

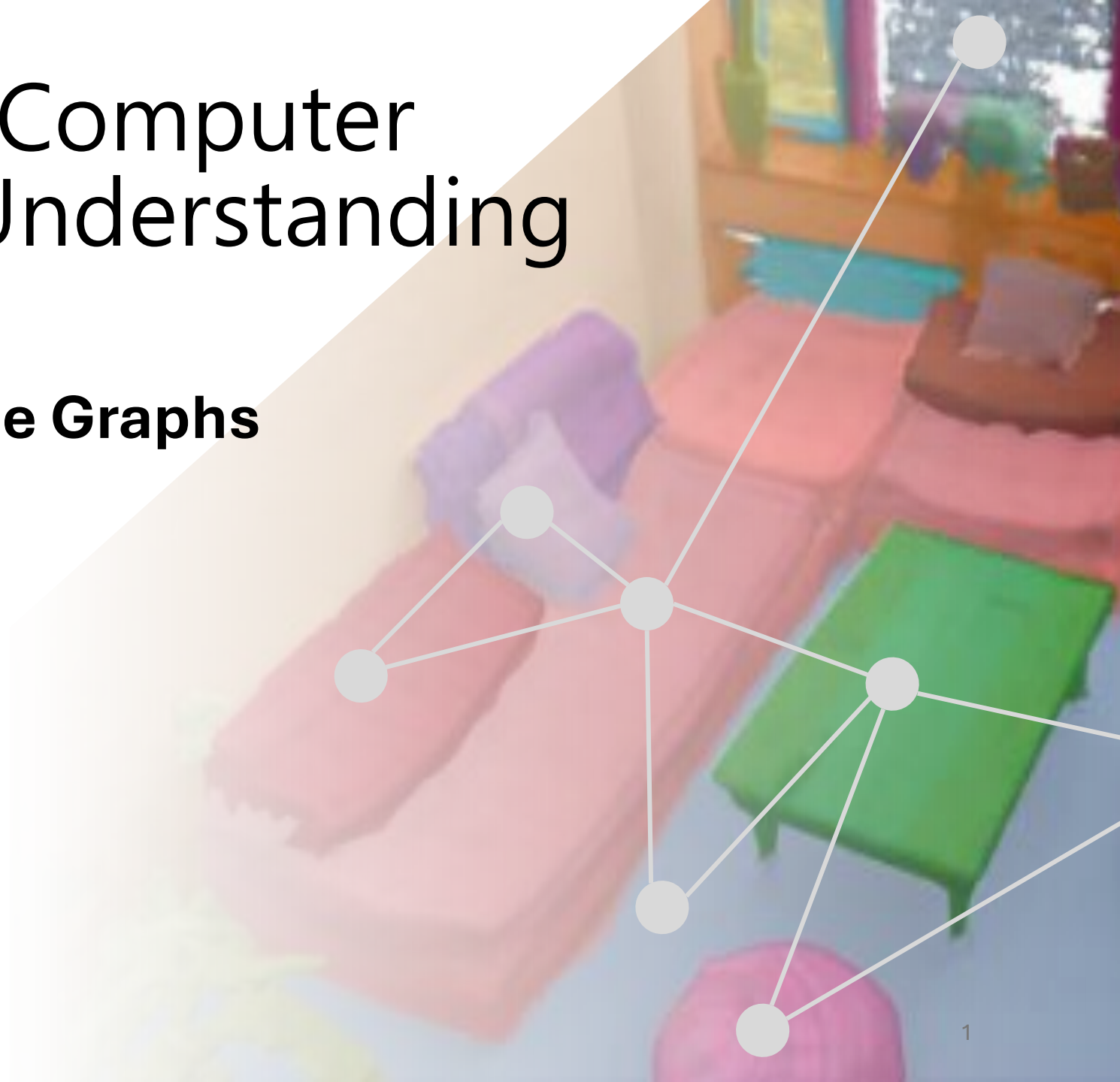
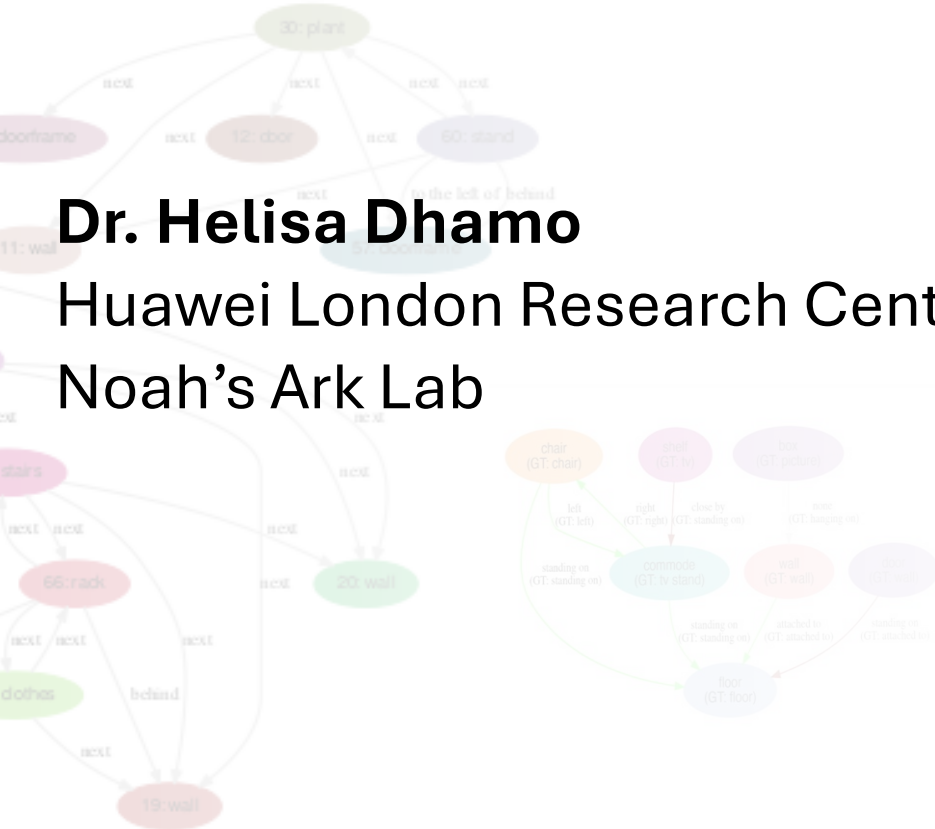
Deep Learning for Computer Vision and Scene Understanding

Lecture 4 – Semantic Scene Graphs

Dr. Helisa Dhamo

Huawei London Research Centre

Noah's Ark Lab



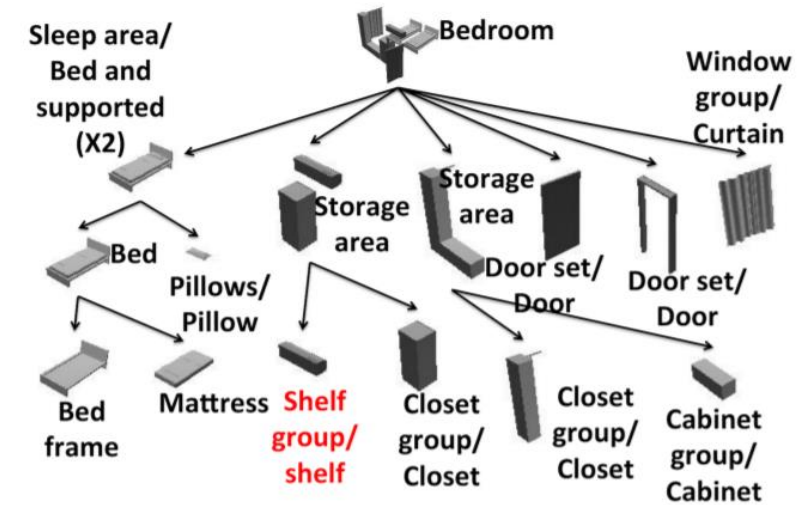
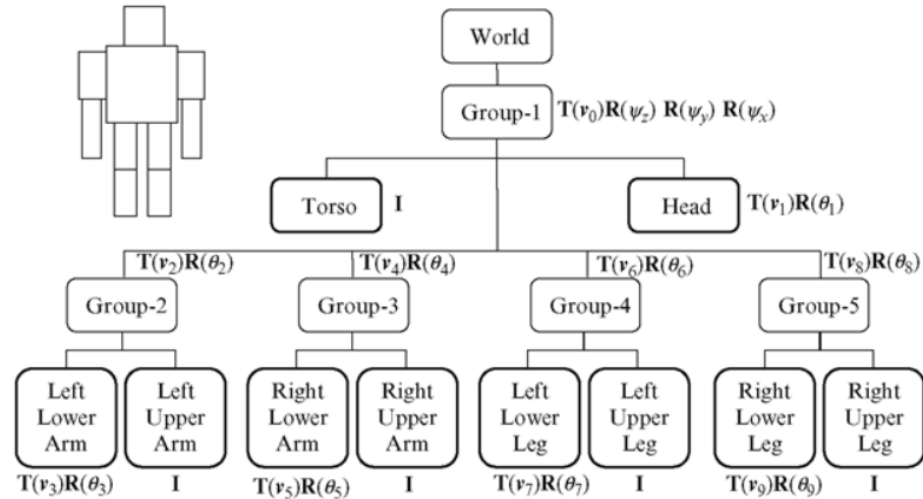


Invited lecture: Novel View Synthesis

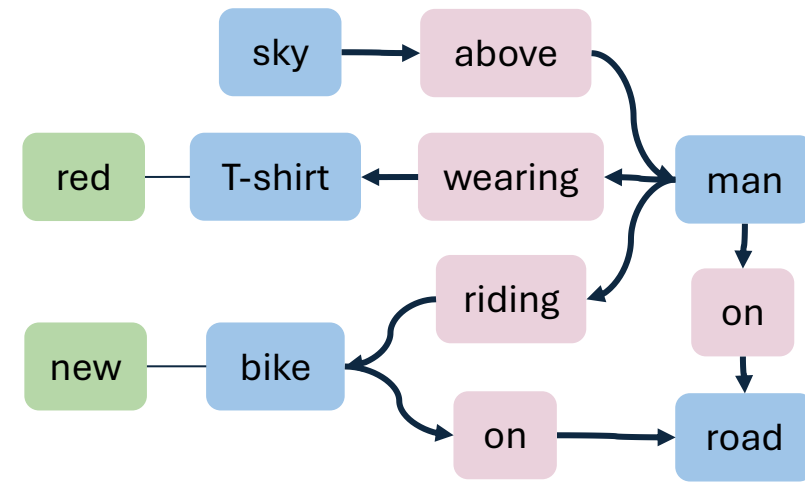
Semantic Scene Graphs

Lecture 4

Scene graph



Scene hierarchy [Liu TOG 2014]

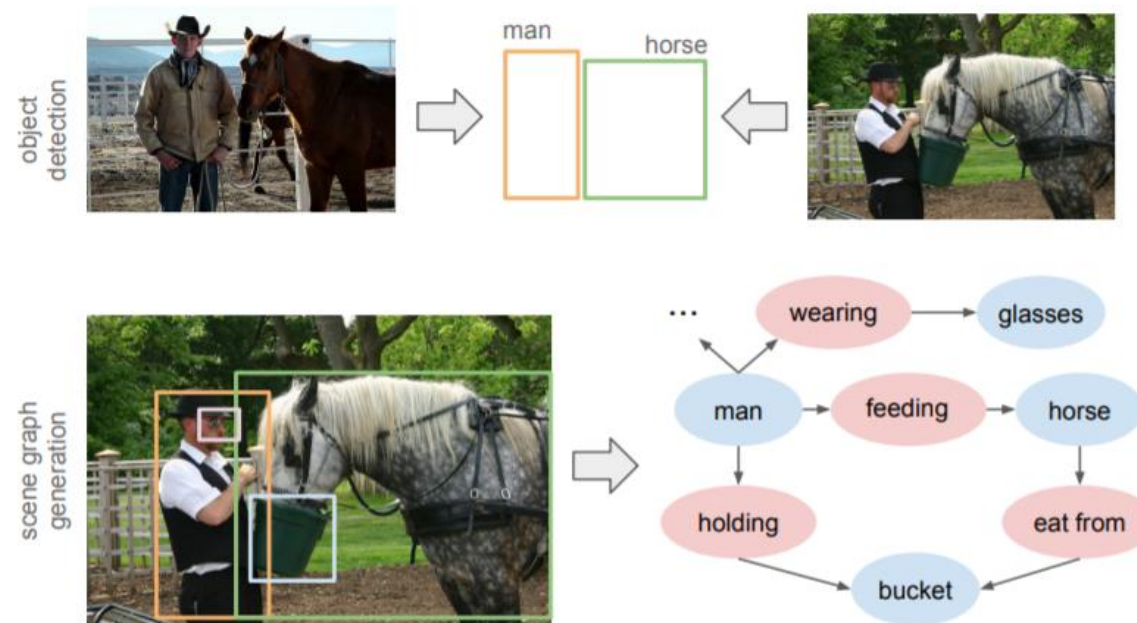


Semantic nodes and edges [Johnson CVPR 2015]

(Recap) Scene Understanding Beyond Objects

Scene Graphs

- **Nodes:** objects in the scene
- **Edges:** relationships between objects (interaction, relative position)



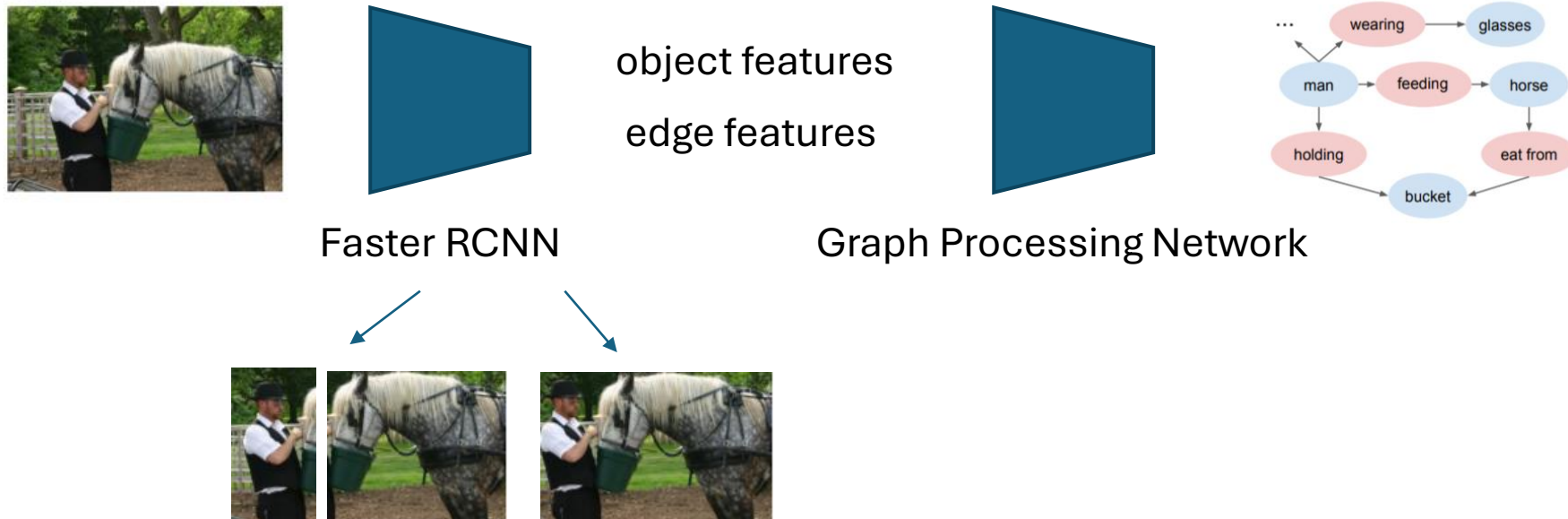
Xu et al., "Scene graph generation by iterative message passing." CVPR'17

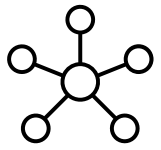
(Recap) Scene Understanding Beyond Objects

Scene Graphs

- **Nodes:** objects in the scene
- **Edges:** relationships between objects (interaction, relative position)

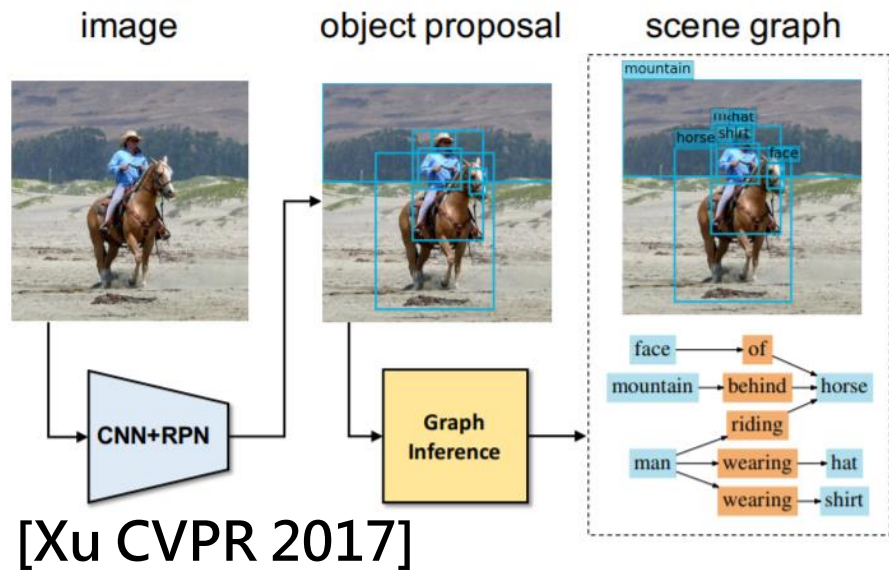
Scene graph generation networks are usually build on top of an object detector





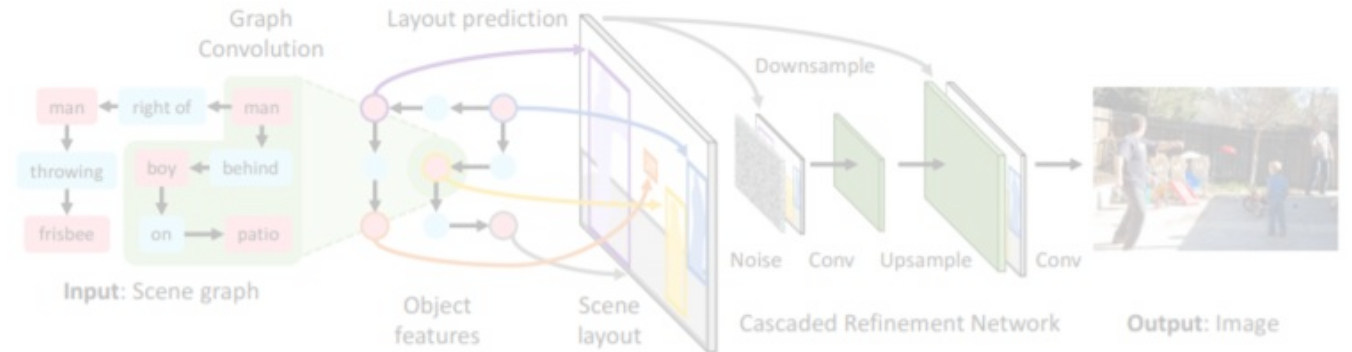
Semantic Scene Graphs and Images

From image to scene graph



Scene Understanding

From scene graph to image

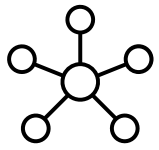


[Johnson CVPR 2018] Purely semantic nodes (object class)

Scene Synthesis

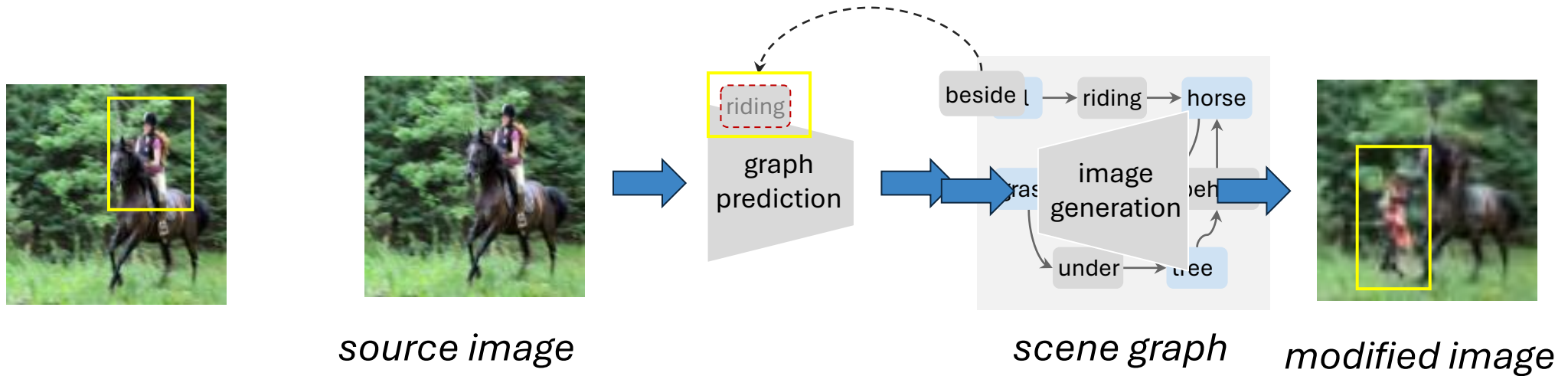
J. Johnson, A. Gupta, and FF. Li. Image generation from scene graphs. **CVPR 2018**.

D. Xu, Y. Zhu, C.r Choi, FF. Li. Scene Graph Generation by Iterative Message Passing. **CVPR 2017**



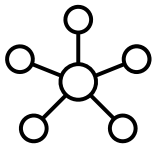
Semantic Image Editing The combined problem...

The generated image is fed back into the graph (and design the graph)



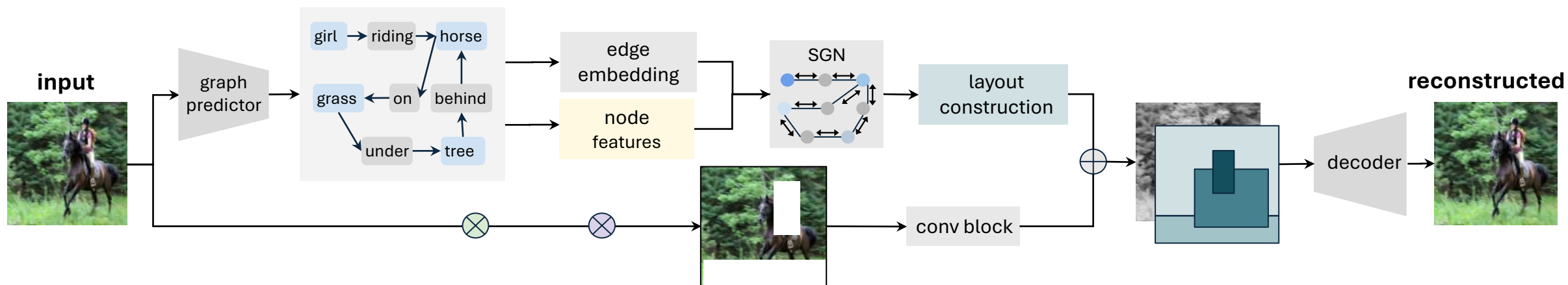
Challenge: No real image pairs with changes!

at test time

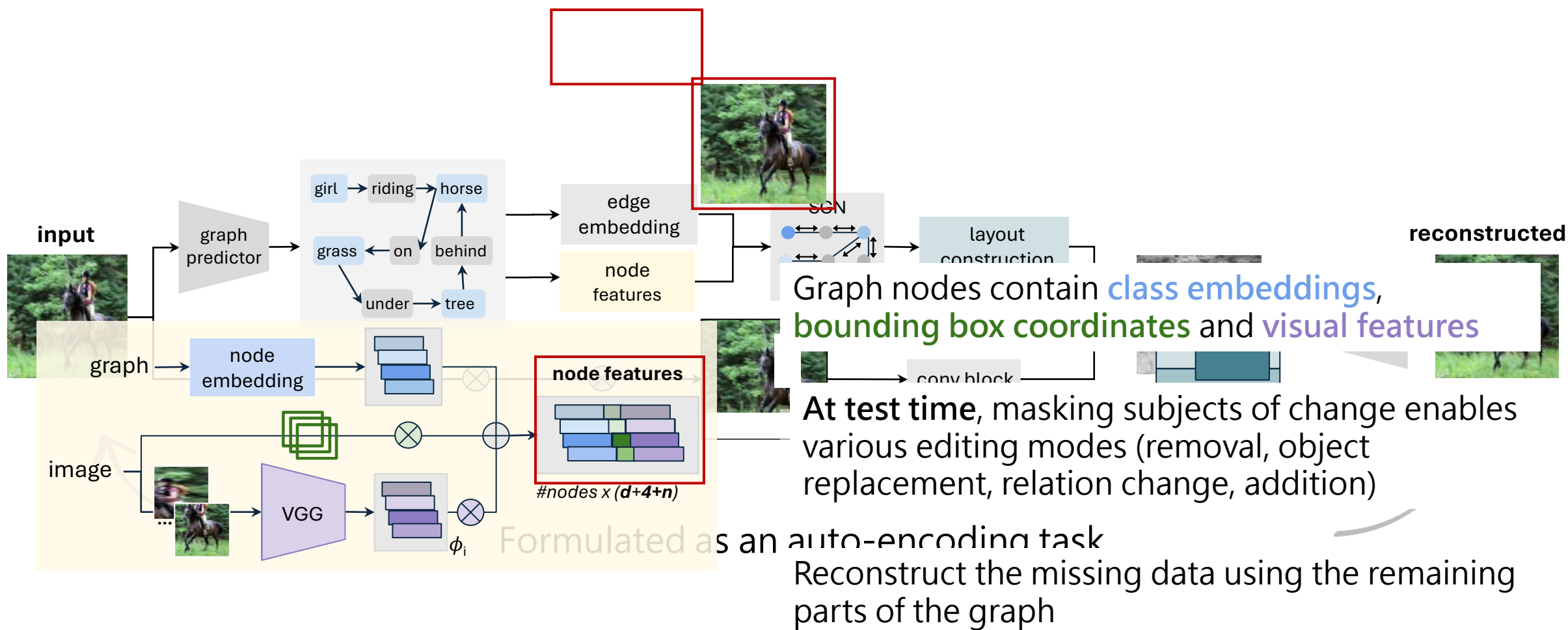


Semantic Image Editing Training

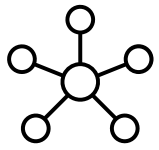
Our training strategy does not require pairs for the editing task



at training time

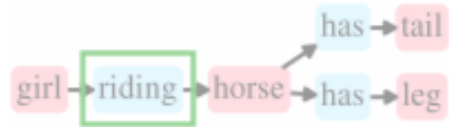


at training time

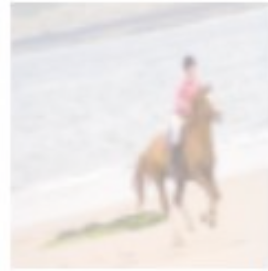


Semantic Image Editing Results

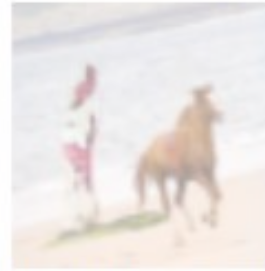
original graph



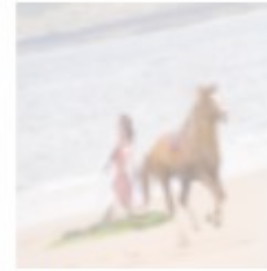
source



ours CRN

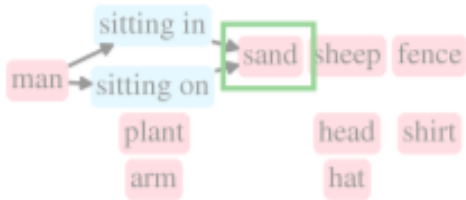


ours SPADE



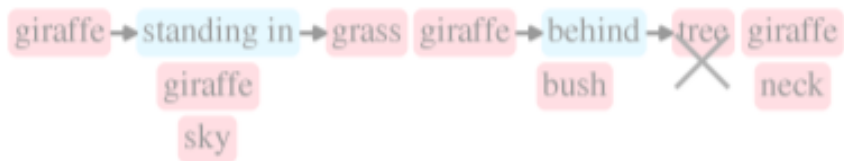
"riding" to "next to"

relationship change



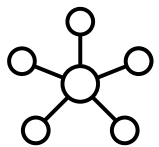
"sand" to "ocean"

object replacement

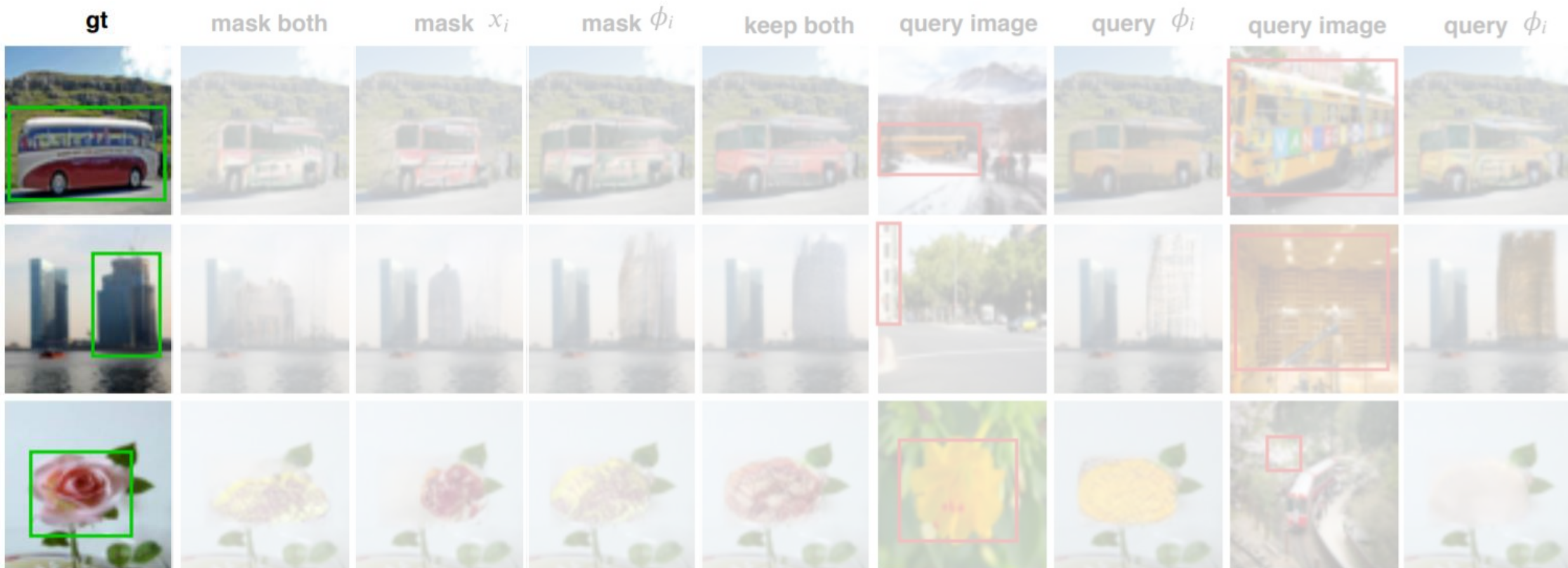


remove "tree"

object removal

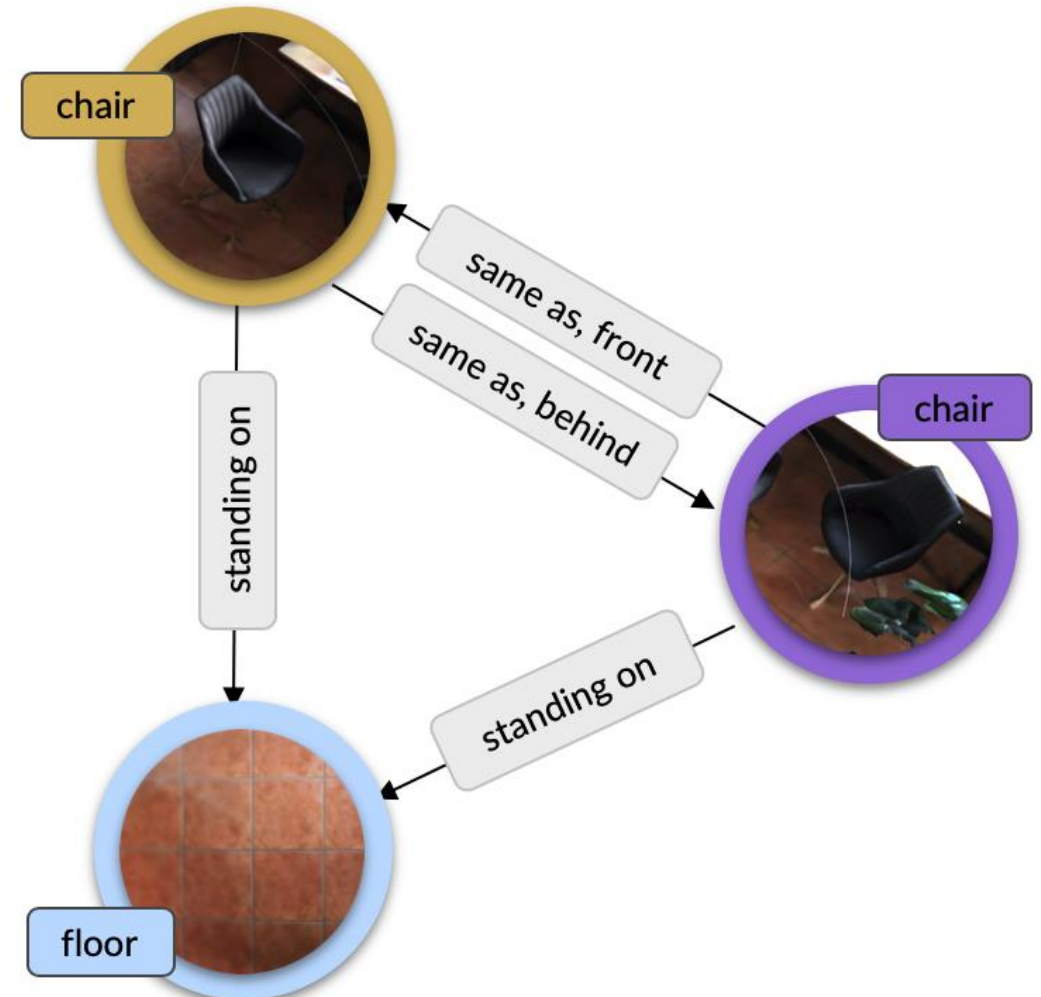


Semantic Image Editing Results



Known bounding coordinates
Known bounding coordinates

3D Semantic Scene Graphs





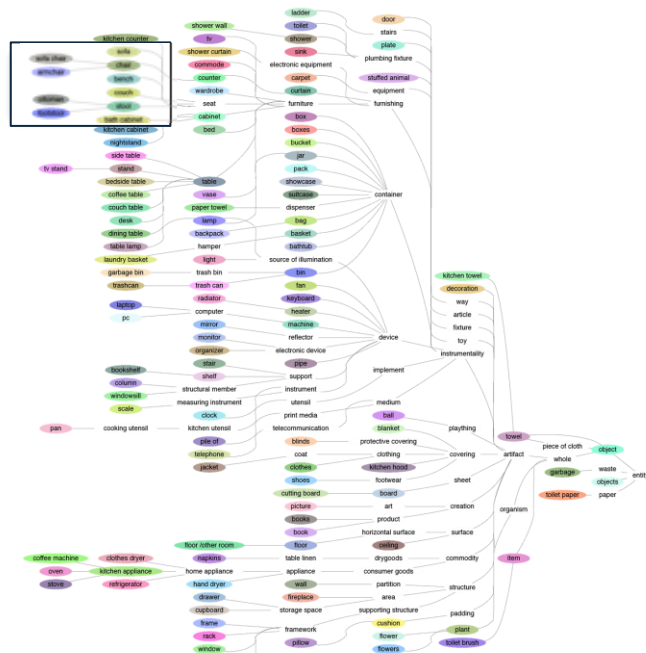
3DSSG

3D Semantic Scene Graphs

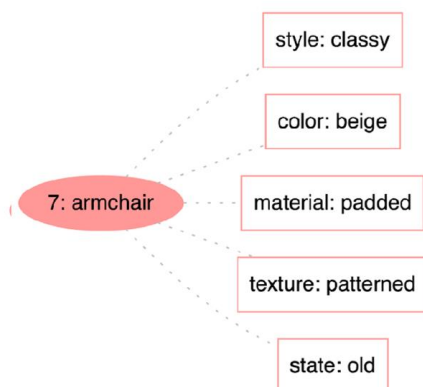
Based on 3RScan Dataset [Wald 2019]

Based on 3RScan Dataset [Wald 2019]

Hierarchy of semantic class labels

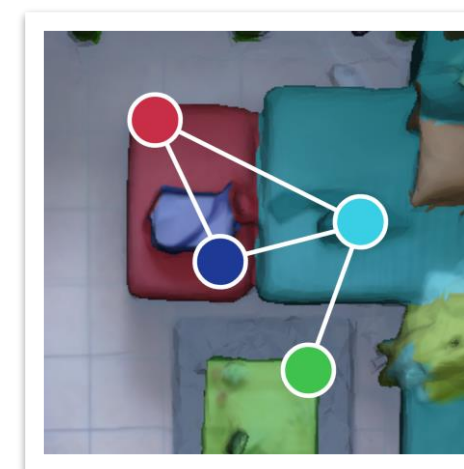


Static (color, material, shape) and **dynamic** attributes (tidy/messy, open/closed) and **affordances** (sitting, eating)



Proximity (left/right, front/behind, close by)
Support (lying in, hanging, leaning against)
Comparative (smaller than, same as)

Proximity (left/right, front/behind, close by)
Support (lying in, hanging, leaning against)
Comparative (smaller than, same as)

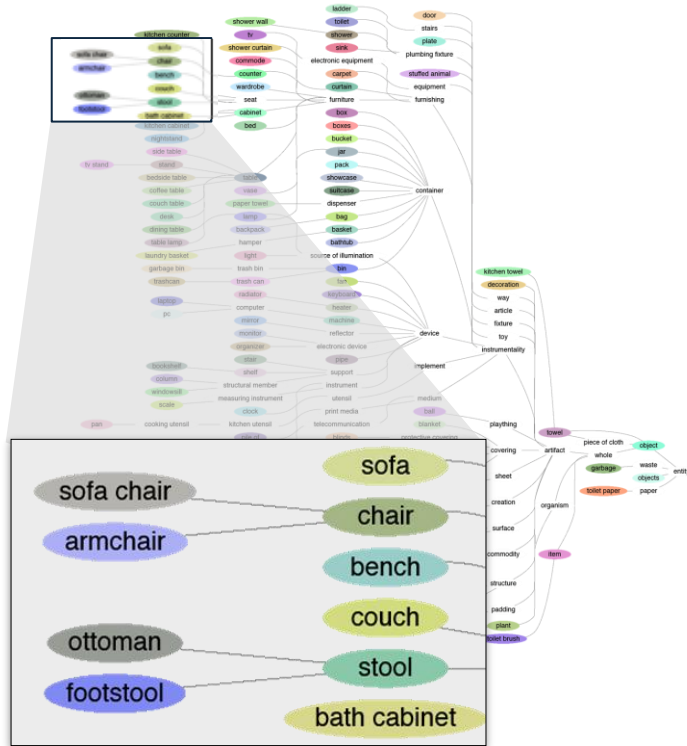


3DSSG Dataset

Based on 3RScan [Wald 2019]

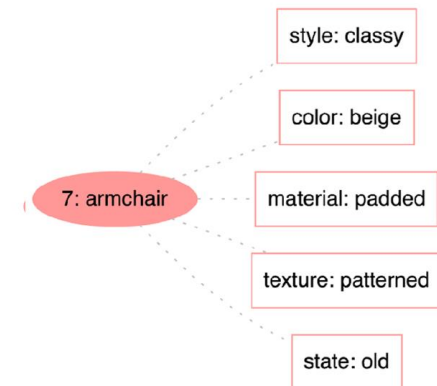
Nodes

Hierarchy of semantic class labels



Attributes

Static (color, material, shape) and **dynamic** attributes (tidy/messy, open/closed) and **affordances** (sitting, eating)

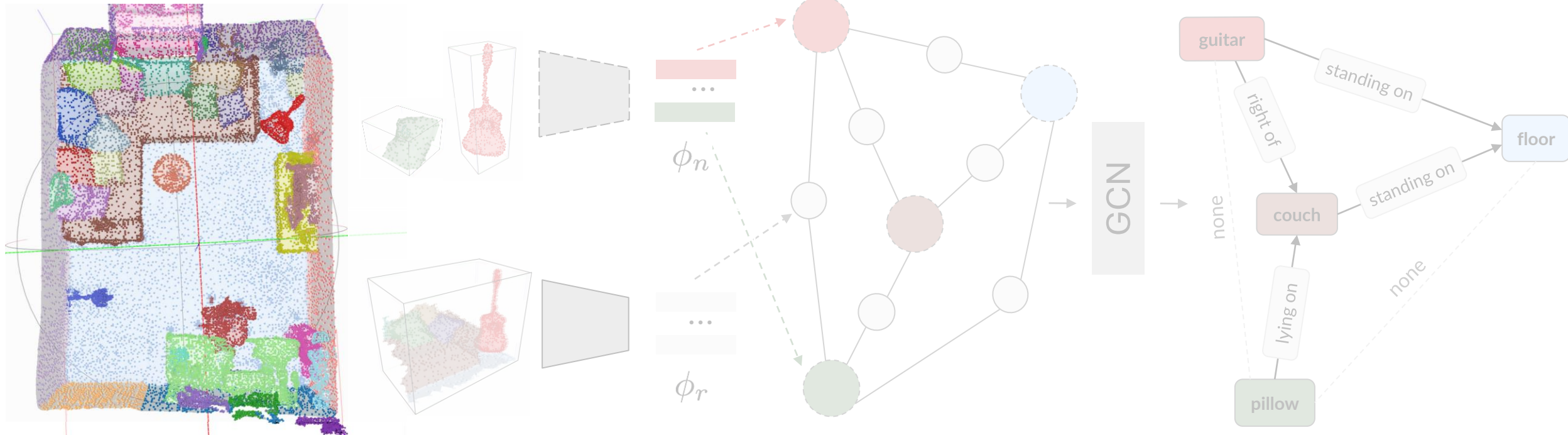


Relationships

Proximity (left/right, front/behind, close by)
Support (lying in, hanging, leaning against)
Comparative (smaller than, same as)



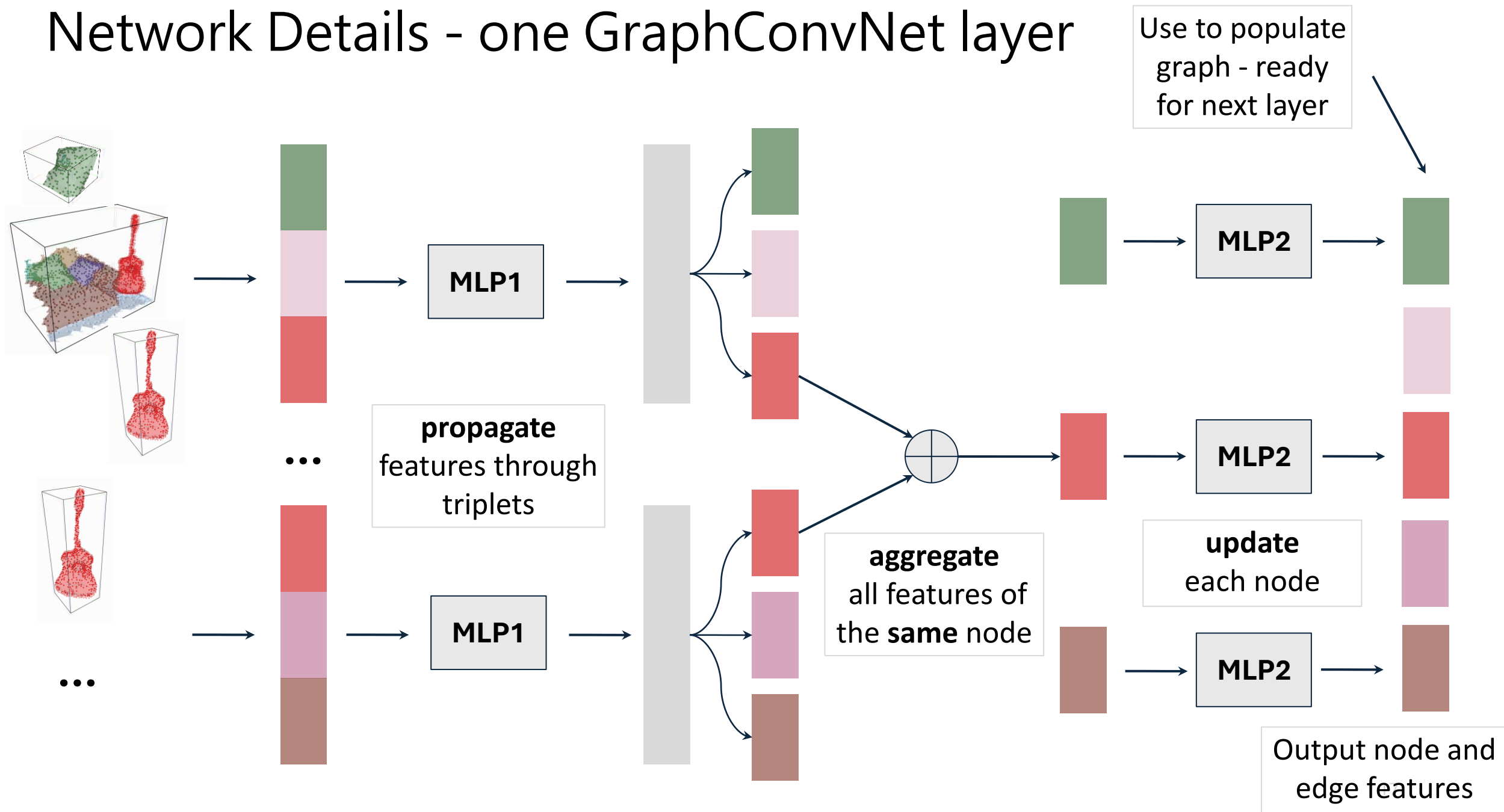
Learning 3D Semantic Scene Graphs



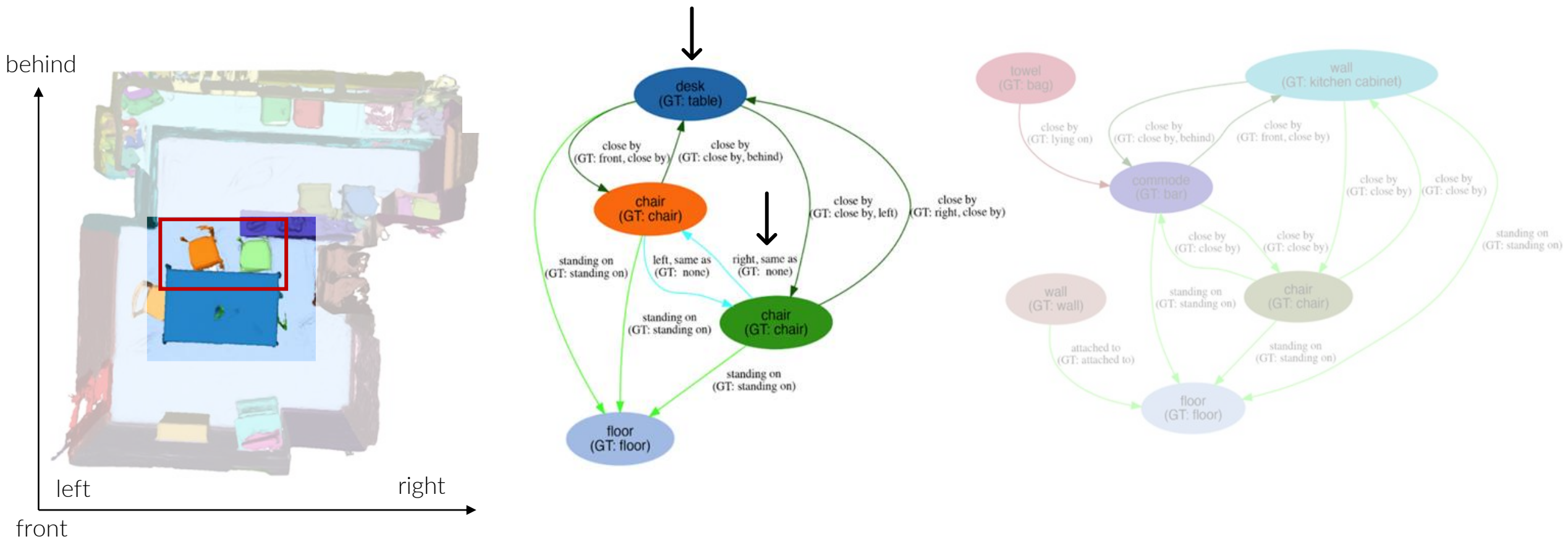
$$\mathcal{L}_{\text{total}} = \lambda_{\text{obj}} \mathcal{L}_{\text{obj}} + \mathcal{L}_{\text{pred}} \quad \mathcal{L} = -\alpha_t (1 - p_t)^\gamma \log p_t$$

none or multiple
predicate pre-
dictions per edge

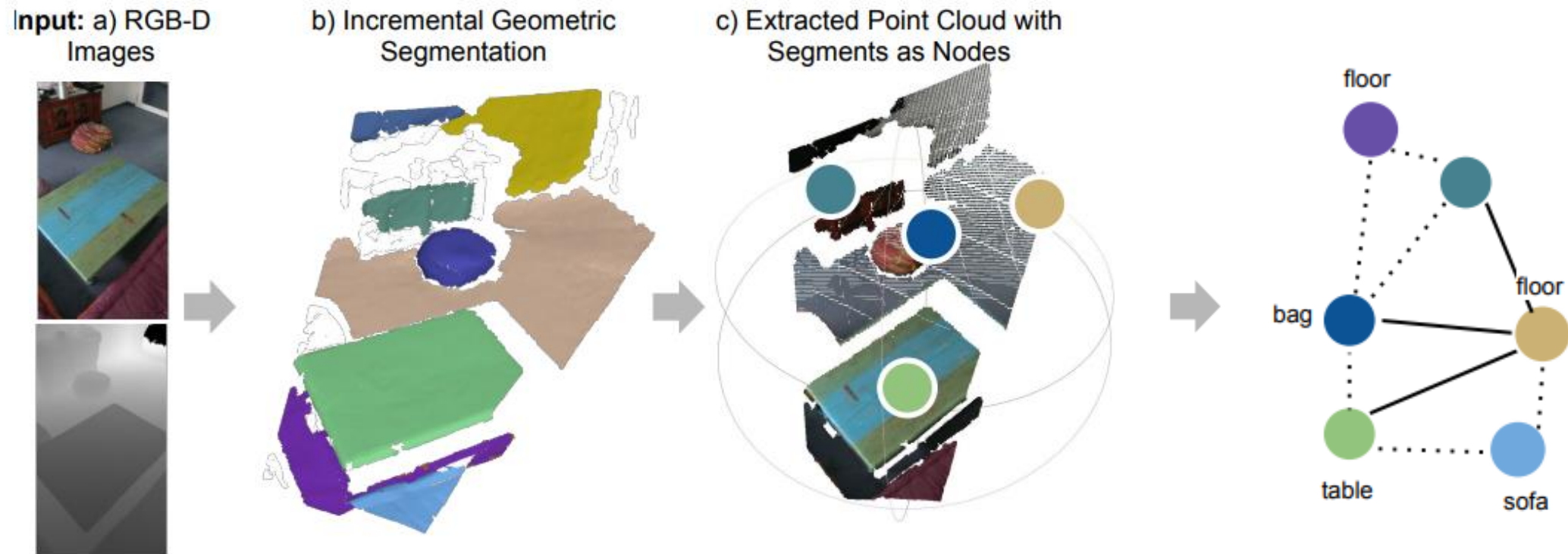
Network Details - one GraphConvNet layer



Learning 3D Semantic Scene Graphs Results



Follow-up research Towards real world requirements

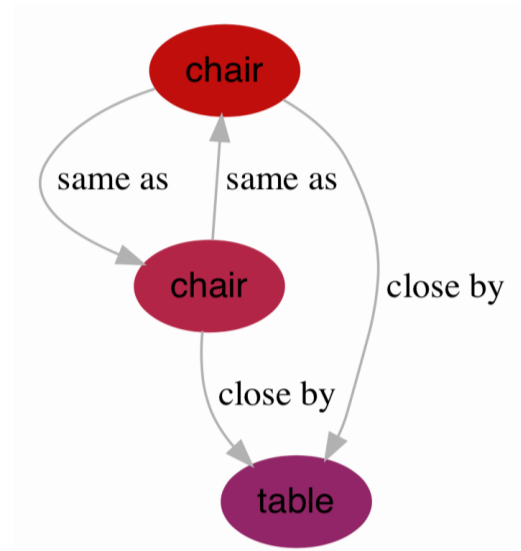
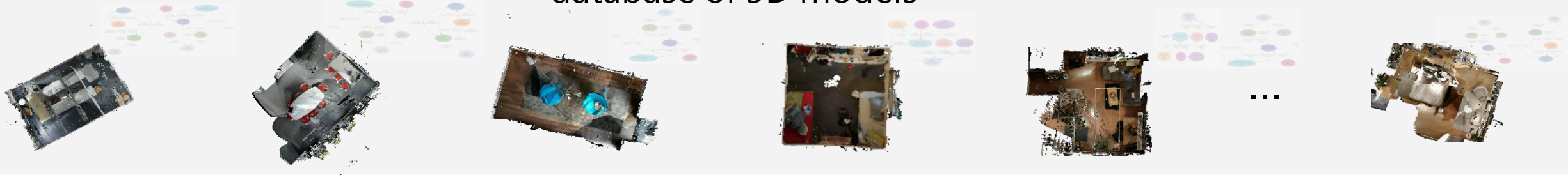


Build scene graph incrementally as the scene is reconstructed
No need for class-agnostic instance segmentation

Applications

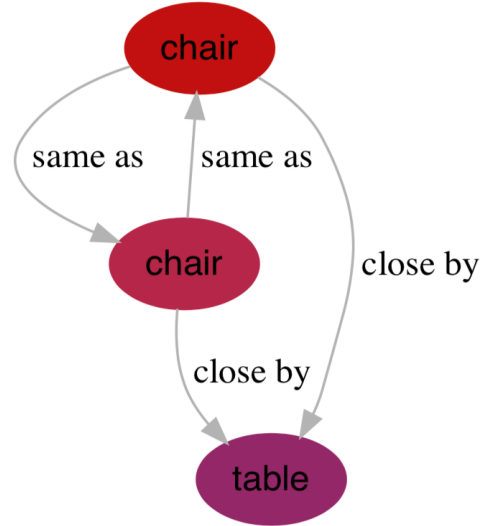
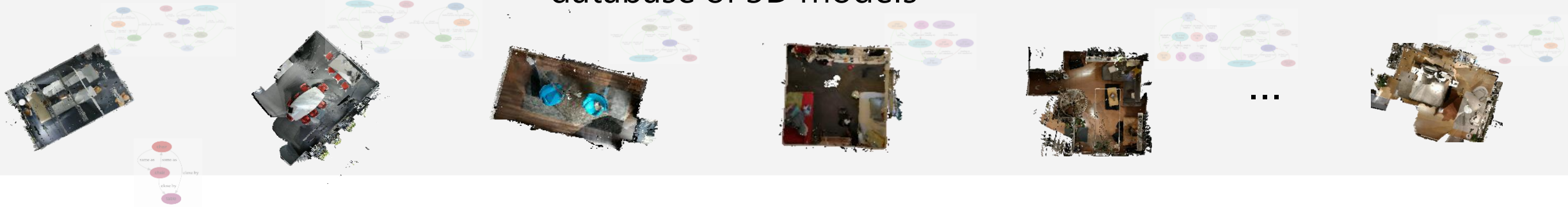
Domain agnostic scene retrieval

database of 3D models



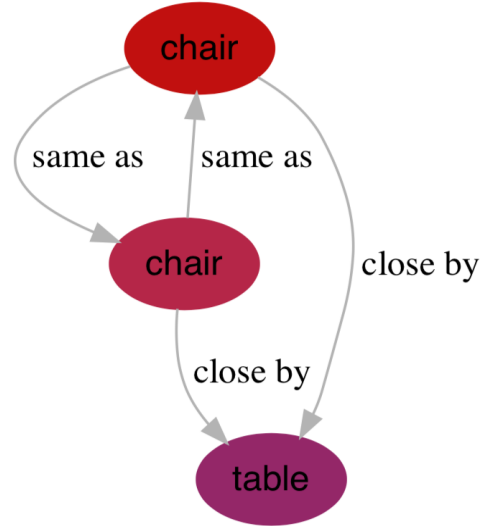
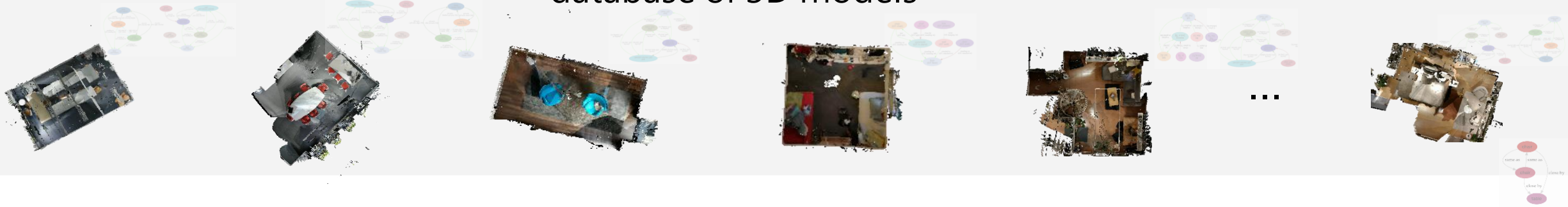
query photo and its scene graph

database of 3D models



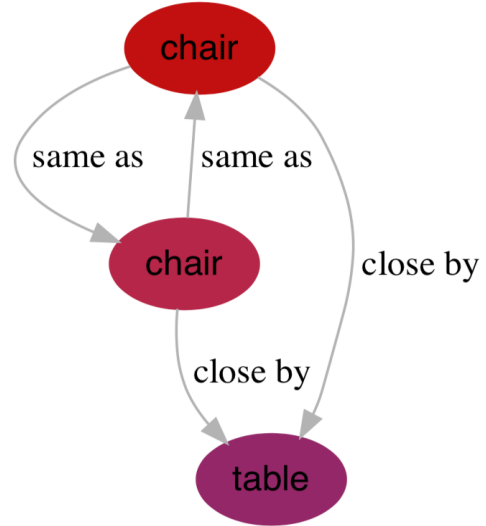
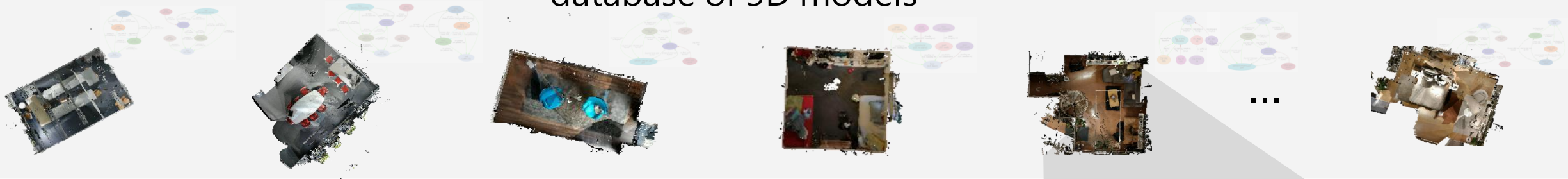
query photo and its scene graph

database of 3D models



query photo and its scene graph

database of 3D models



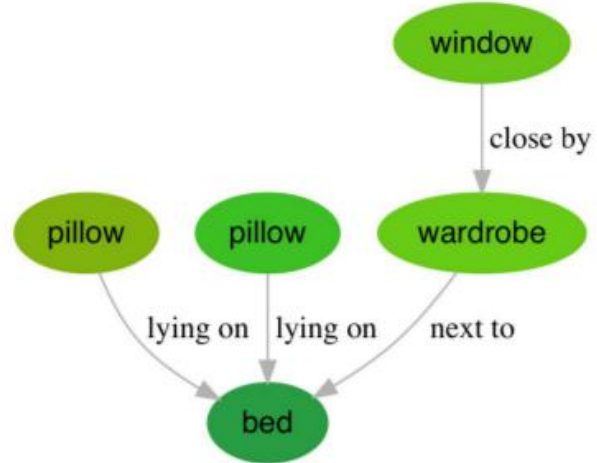
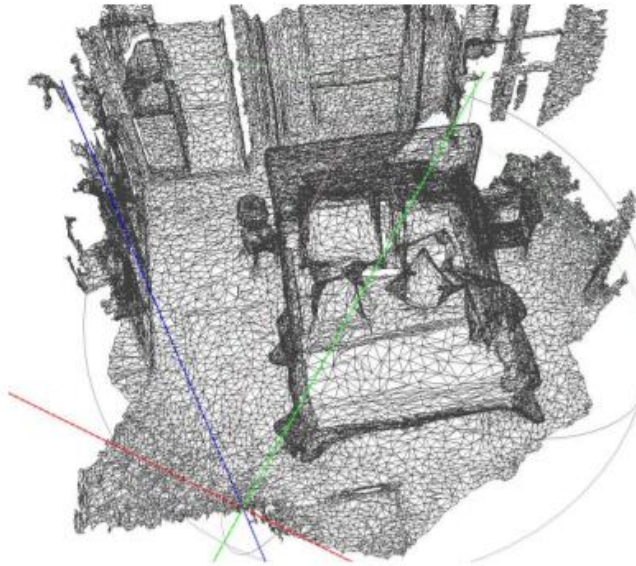
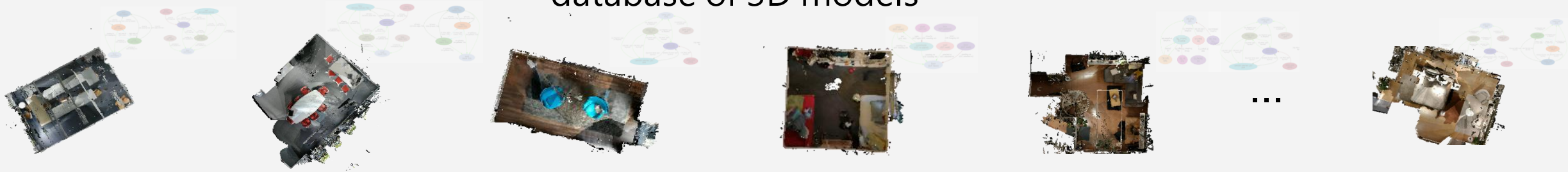
Query based on
node set and triplet
set



query photo and its scene graph

best match

database of 3D models



query 3D scene and its scene graph

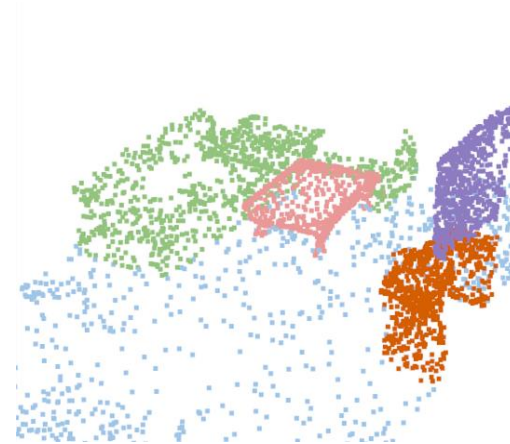
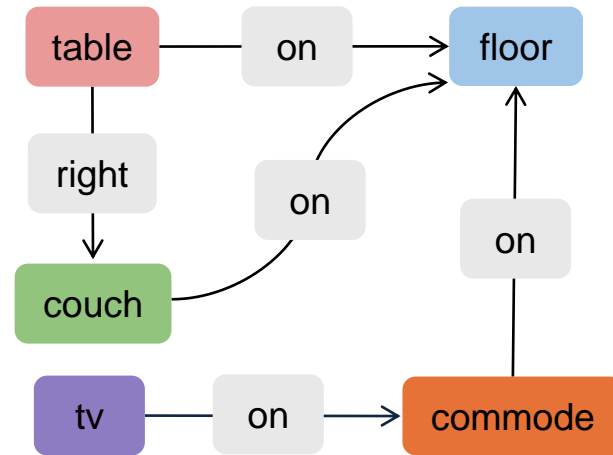


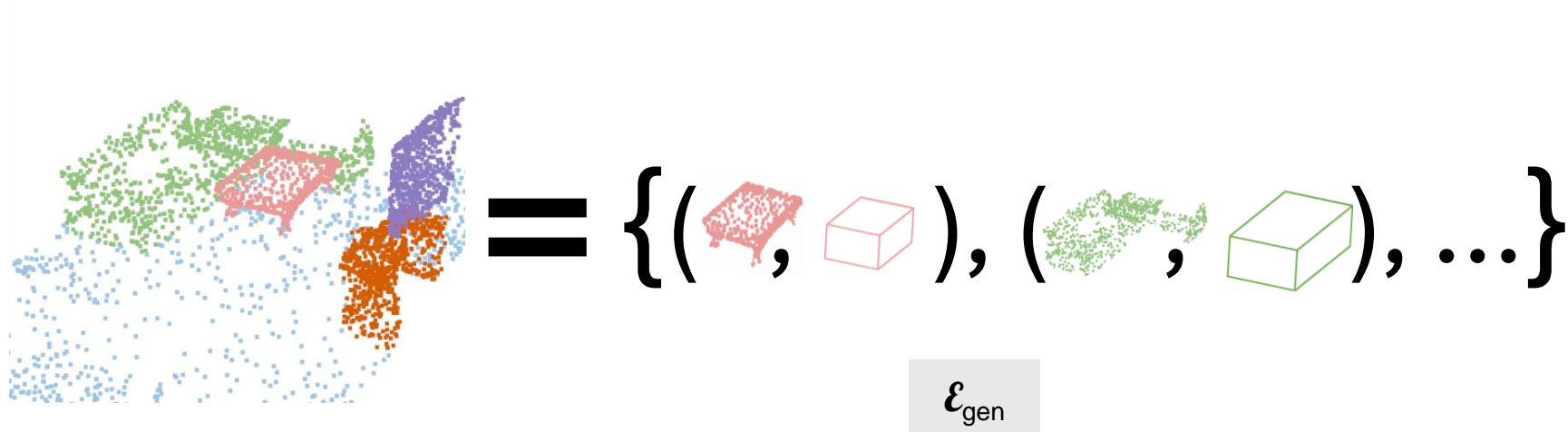
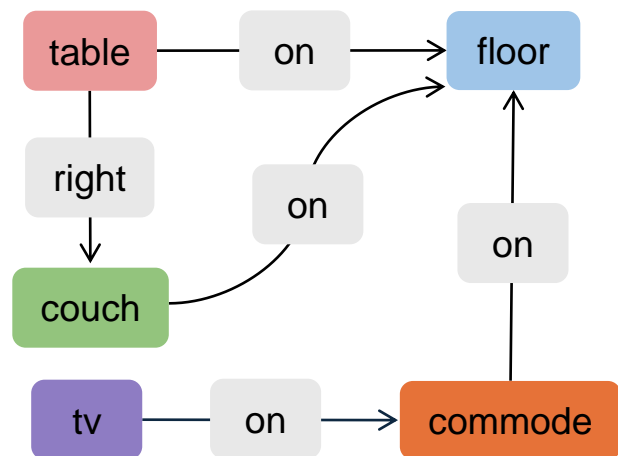
best match

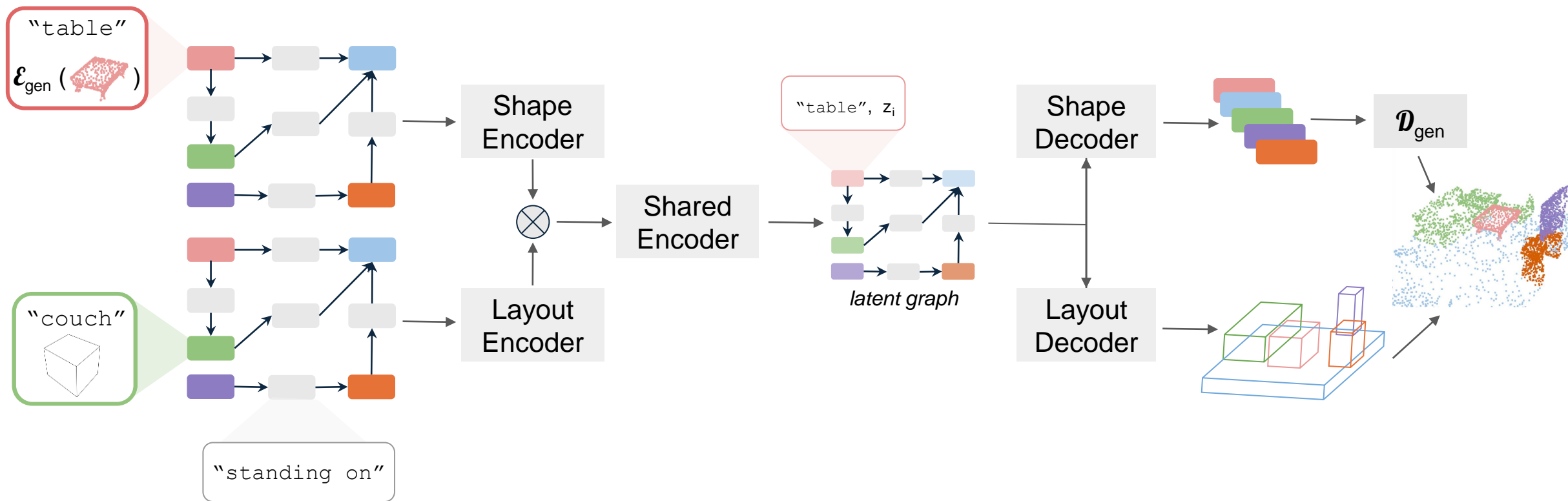
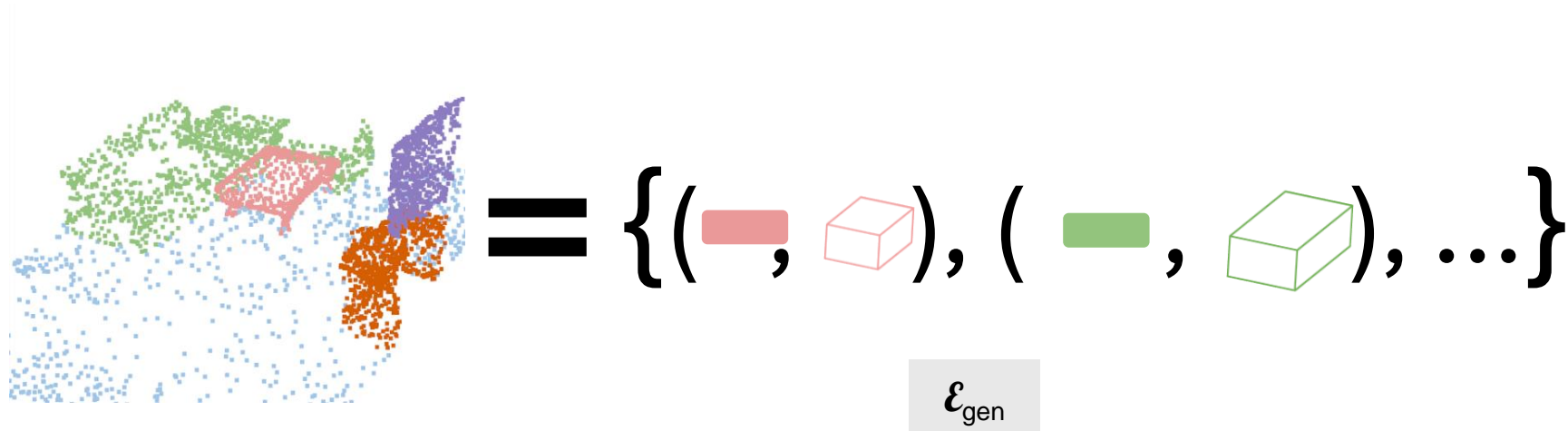
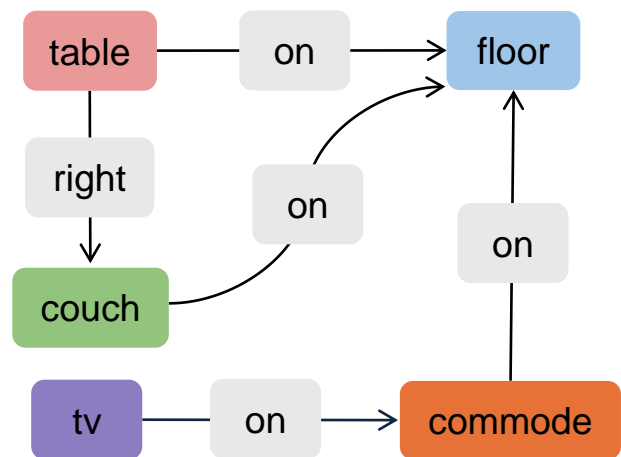
Applications

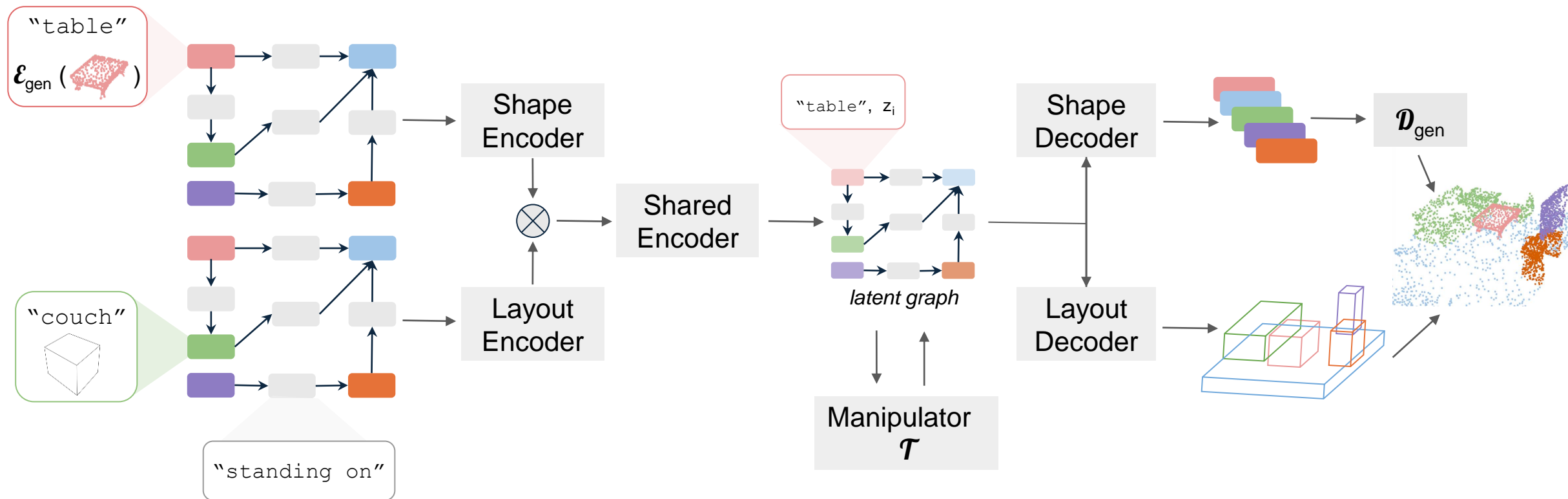
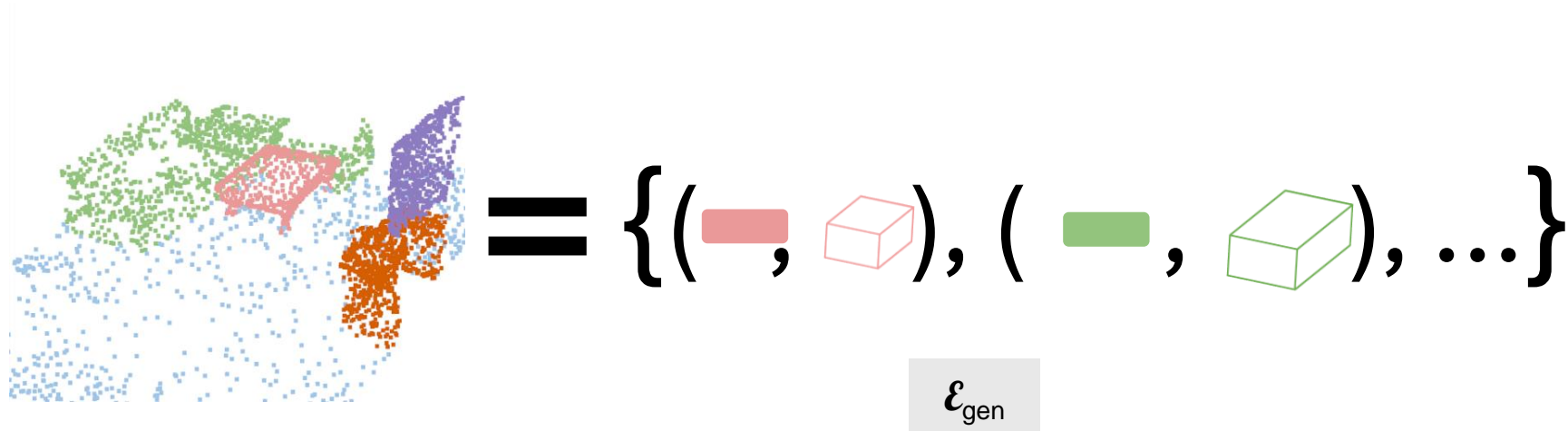
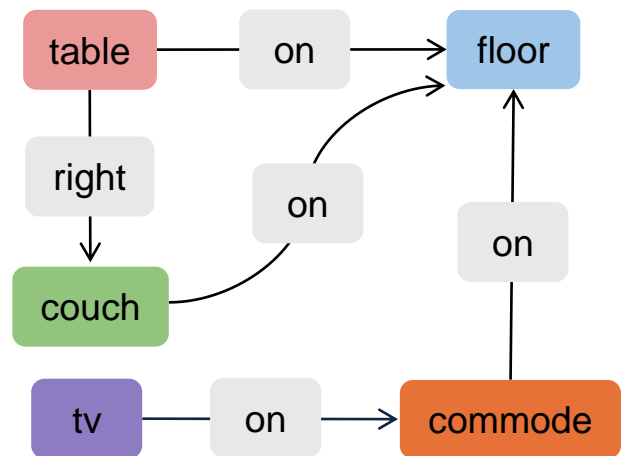
Scene generation and editing

$$= \{ (\text{point cloud}, \text{box}), (\text{point cloud}, \text{box}), \dots \}$$

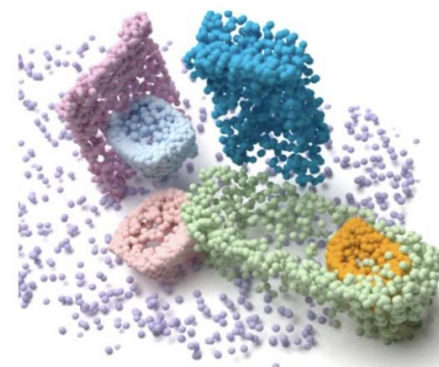
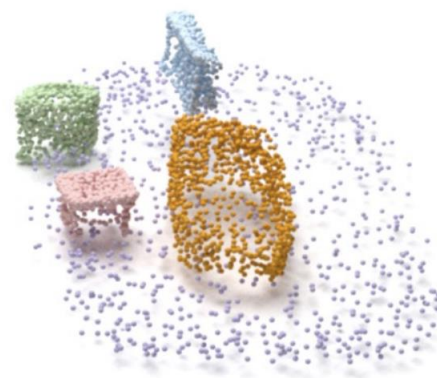
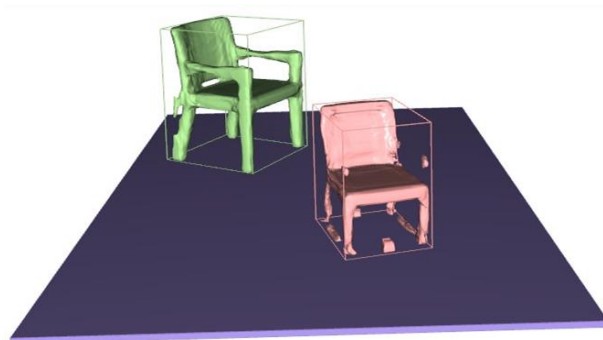
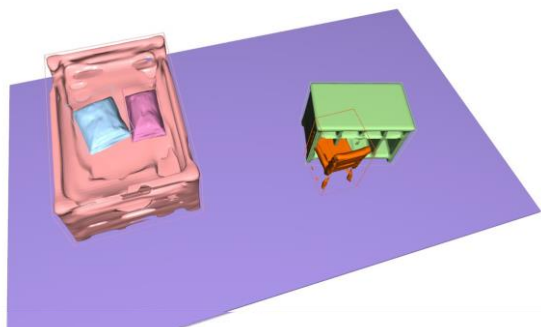
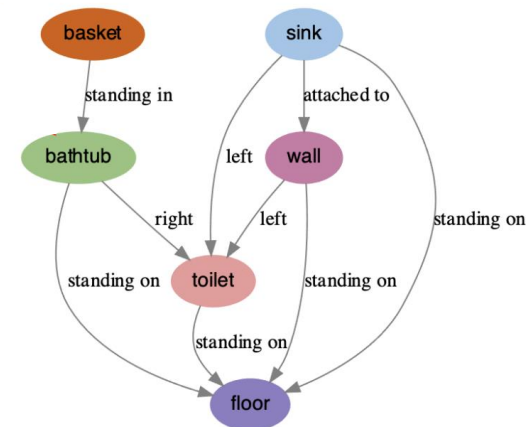
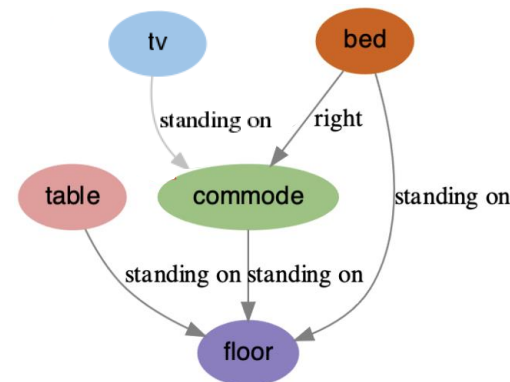
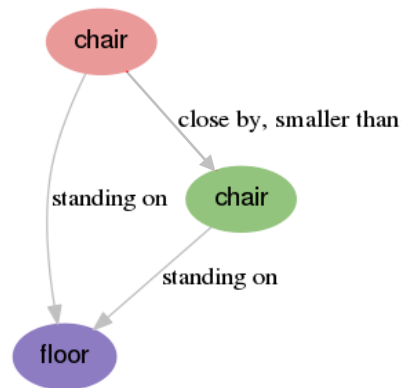
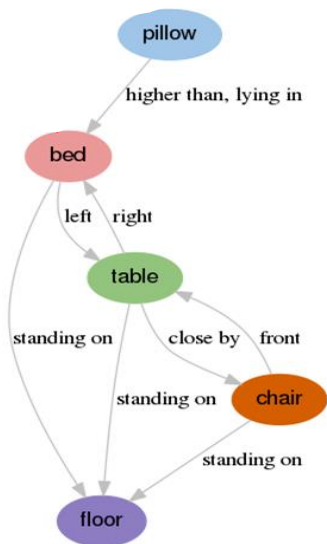








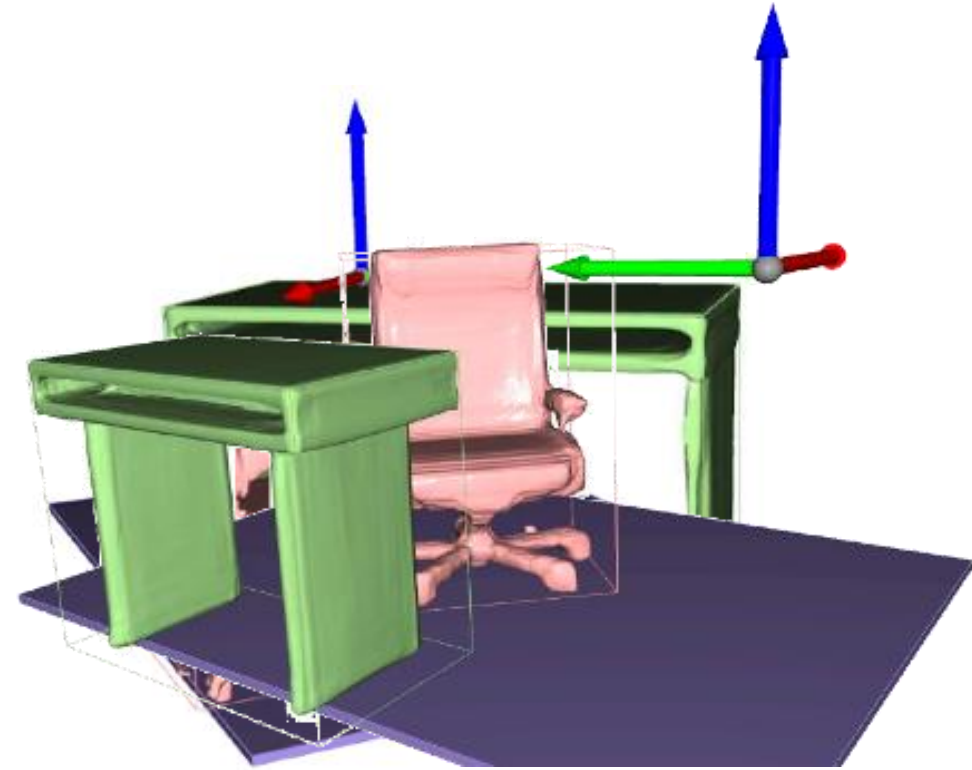
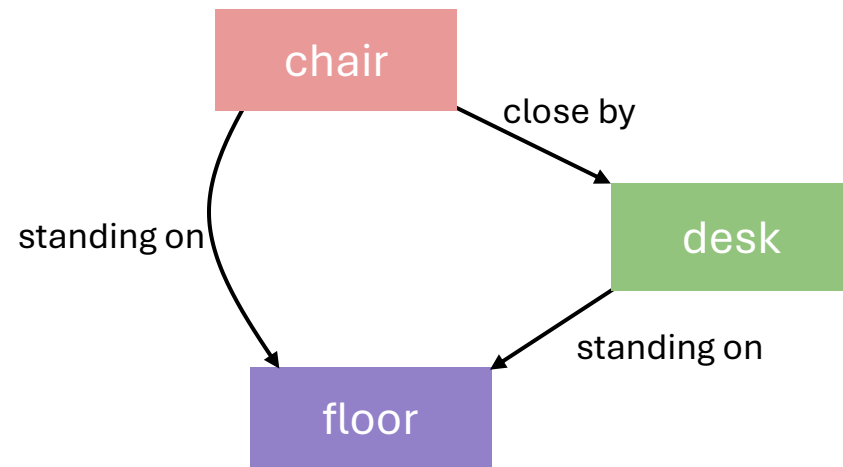
Graph-to-3D Results



Meshes from DeepSDF encoding

Point clouds from AtlasNet encoding

Graph-to-3D Context learning results

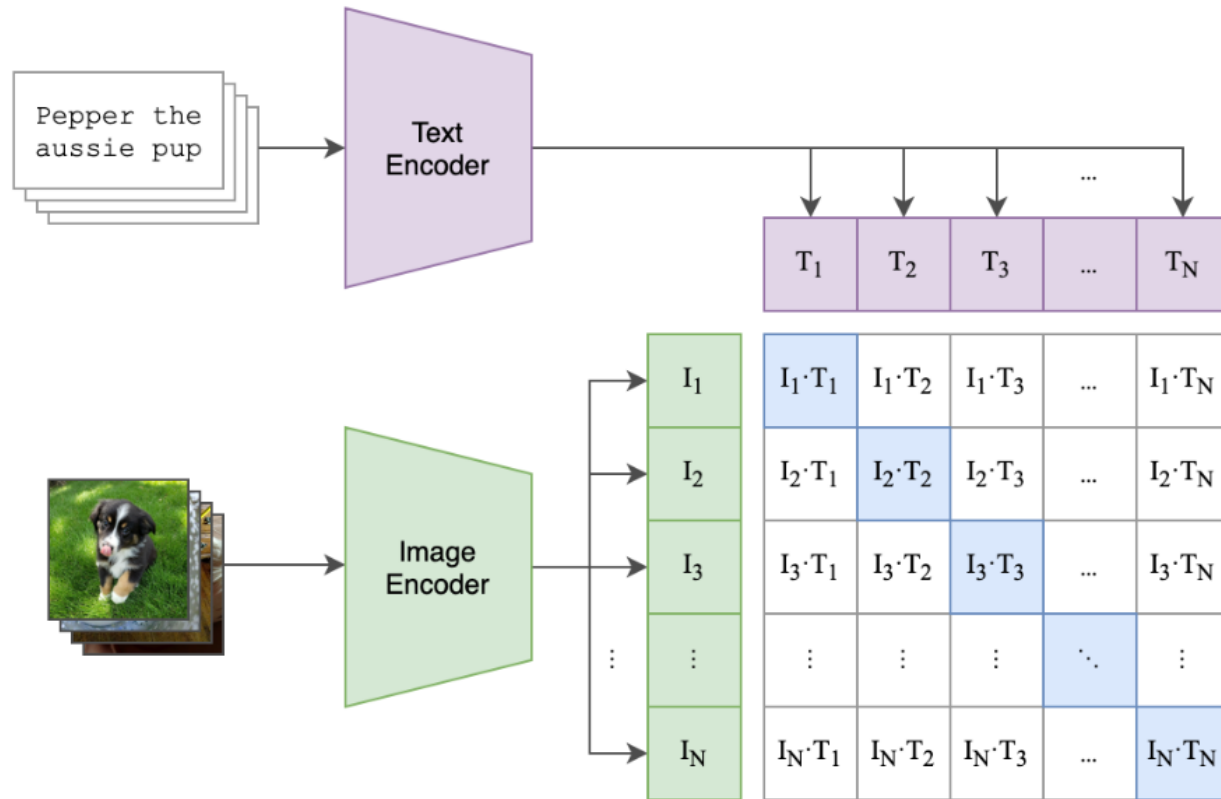


Open problem

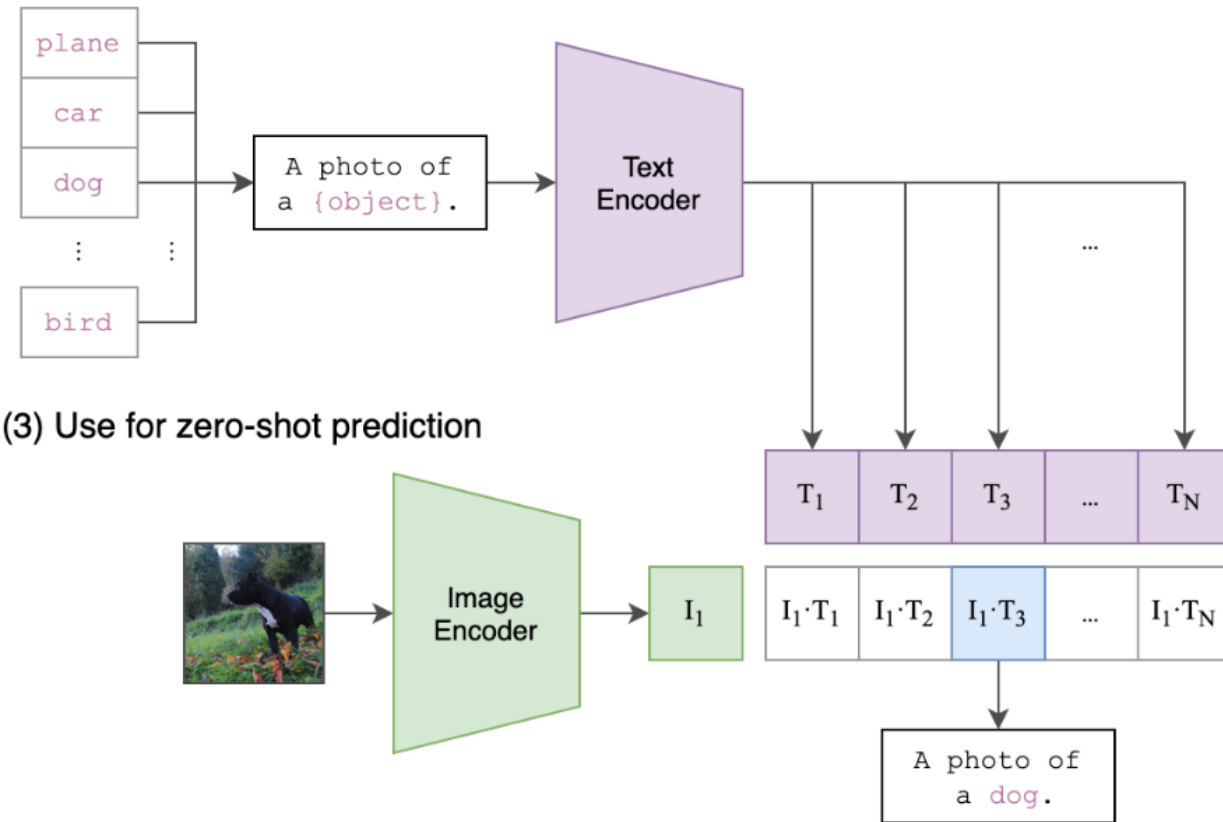
- **Fixed vocabulary:** Pre-defined set of semantic class categories for objects and relationships
- Why is that a problem?
- Solution **Open vocabulary** 3D Scene Graphs

CLIP: Contrastive Language-Image Pretraining

(1) Contrastive pre-training

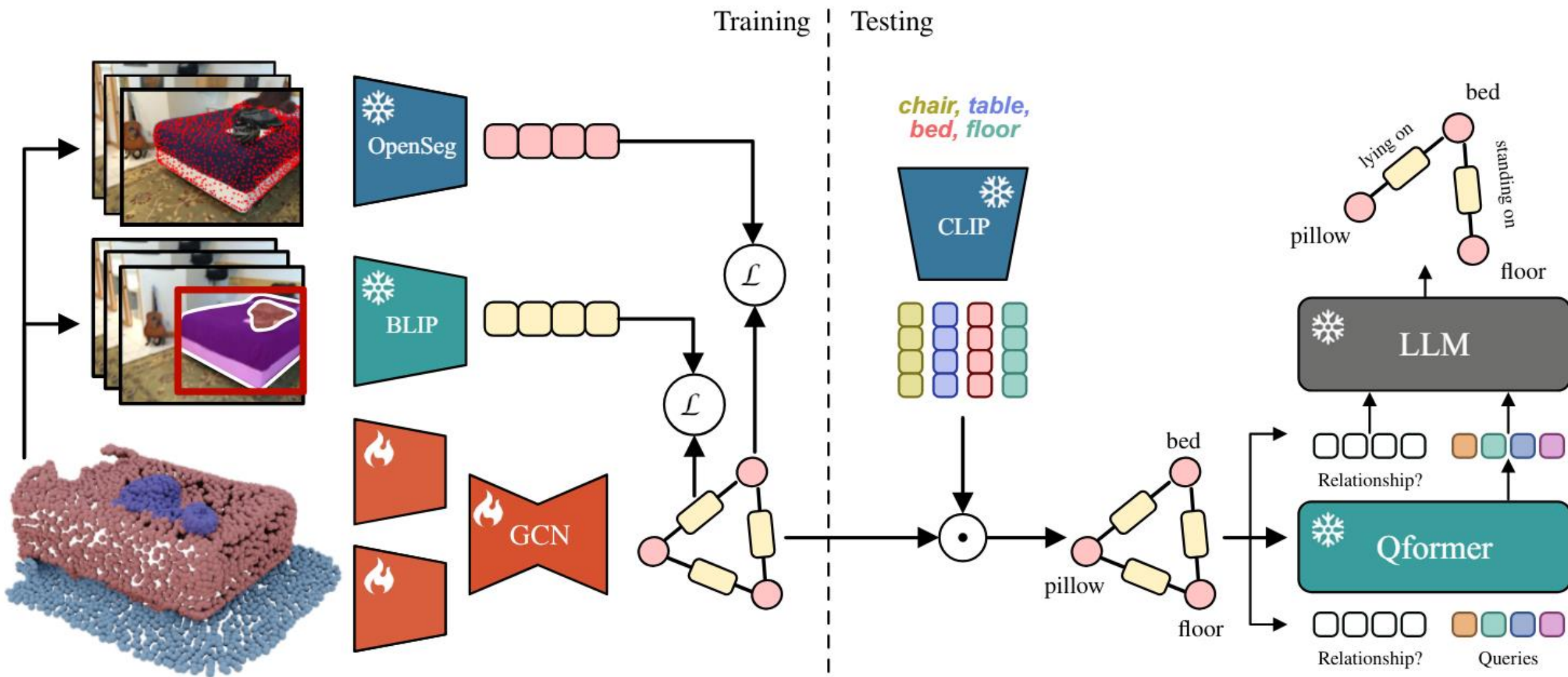


(2) Create dataset classifier from label text



(3) Use for zero-shot prediction

Open 3DSG



Open 3DSG

