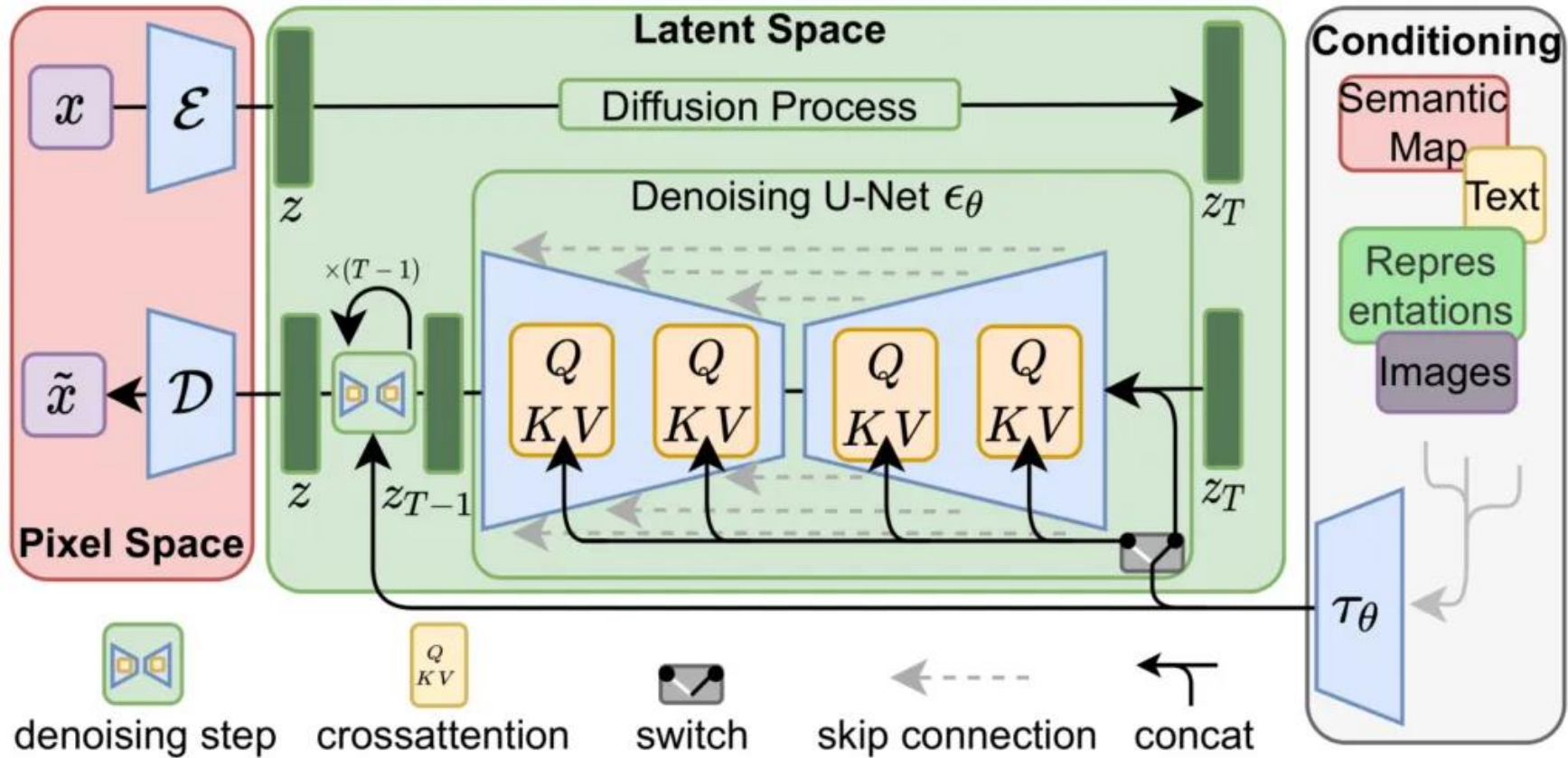
Seung Ju Cha, Kwanyoung Lee, Ye-Chan Kim, Hyunwoo Oh, Dong-Jin Kim

# Contents



1. Backgrounds
2. Related Works
3. Research Question
4. Method
5. Experiments
6. Conclusion

# Backgrounds



LDM<sup>[1]</sup>

Stable Diffusion(SD)<sup>[1]</sup> takes natural prompts as input and can generate photo-realistic images

# However...

SD sometimes struggles to generate accurate **interactions** between humans and **objects**<sup>[2]</sup>



Exiting, Train



Walking, Bicycle



Holding, Backpack



Pouring, Bottle

**The generated images do not match with our perceptual interactions**

**SD cannot reflect the intended interactions** as these are detailed semantics when depicting images<sup>[2]</sup>

# Why?

CLIP has a **strong bias towards objects or backgrounds**<sup>[3]</sup>



a girl **skateboarding** in a public place

a girl **dancing** in a public place

a girl **running** in a public place

a girl **singing** in a public place

a girl **sitting on her skateboard** in a public place

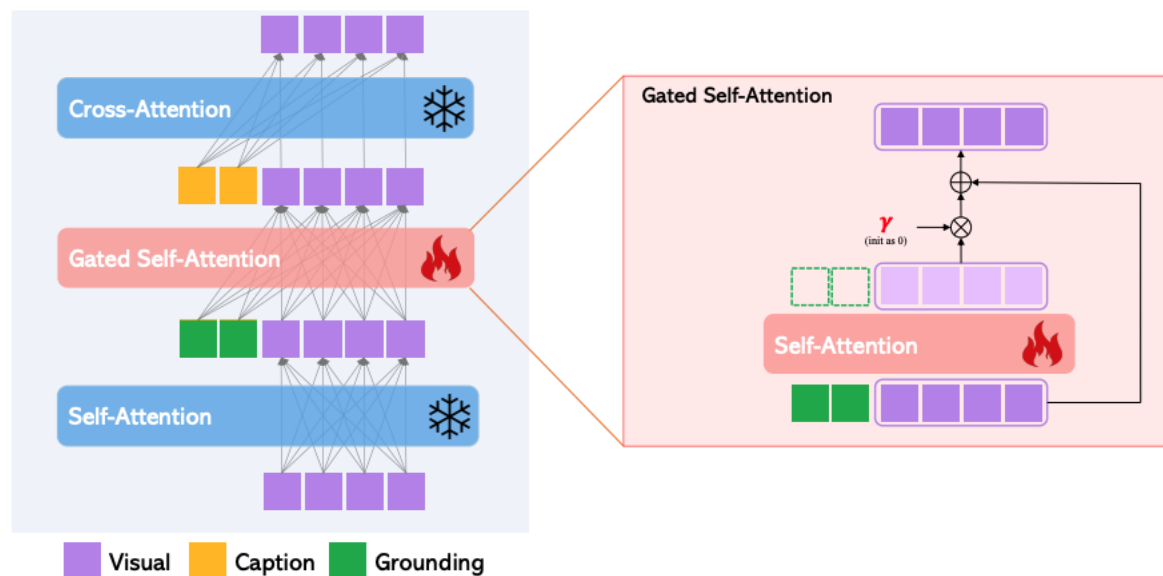
a girl **falling off her skateboard** in a public place

CLIP has a **lack of verb understandability** → affects the diffusion model

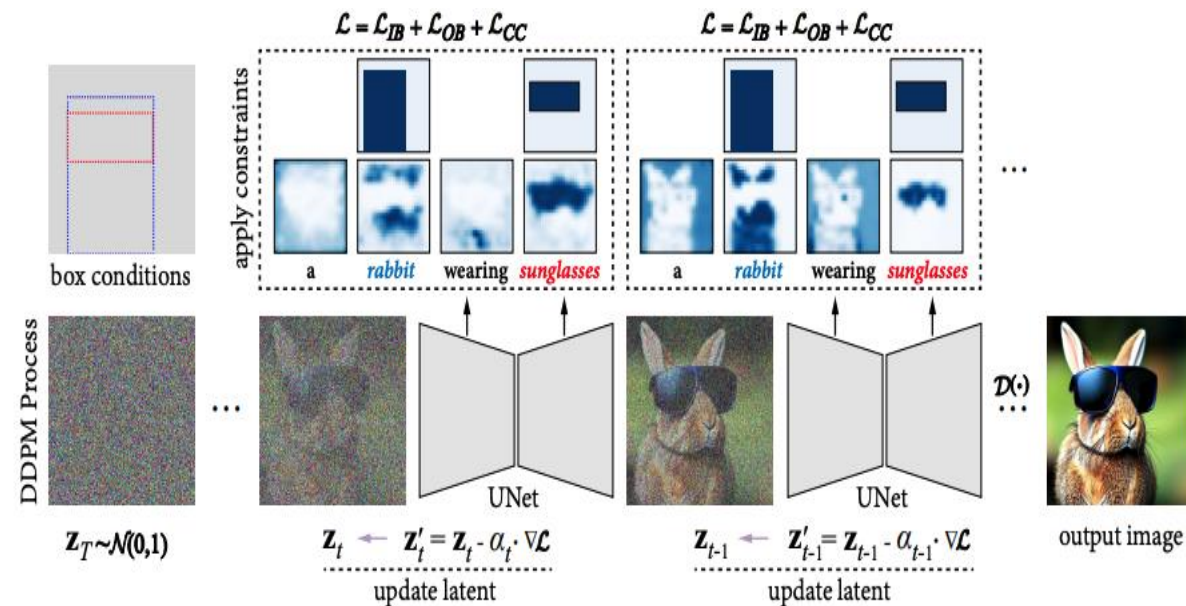
# Related Works

To enhance the text understandability

## 1. Layout (bounding box) based method



GLIGEN<sup>[4]</sup>



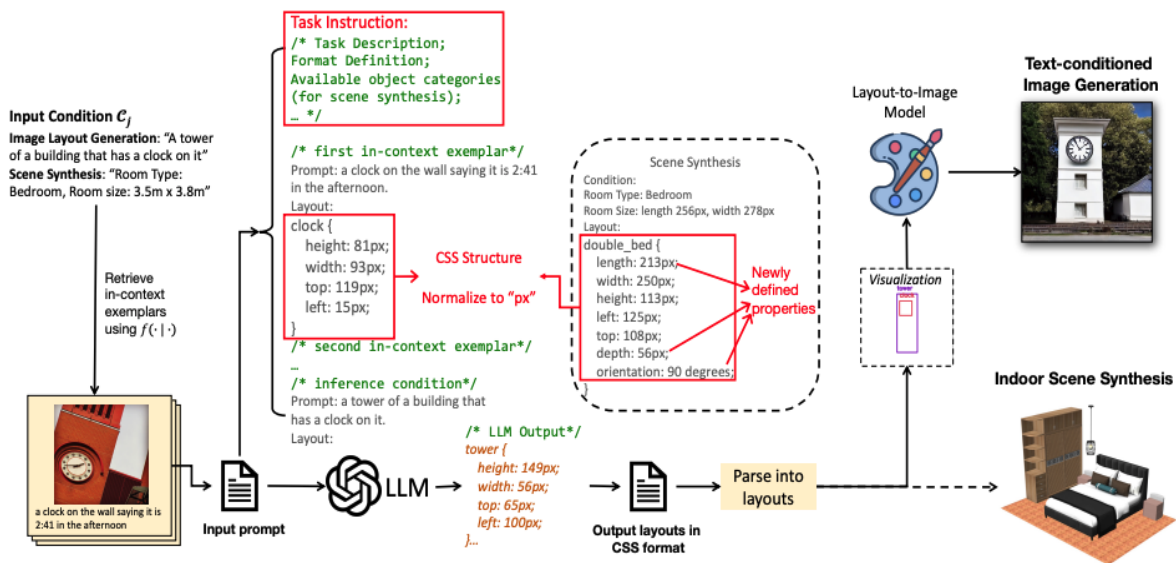
BoxDiff<sup>[5]</sup>

[4] Li, Yuheng, et al. "Gligen: Open-set grounded text-to-image generation." CVPR (2023)

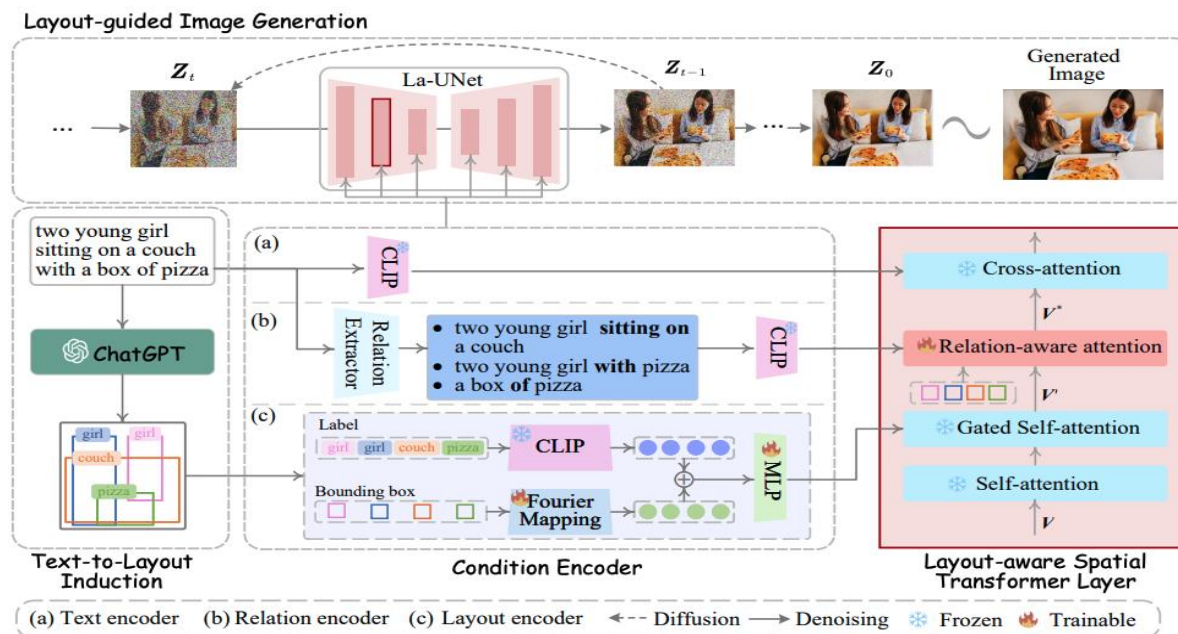
[5] Xie, Jinheng, et al. "Boxdiff: Text-to-image synthesis with training-free box-constrained diffusion." ICCV (2023)

# Related Works

## 2. LLM based method



LayoutGPT<sup>[6]</sup>



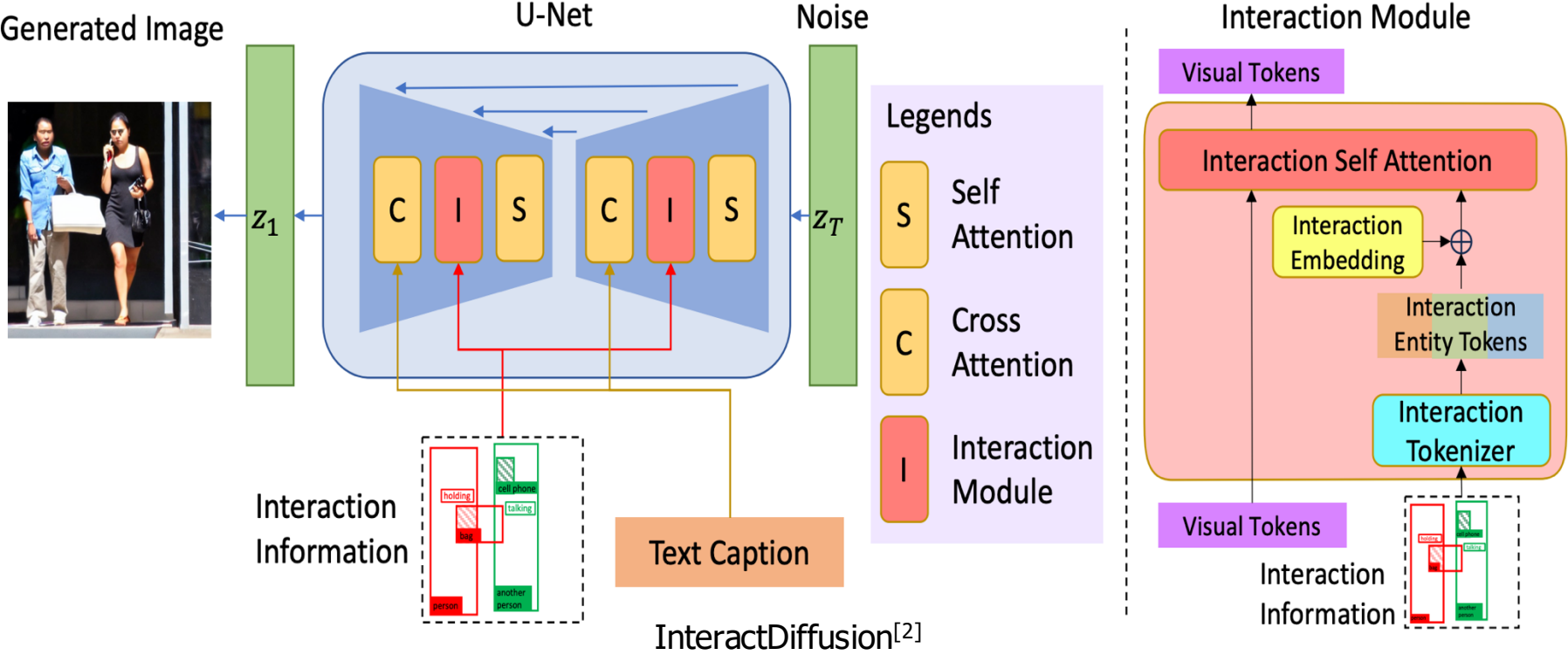
LayoutLLM-T2I<sup>[7]</sup>

These methods **focus on the spatial location of object relation not the interaction**

[6] Feng, Weixi, et al. "Layoutgpt: Compositional visual planning and generation with large language models." NeruIPS (2024)

[7] Qu, Leigang, et al. "Layoutllm-t2i: Eliciting layout guidance from llm for text-to-image generation." ACM MM (2023)

# Related Works

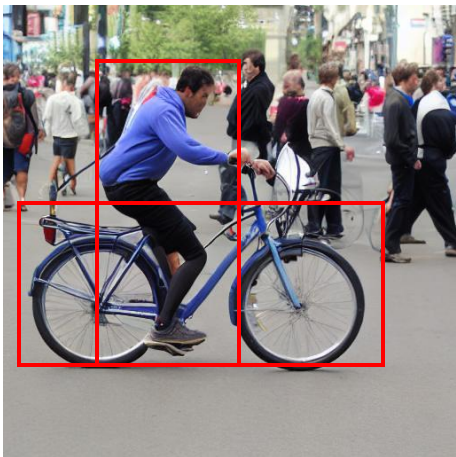


Leverage Bounding Box corresponding to Human, Object

Training Interaction Self-Attention Layer

# Research Question

Walking, Bicycle



Carrying, Kite



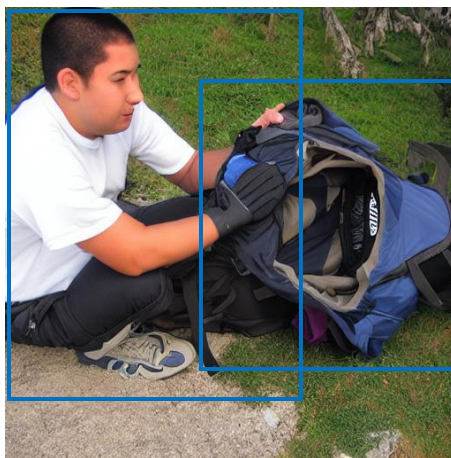
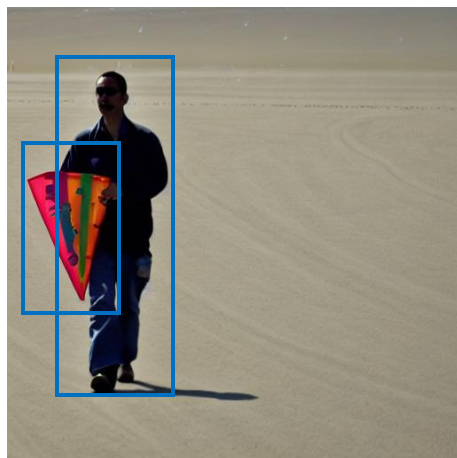
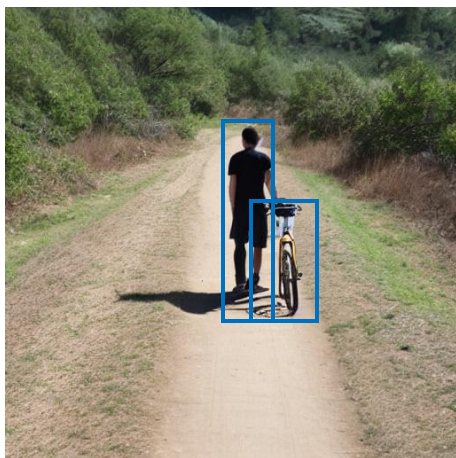
Holding, Backpack



When given bounding boxes that **multiple interaction shares**  
e.g., walking/riding bicycle



Cannot distinguish the interaction word semantic difference



**Heavily relies on the accurate bounding box**



considering the **precise boxes for interaction are labor-intensive**

Success and failure cases of InteractDiffusion<sup>[2]</sup>

Problem 1. Still has problem in **distinguishing the semantic differences between interaction words (verbs)**

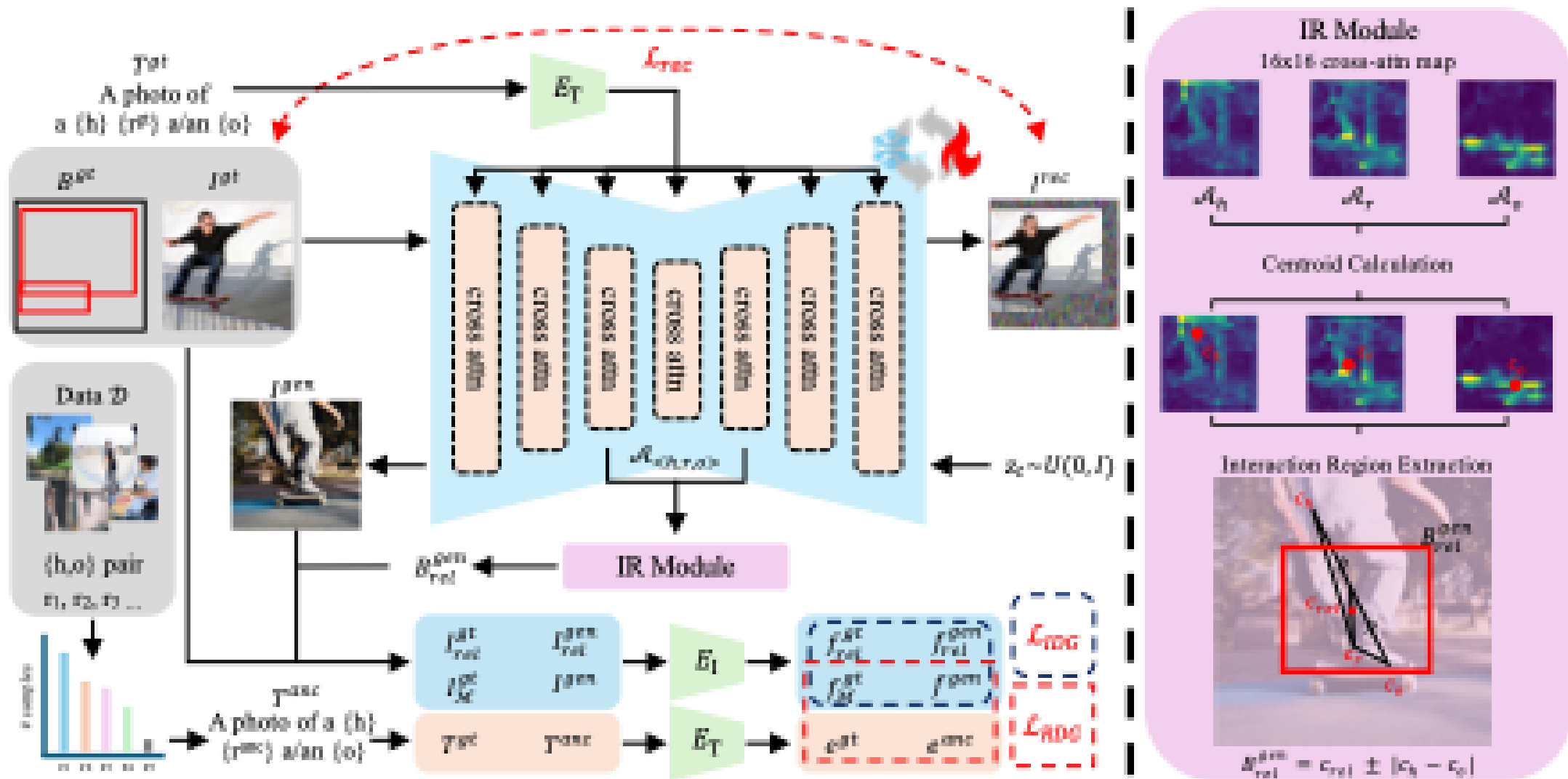
Problem 2. The generated images **heavily rely on precise bounding boxes which is labor-intensive** to provide

Can we enhance the **understandability of interaction words** in SD in a **more generalized way**?

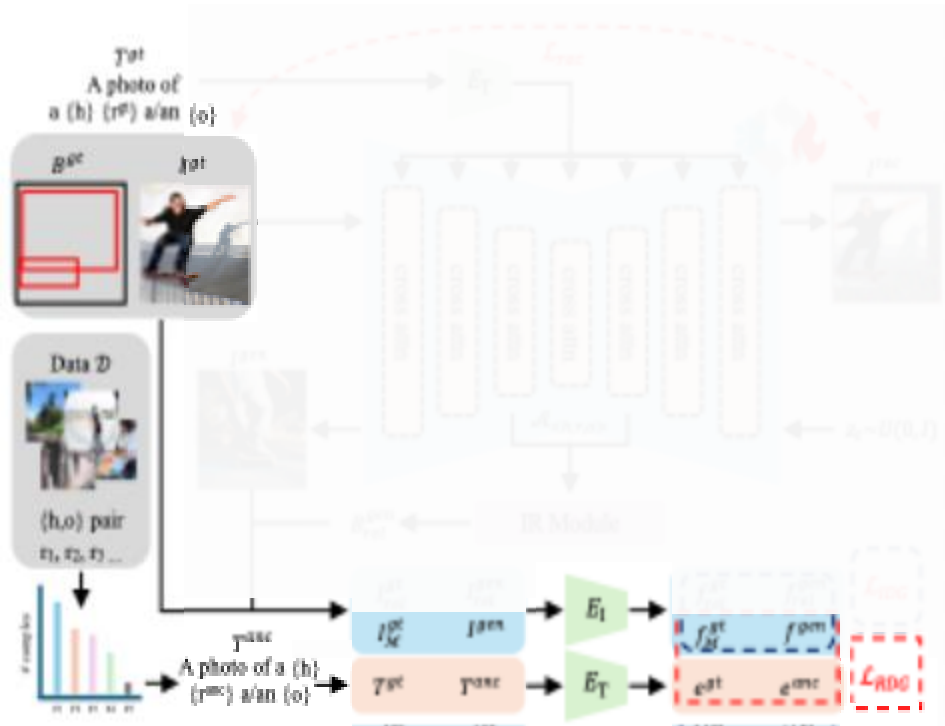
Propose VerbDiff with two additional guidances

1. **Relation Disentanglement Guidance** based on **frequency-based anchor text**
2. **Interaction Direction Guidance** with **IR modules** that capture detailed interaction regions

# Methods



# Relation Disentanglement Guidance



$\mathcal{L}_{triplet}$  distinguish the semantic difference between interaction words  
disentangles the input words from the frequency-based anchor word

$$\mathcal{L}_{triplet} = \max(0, m + \text{sim}(f^{gen}, e^{gt}) - \text{sim}(f^{gen}, e^{anc}))$$

$\mathcal{L}_{align}$  aims to control the accurate interaction in image level

$$\mathcal{L}_{align} = 1 - \frac{f_{\mathcal{M}}^{gt} \cdot f^{gen}}{|f_{\mathcal{M}}^{gt}| |f^{gen}|}$$

# Frequency-based anchor interaction word

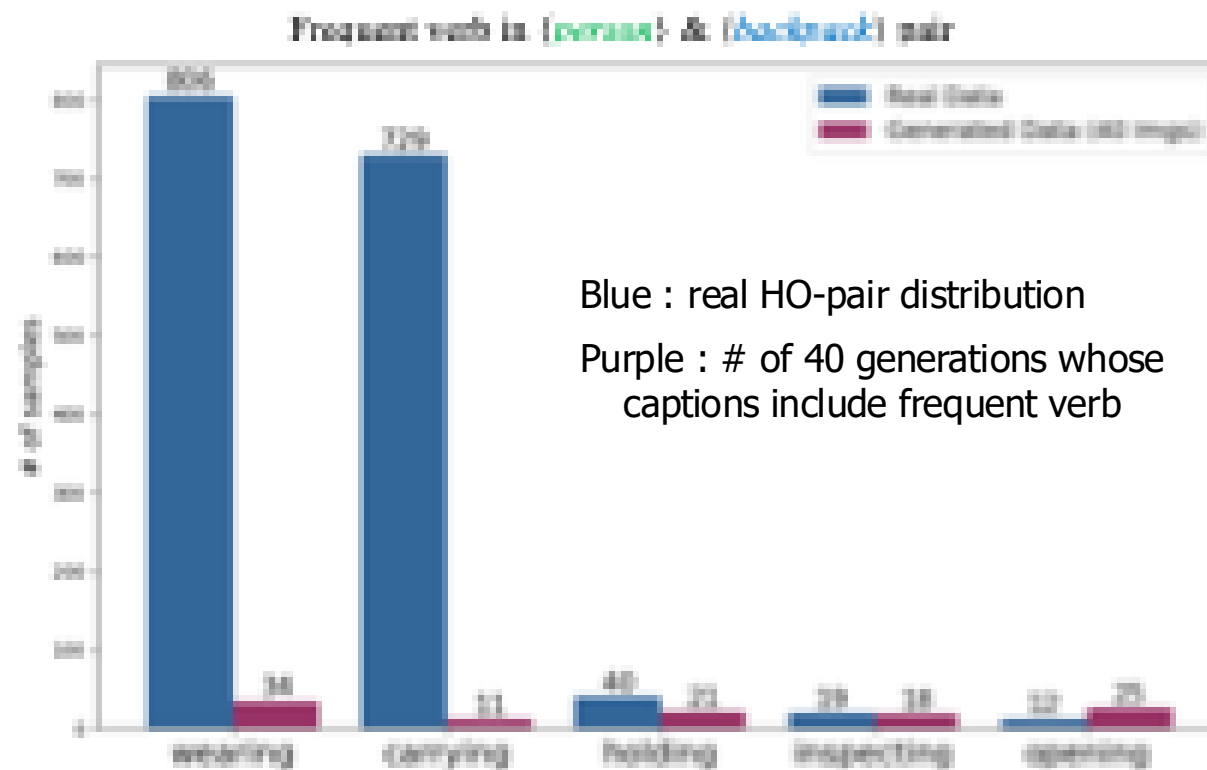
"a man **holding** a backpack"

SD



Captioner  
(InstructBLIP)

"man **wearing** a backpack on his back"



**SD tends to generate the images that matches with the most frequent verb!**

Define the anchor interaction word as the most frequent verb along each h,o pair

$$r^{anc} = \arg \max_{r \in R_o} \mathcal{C}(r|o)$$

# Adaptive Interaction Modification Number

Each h,o pair has different kinds of verbs and different number of images

The more the sample, the more the model effected

e.g., carrying backpack >> opening backpack

To **balance the modification extent**

$$\alpha(k) = \frac{1 - \beta^{n_k}}{1 - \beta}$$

$$\beta = \frac{N-1}{N}$$

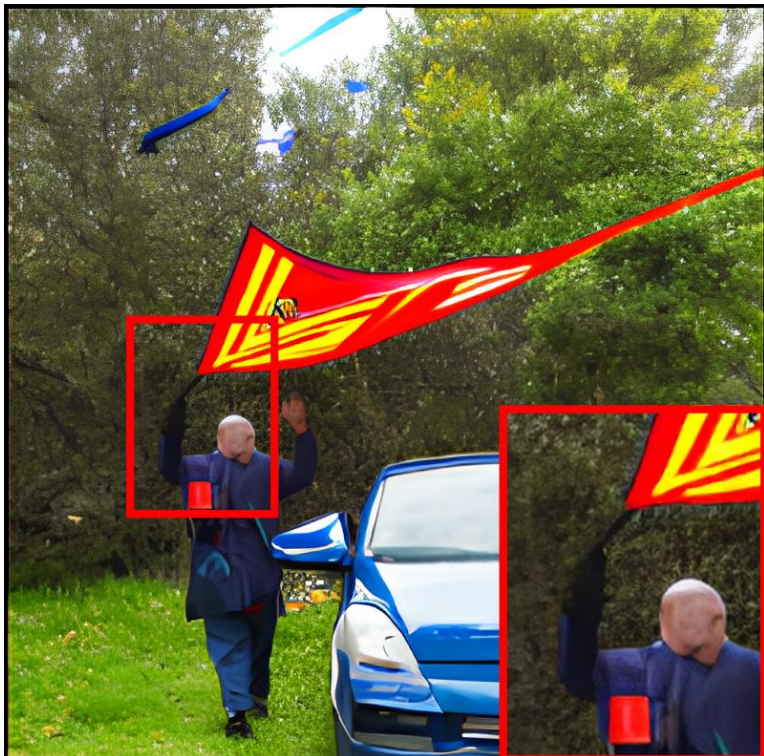
$n_k$  number of samples in class  $k$

$N$ : total number of samples in dataset

$$\mathcal{L}_{RDG} = \alpha \cdot (\mathcal{L}_{tiple} + \mathcal{L}_{align})$$

$\alpha$  : adaptive interaction modification number

# Interaction Direction Guidance



w/ Relation Disentanglement Guidance (RDG)

The images **misalign with human expectations**, especially in the **local region**

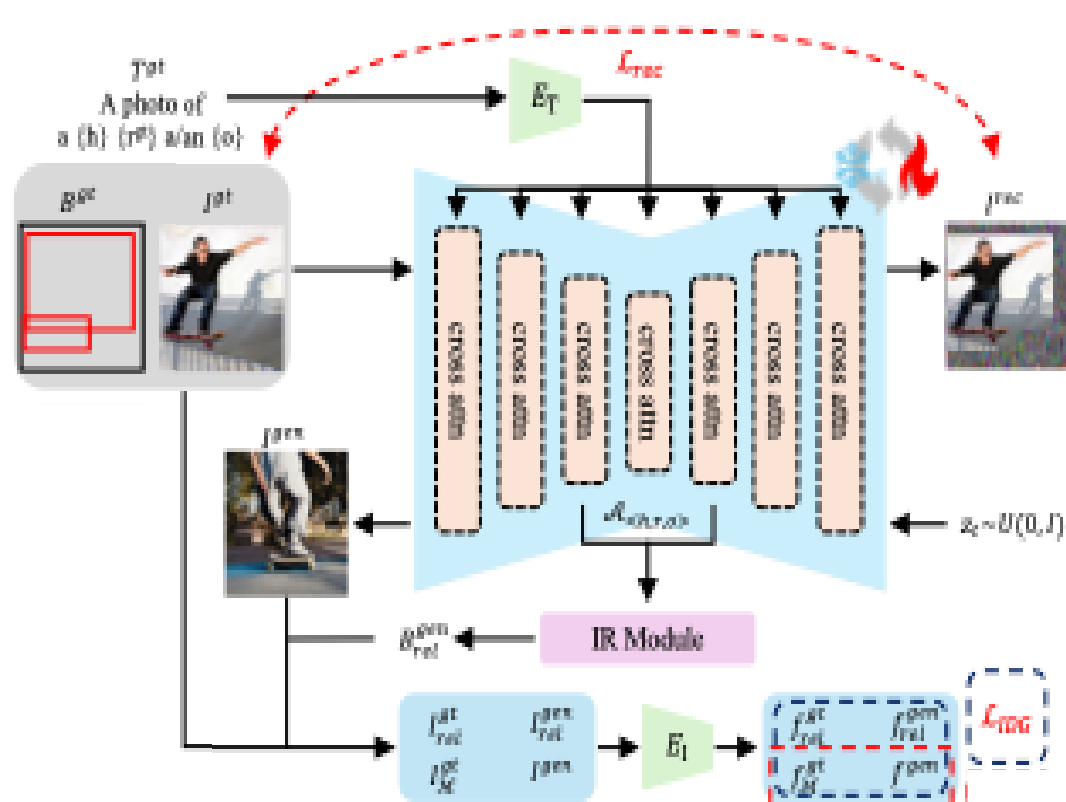
RDG focus on aligning global interaction feature

**To focus on the finer, localized interaction region**

Apply Interaction Direction Guidance(IDG)

Leveraging the interaction region extracted from Interaction Region module

# Interaction Direction Guidance



Modify the CLIP direction loss into the interaction region loss

$$L_{IDG} = 1 - \frac{(f_M^{gt} - f^{gen}) \cdot (f_{rel}^{bias})}{|f_M^{gt} - f^{gen}| |f_{rel}^{bias}|}$$

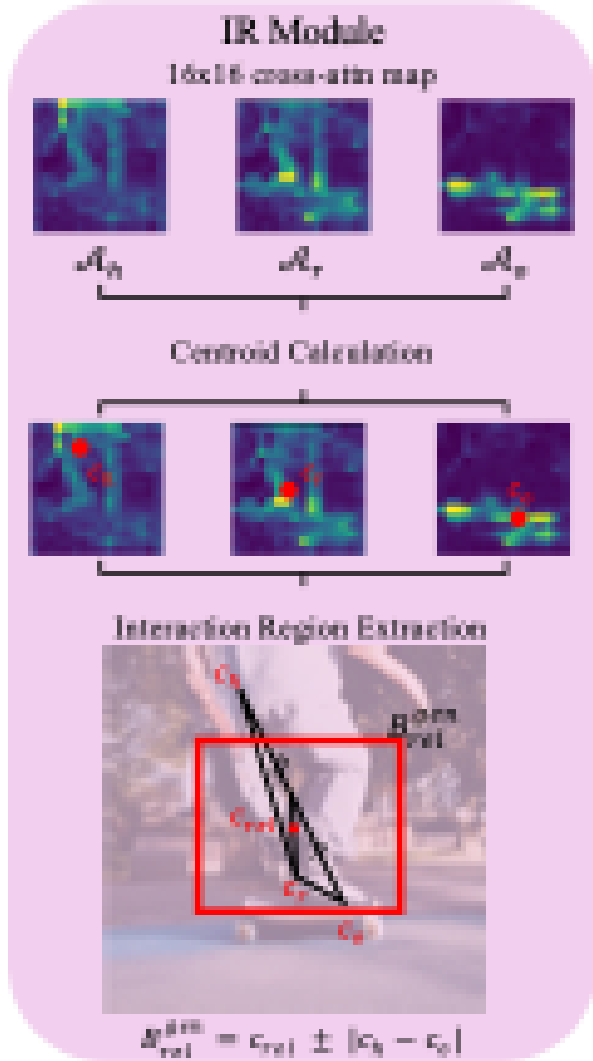
$$f_{rel}^{bias} = f_{rel}^{gt} - f_{rel}^{gen}$$

Can extract interaction region from ground-truth image

How to extract interaction regions from generated images?

Propose Interaction Region module (IR module)

# Interaction Region Module



Cross-attention maps in SD reflect the existence of each token in prompts<sup>[9]</sup>

Extract 16x16 resolution map corresponding to each **human**, **relation**, **object** token

Calculate the centroid  $c_{h,r,o}$  of each aggregated  $\mathcal{A}_{h,r,o}$

$$c = \frac{1}{\sum_{h,w} \mathcal{A}} \begin{bmatrix} \sum_{h,w} w \cdot \mathcal{A} \\ \sum_{h,w} h \cdot \mathcal{A} \end{bmatrix}$$

Define  $c_{rel}$  to be the centroid of the triangle defined by three centroids

Extract interaction region  $B_{rel}^{gen}$

$$B_{rel}^{gen} = c_{rel} \pm |c_h - c_o|_2$$

Training only the cross-attention layers in SD

$$\mathcal{L}_{rec} = \mathbb{E}_{z, \epsilon \sim \mathcal{N}(0,1), t} [||(\epsilon \odot \mathcal{M} - \epsilon_{\theta}(z_t, t, T) \odot \mathcal{M})||]$$

$$\mathcal{L}_{total} = \lambda_1 \cdot \mathcal{L}_{rec} + \lambda_2 \cdot \mathcal{L}_{RDG} + \lambda_3 \cdot \mathcal{L}_{IDG}$$

## Metric

1. CLIP & Sentence-BERT<sup>[10]</sup> similarity
2. HOI Classification Accuracy
3. VQA-Score<sup>[11]</sup>



























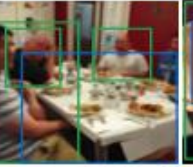






















1. Extract caption from generated image (InstructBLIP<sup>[12]</sup>) and calculate the embedding cosine similarity
2. HOI classification accuracy (K.O. : assume the object is correct, only compare the verb, Def. : both object and verb must be accurate)
3. The average of probability that the model answer yes to a given question

[10] Reimers, N. "Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks." *arXiv preprint arXiv:1908.10084* (2019).

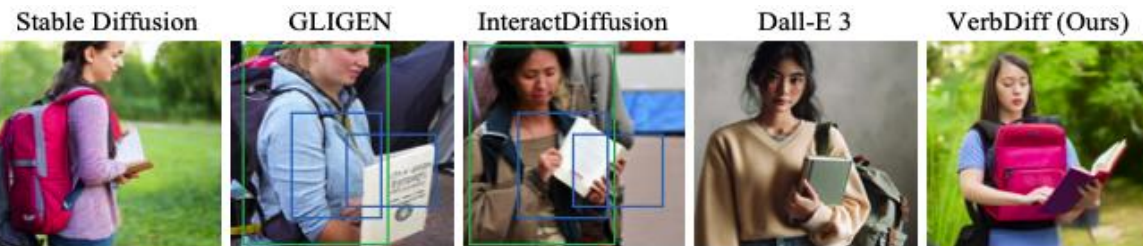
[11] Lin, Zhiqiu, et al. "Evaluating text-to-visual generation with image-to-text generation." *European Conference on Computer Vision*. Springer, Cham, 2025.

[12] Lin, Zhiqiu, et al. "Evaluating text-to-visual generation with image-to-text generation." *European Conference on Computer Vision*. Springer, Cham, 2025.

# Single Interaction

	(a) <i>Man, Exiting Train</i>	(b) <i>Man, Walking Bicycle</i>	(c) <i>Man, Holding Backpack</i>	(d) <i>Woman, Blowing Cake</i>	(e) <i>Woman, Drinking Bottle</i>	(f) <i>People, Eating at Dining table</i>	(g) <i>Man, Pouring Bottle</i>
$\{H, R, O\}$							
GT							
SD							
GLIGEN							
InteractDiffusion							
Dall-E 3							
VerbDiff (Ours)							
							

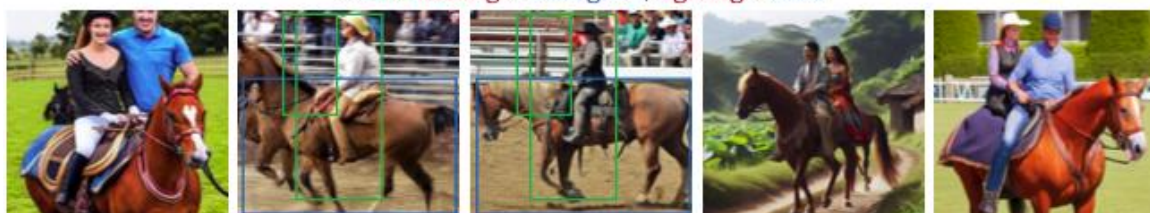
# Multi-Interactions



"A woman holding a backpack and a book"



"A man holding a wine glass, lighting a cake"



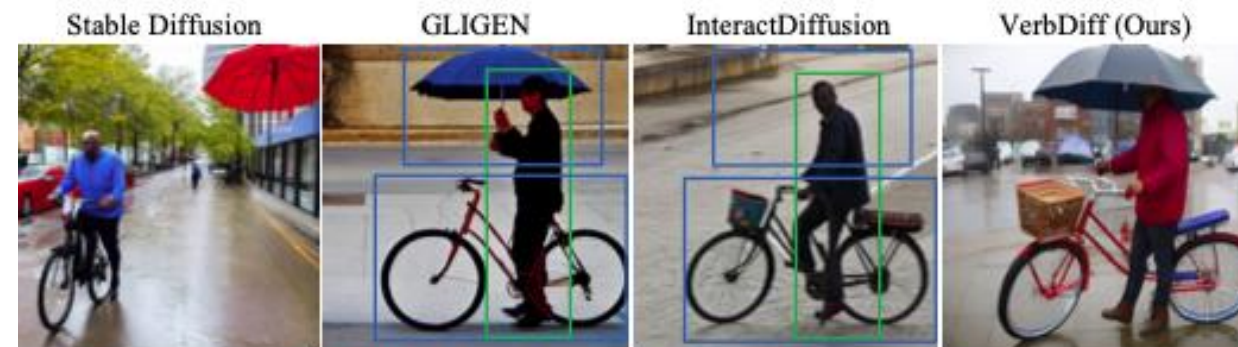
"A man and a woman riding a same horse"



"An old man and a young girl riding a tandem bicycle"



"A man lying on a bench, reading a book"



"A photo of a man walking a bicycle on the street, carrying an umbrella"



"A basketball player jumping and throwing a basketball and a man blocking the ball"




# Experiments

Models	CLIP		S-BERT	VQA-Score	
	T2T	T2I	T2T	I2T	VQA
SD [25]	<u>0.725</u>	<b>0.251</b>	<u>0.620</u>	<u>0.769</u>	<u>0.765</u>
GLIGEN [15]	0.683	0.238	0.554	0.679	0.674
InteractDiffusion [10]	0.703	0.232	0.575	0.728	0.734
VerbDiff (Ours)	<b>0.733</b>	<u>0.242</u>	<b>0.633</b>	<b>0.771</b>	<b>0.766</b>

Table 1. **Similarity comparison between VerbDiff and other models.** We evaluate scores on CLIP, S-BERT [24] and a large vision-language alignment benchmark VQA-Score [17].

Models	SOV-STG-S (Acc $\uparrow$ )				SOV-STG-Swin-L (Acc $\uparrow$ )			
	Def.		KO.		Def.		KO.	
	Full	Rare	Full	Rare	Full	Rare	Full	Rare
SD [25]	16.09	4.59	18.22	4.85	20.08	8.07	21.69	8.66
GLIGEN [15]	15.88	4.85	17.91	5.24	17.83	7.00	19.35	7.57
InteractDiffusion [10]	<u>19.67</u>	<u>7.00</u>	<u>21.31</u>	<u>7.69</u>	<u>23.53</u>	<u>10.27</u>	<u>24.86</u>	<u>11.18</u>
VerbDiff (Ours)	<b>22.59</b>	<b>7.62</b>	<b>24.79</b>	<b>7.83</b>	<b>27.05</b>	<b>12.60</b>	<b>28.43</b>	<b>13.18</b>

Table 2. **HOI accuracy comparison between VerbDiff and previous methods.** Def. and KO. refer to Default and Known Object.

(a)	(b)	(c)
		
walking : 0.6846	walking : 0.7490	walking : 0.3003
riding : 0.3947	riding : 0.6954	riding : 0.1658
repairing : 0.1821	repairing : 0.1711	repairing : <b>0.9915</b>
jumping : 0.1038	jumping : 0.1158	jumping : 0.1140




(d)	(e)	(f)
		
walking : 0.0564	walking : 0.0821	walking : 0.0808
riding : 0.8756	riding : <b>0.9897</b>	riding : <b>0.9852</b>
repairing : 0.0483	repairing : 0.0759	repairing : 0.0639
jumping : <b>0.9799</b>	jumping : 0.0808	jumping : 0.1194

Figure 3. **VQA-score result examples.** We measure the probability that the VQA model answers “yes” for questions based on four verbs associated with the “bicycle” class. The verb with the highest score is highlighted in red.

# Ablations

VerbDiff	$\mathcal{L}_{rec}$	$\mathcal{L}_{triple}$	$\mathcal{L}_{align}$	$\mathcal{L}_{IDG}$	CLIP T2T	S-BERT T2T	HOI Acc	
							Def.	KO.
(a)	✓				0.691	0.582	19.38	20.89
(b)	✓	✓	✓		0.700	0.589	20.32	21.87
(c)	✓			✓	0.699	0.588	20.21	21.67
(d)	✓	✓		✓	0.710	0.610	23.39	24.51
All (Ours)	✓	✓	✓	✓	<b>0.733</b>	<b>0.633</b>	<b>27.05</b>	<b>28.43</b>

Table 3. **Ablation of the VerbDiff guidance settings.** We score the similarity and HOI accuracy on the Full setting. We use the SOV-STG-Swin-L model. Combining all the proposed loss functions shows the best performance.

