



OneHOI: Unifying Human-Object Interaction Generation and Editing

Jiun Tian Hoe · Weipeng Hu · Xudong Jiang · Yap-Peng Tan · Chee Seng Chan

Motivation

Human-Object Interaction (HOI) = $\langle \text{person, action, object} \rangle$

Current Limitations in HOI content creation

HOI Editing

- Only text-driven edits
- Fails to decouple pose from physical contact
- Limited to single HOI
- Lack spatial control
- Limited HOI semantic understanding

HOI Generation

- Hard to accept arbitrary shape
- Cannot mix HOI with object-only entities
- Input conditions not flexible

Key Insight
Generation and editing are the same conditional denoising process
 Joint learning creates strong **synergy**; broad semantics from generation enable more plausible and physically-aware edits.

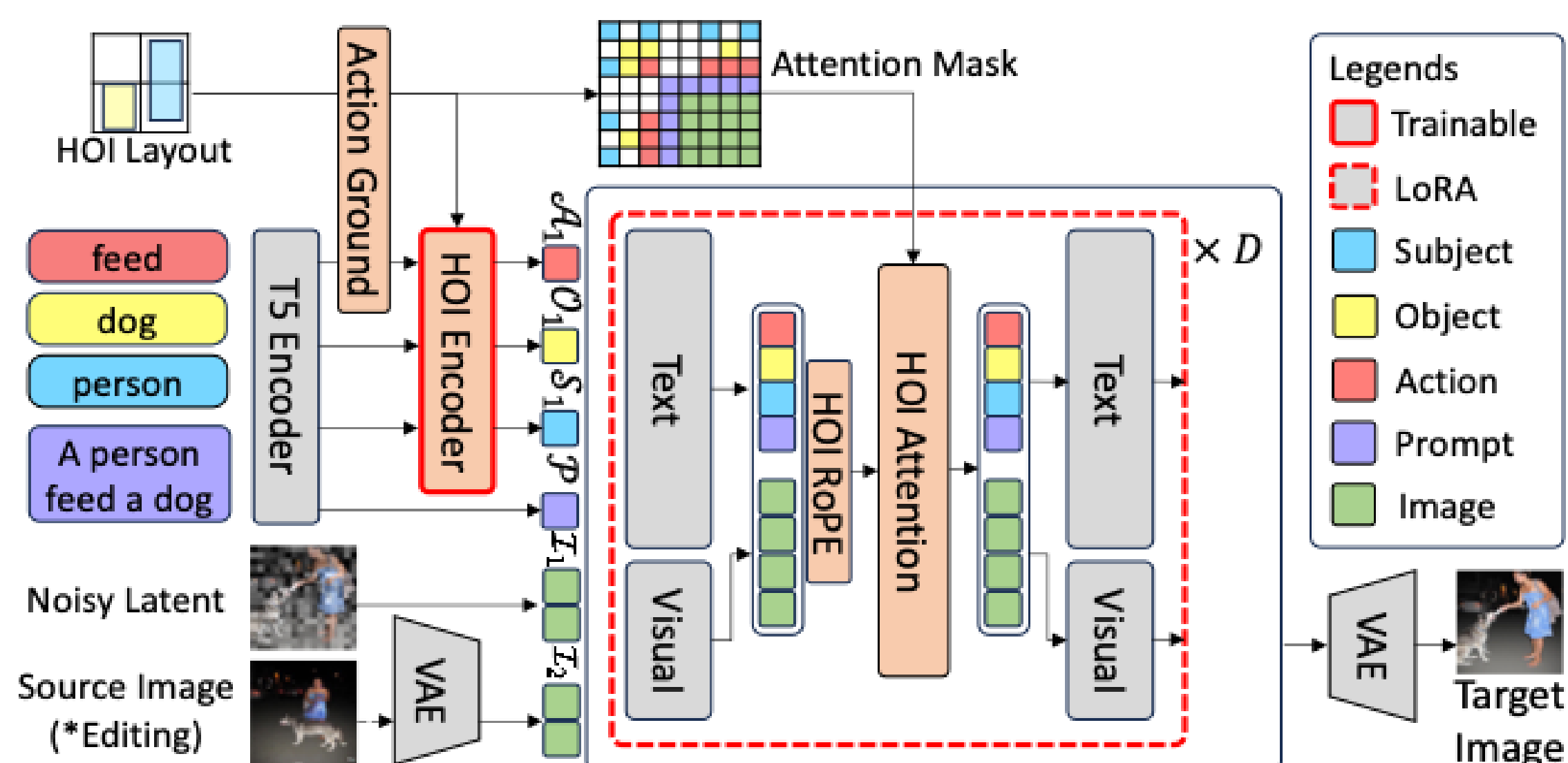
Our Contributions

- 👍 First unified framework for HOI generation & editing
- 👍 Relational Diffusion Transformer (R-DiT) with explicit interaction modeling
- 👍 Supports layout-guided, layout-free, arbitrary-mask, and mixed condition control
- 👍 Enable multi-HOI editing
- 👍 HOI-Edit-44K dataset, first for HOI editing pairs
- 👍 State-of-the-art result across all HOI synthesis task

TL;DR

Can HOI generation and editing be unified?
YES - and they actually make each other better.

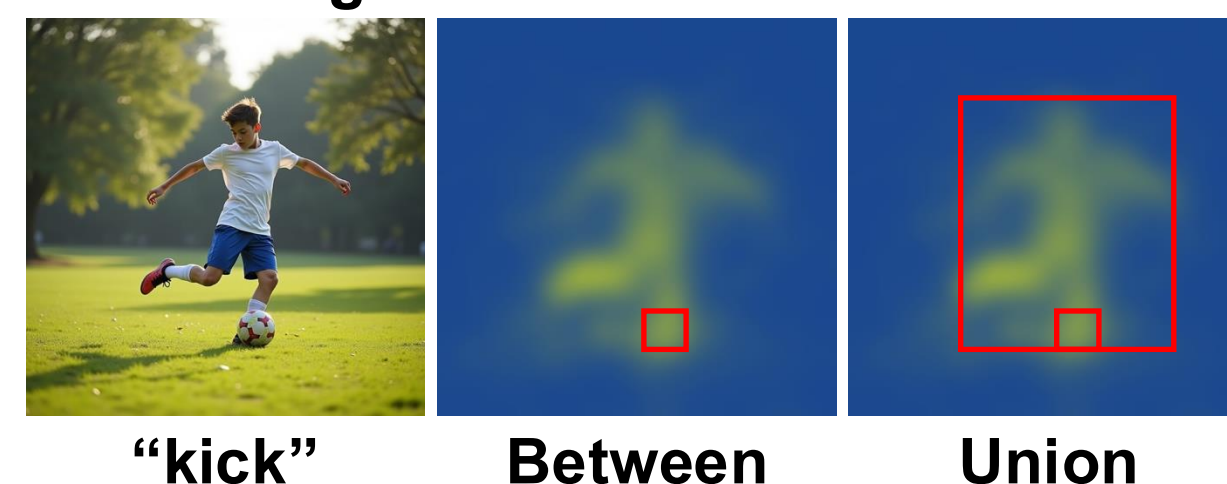
Overall Framework



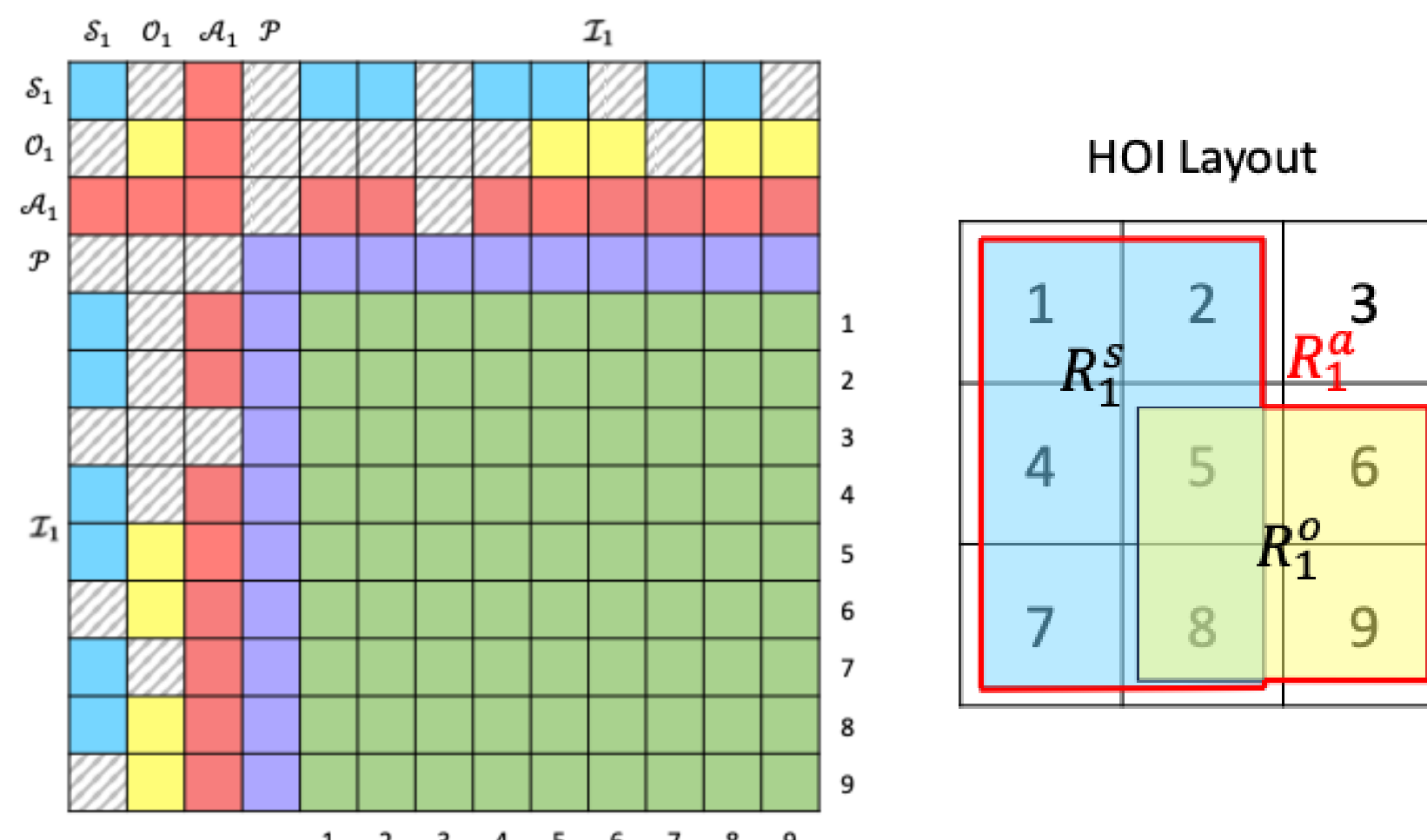
Train Strategy

- Joint Training Generative + Editing
- Modality Dropout Allow flexible condition combinations
- Multi-source Data HOI-Edit-44K and HOI generation data

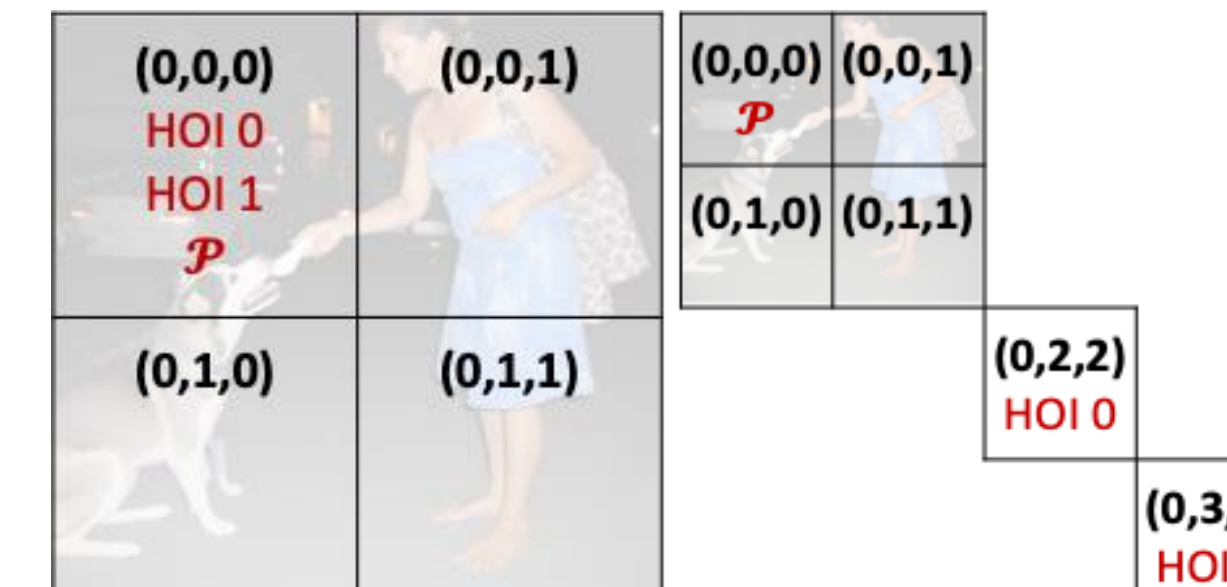
1 Action Region (Union)



2 Structured HOI Attention



3 HOI RoPE



1. Our "union" better aligns with action attention maps in DiT
2. Verb-mediated attention allows explicit interaction modelling
3. HOI RoPE assigns distinct positional slots to distinguish multiple HOIs

Results

The Results section displays four main categories of HOI synthesis: HOI Generation, HOI Editing (Layout-free), HOI Editing (Add), and Attribute Editing (Object Change). Each category includes a text prompt, a visual representation of the HOI, and the resulting generated image. Examples include an astronaut holding a flag, a person holding a sword, and a person walking a dog.

Method	HOI Editing		Image Quality		
	Editability-Identity	HOI Editability	PickScore	HPS	ImageReward
Null-Text Inversion	0.443	0.390	20.81	0.2483	-0.3329
InstructPix2Pix	0.380	0.269	20.28	0.2178	-0.7717
Flux.1 Kontext	0.471	0.328	20.45	0.2427	-0.5137
Qwen Image Edit	0.580	0.460	20.81	0.2585	0.0748
InteractEdit	0.573	0.514	<u>21.08</u>	<u>0.2640</u>	0.1630
Nano Banana (close)	0.623	0.530	20.97	0.2544	0.1810
Ours	0.638	0.596	21.26	0.2805	0.4713

Method	HOI Generation		Layout-guided HOI Editing				
	Controllability	Image Quality	Layout-guided HOI Editing		Image Quality		
	Spatial	HOI	Edit-Identity	HOI Edit	PickScore	HPS	ImageReward
GLIGEN	0.5150	0.3344	0.559	0.520	0.749	0.2418	-0.3072
InstanceDiffusion	0.5228	0.3476	0.559	0.520	0.749	0.2418	-0.3072
MIGC++	0.5331	0.3616	0.559	0.520	0.749	0.2418	-0.3072
Eligen	0.4371	0.3061	0.559	0.520	0.749	0.2418	-0.3072
InteractDiffusion	0.5768	0.4505	0.559	0.520	0.749	0.2418	-0.3072
Ours	0.6104	0.4528	0.638	0.570	0.822	0.2678	0.2897