

# Scene Graph Disentanglement and Composition for Generalizable Complex Image Generation

Yunnan Wang<sup>1,2</sup>   Ziqiang Li<sup>1,2</sup>   Wenyao Zhang<sup>1,2</sup>   Zequn Zhang<sup>2,3</sup>  
Baao Xie<sup>2</sup>   Xihui Liu<sup>4</sup>   Wenjun Zeng<sup>2</sup>   Xin Jin<sup>2</sup>

<sup>1</sup>MoE Key Laboratory of Artificial Intelligence, Shanghai Jiao Tong University

<sup>2</sup>Ningbo Institute of Digital Twin, Eastern Institute of Technology

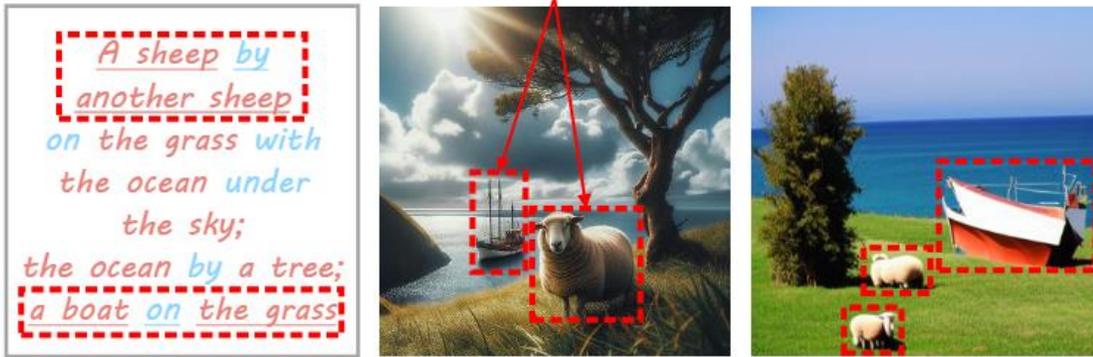
<sup>3</sup>School of Information Science and Technology, University of Science and Technology of China

<sup>4</sup>Department of Electrical and Electronic Engineering, The University of Hong Kong



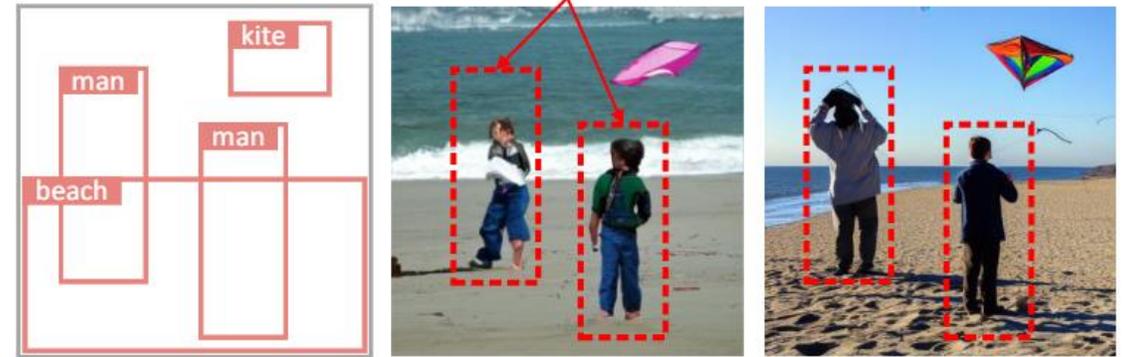
# Introduction

## Relationship and Quantity Confusion of T2I



(a) DALL-E 3 (middle) and Ours (right).

## Non-Spatial Interaction Dilemma of L2I



(b) LayoutDiffusion (middle) and Ours (right).

## Independent Nodes Missing of Semantics-based SG2I



(c) R3CD (middle) and Ours (right).

## 1. Challenges of Text-to-Image (T2I)

- Relationships and quantities of objects

## 2. Challenges of Layout-to-Image (L2I)

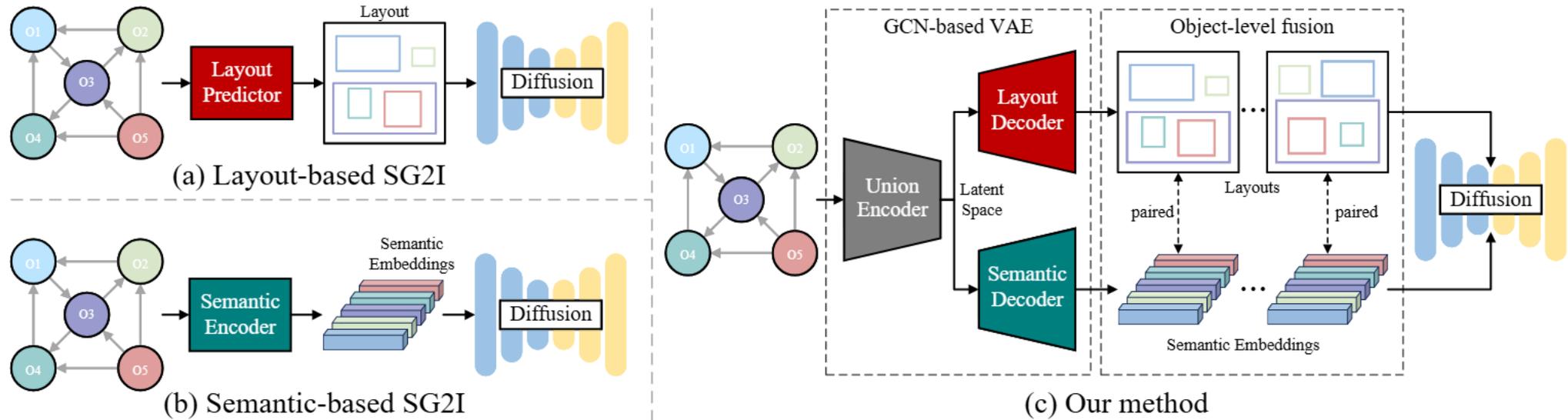
- non-spatial interactions between objects

## 3. Challenges of Scene-Graph-to-Image (SG2I)

- absence of independent



# Introduction

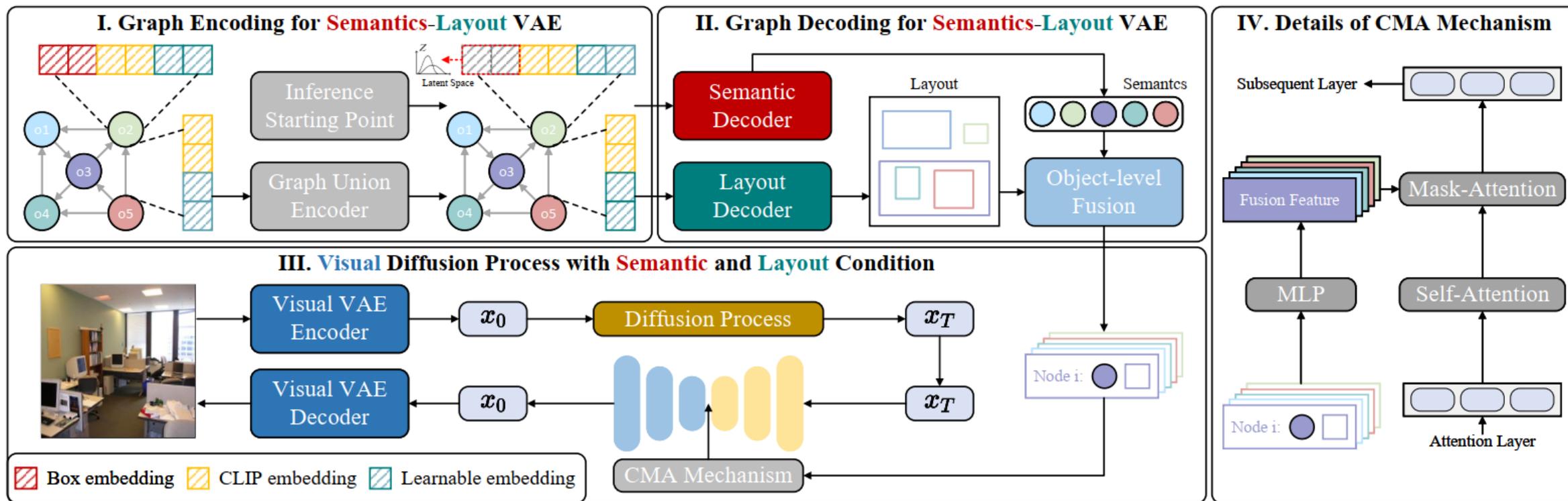


## Key contributions of our DisCo:

- *Semantics-Layout Variational AutoEncoder* (SL-VAE) that disentangles diverse spatial layouts and interactive semantics from the scene graph;
- *Compositional Masked Attention* (CMA) that injects extracted object-level graph information with fine-grained attributes into the diffusion model;
- *Multi-Layered Sampler* (MLS) that leverages the diverse conditions produced by SL-VAE to implement object-level graph manipulation.



# Framework Overview



**Scene Graph Representation: Nodes (objects) + Edges (relationships)**

$$G = (O, E) = O = \{o_i\}_{i=1}^{N_o} + E = \{e_{ij}\}_{1 \leq i, j \leq N_o, i \neq j}$$



# Experiments

A sheep by another sheep on the grass with the ocean under the sky; the ocean by a tree; a boat on the grass



A building with a window on side of a bus has a tire



Three turkeys on top of the pasture



SD-XL      DALL·E 3      Imagen 2      Ours

(a) Comparison with T2I methods in spatial relationships and object quantities.

Two men standing on a beach playing a kite



A dog chasing a cat on the grass



A building with a window on side of a bus has a tire

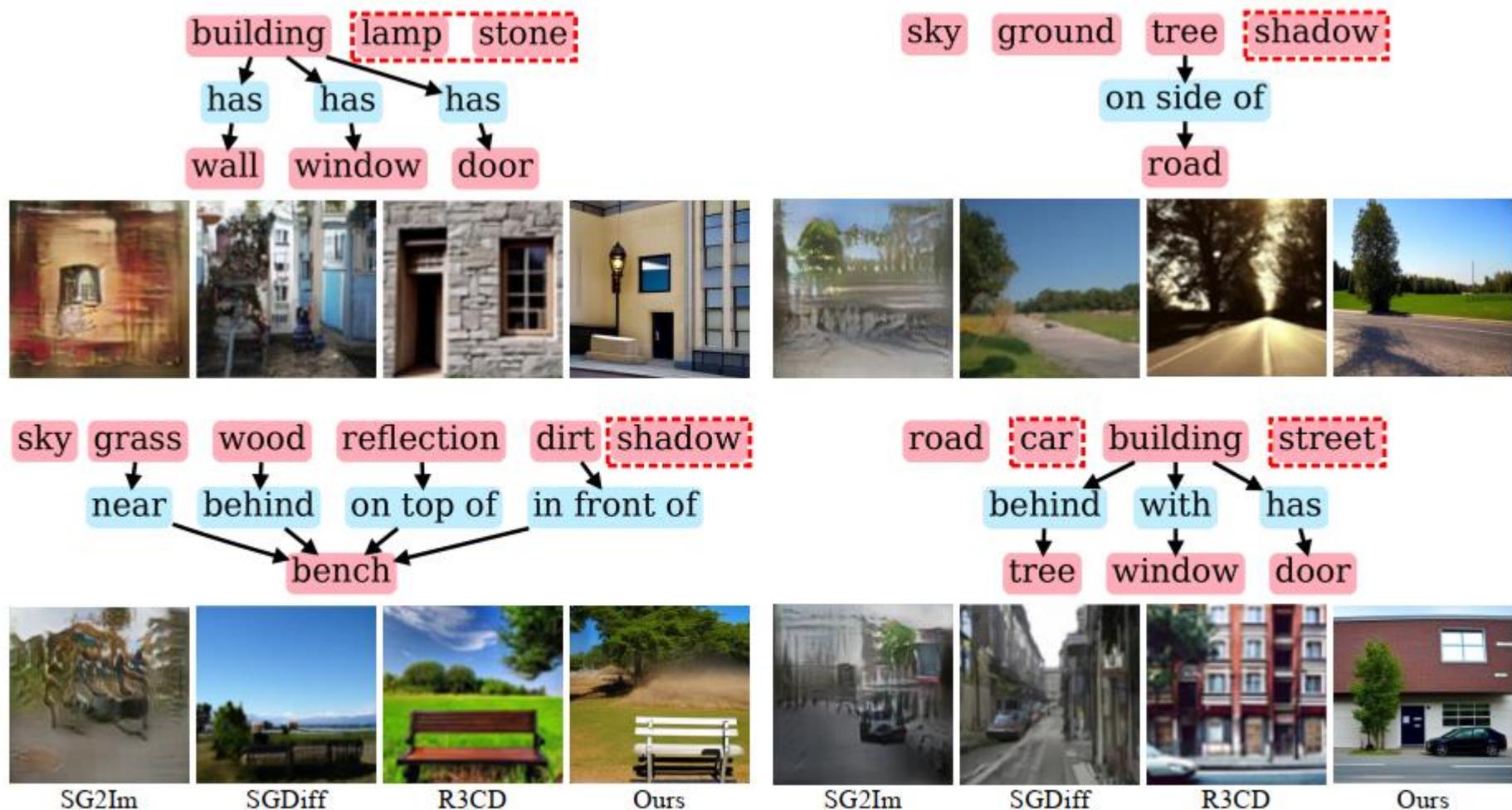


GLIGEN      LayoutDiffusion      MIGC      Ours

(b) Comparison with L2I methods in non-spatial interactions and rationality.



# Experiments

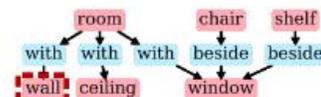
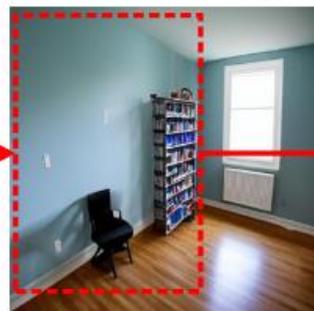
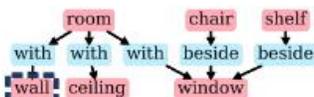
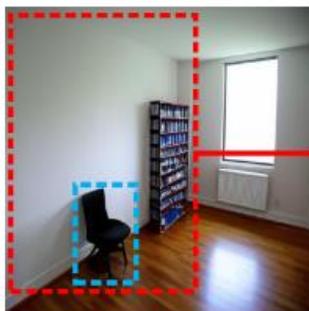
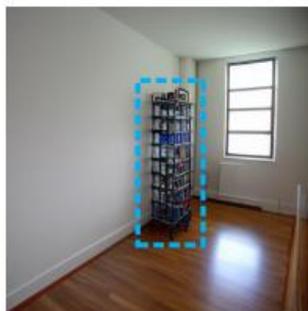
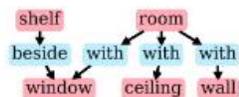
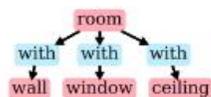
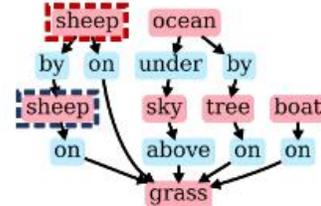
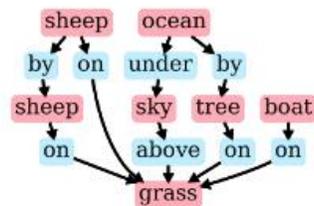
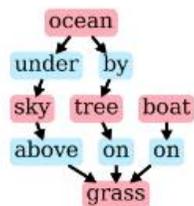
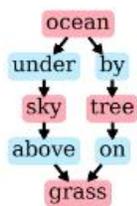


(c) Comparison with SG2I methods in *independent node inference* and *generation quality*.



# Experiments

Graph Manipulation (Node Addition and Attribute Control)



# Thanks for listening!

