

Auto-Encoding Scene Graphs with Image Reconstruction

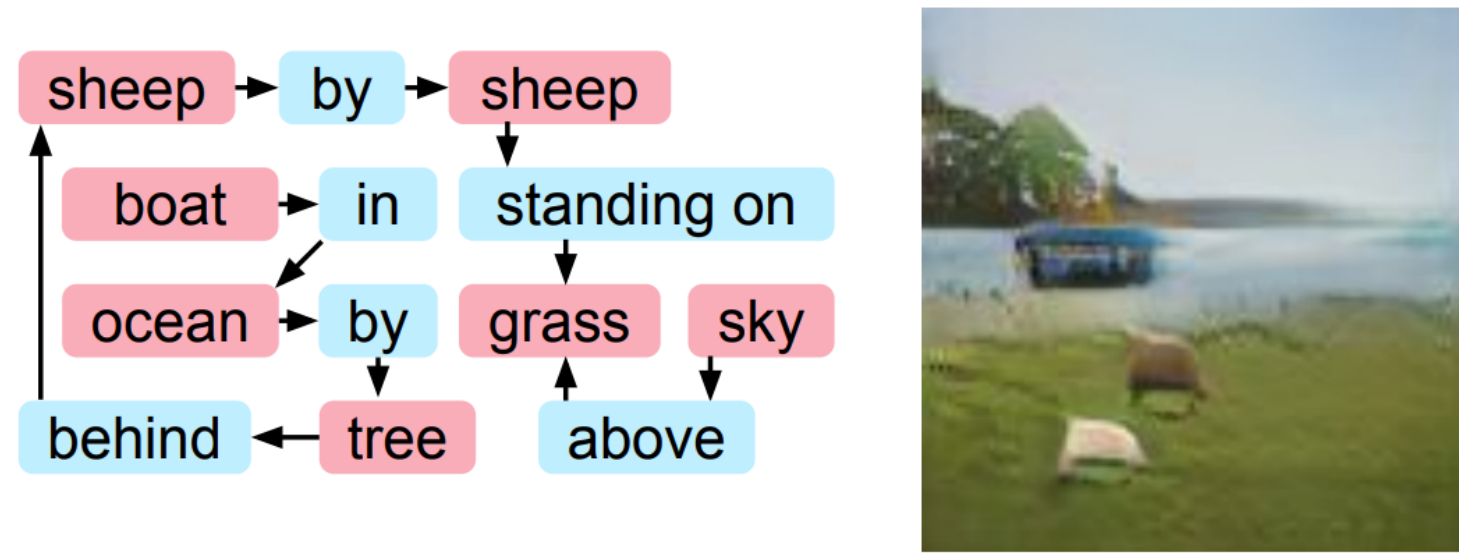
Kevin Tan

kevintan@cs.stanford.edu

Stanford
Computer Science

Motivation

- Goal: Learn to generate scene graphs from images and reconstruction by conditional image synthesis
- To fully understand the visual world requires not only recognizing individual objects, but also inferring the relationships and interactions between them
- Modern deep generative models excel in producing high-quality samples, but fall short in capturing compositional, object-centric semantic structure in visual scenes



Contributions

- We propose a method for scene graph generation in an end-to-end trainable fashion via *image reconstruction supervision*
- Auto-Encoding Scene Graphs consists of two parts:
 - (1) Encoder: given an input image, proposes object regions by a region proposal network, prunes connections with relational proposal network, and aggregates contextual information via graph convolution, and outputs a scene graph
 - (2) Decoder: Given a scene layout, feed noise into a cascaded refinement network to perform conditional image synthesis

Problem Statement

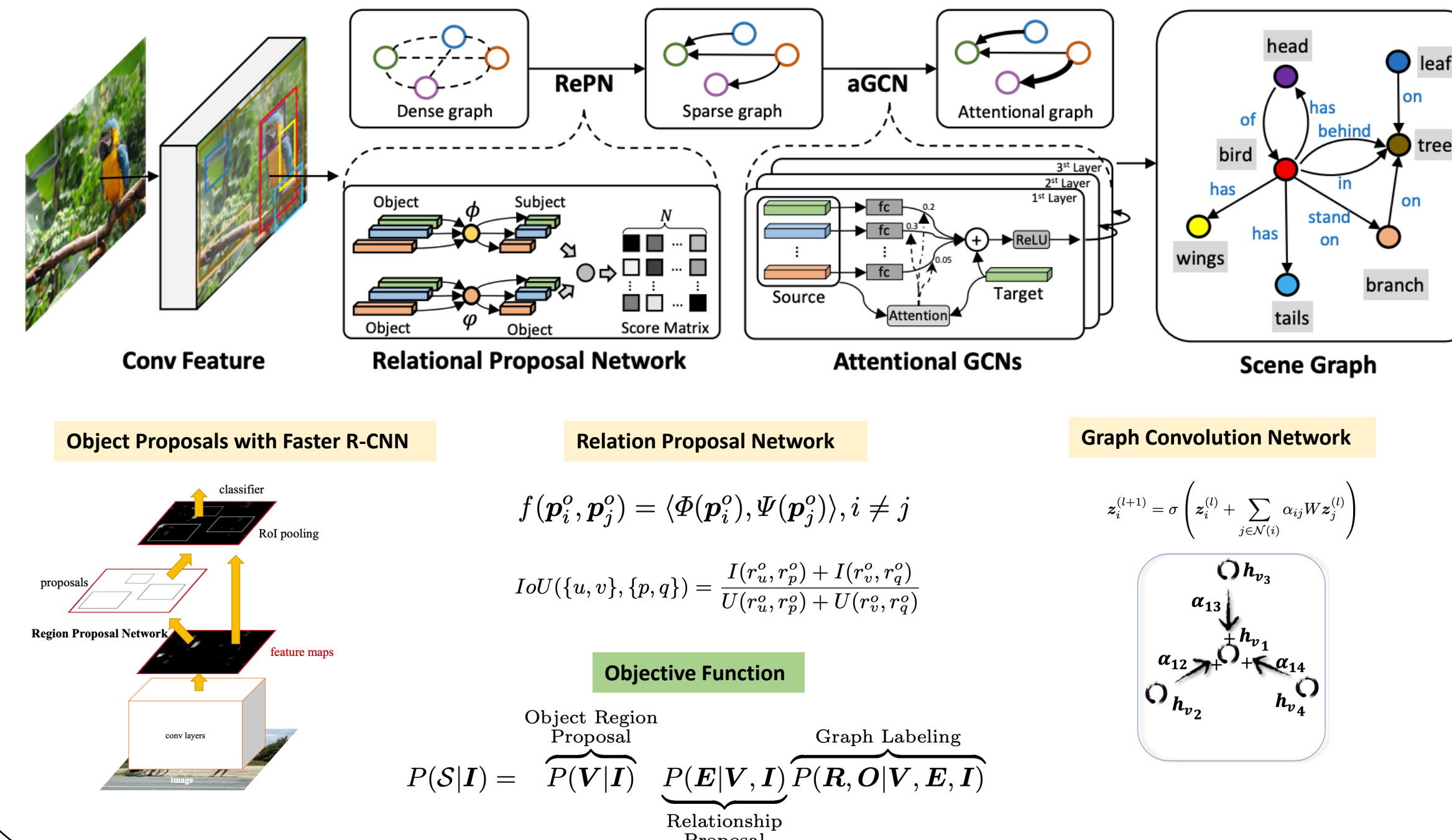
More formally, a *scene graph* can be defined by a 3-tuple set $G = \{B, O, R\}$:

- $B = \{b_1, b_2, \dots, b_n\}$ is the region candidate set, with elements $b_i \in \mathbb{R}^4$ denoting the bounding box of the i^{th} region
- $O = \{o_1, o_2, \dots, o_n\}$ is the object set, with element $o_i \in \mathbb{N}$ denoting the corresponding class label of region b_i
- $R = \{r_1, r_2, \dots, r_m\}$ of pairwise relationships between those objects, where r_k denotes a triplet of a start node $(b_i, o_i) \in B \times O$, an end node $(b_j, o_j) \in B \times O$, and a relationship label $x_{i \rightarrow j} \in \mathcal{R}$, where \mathcal{R} is the set of all possible predicate types.

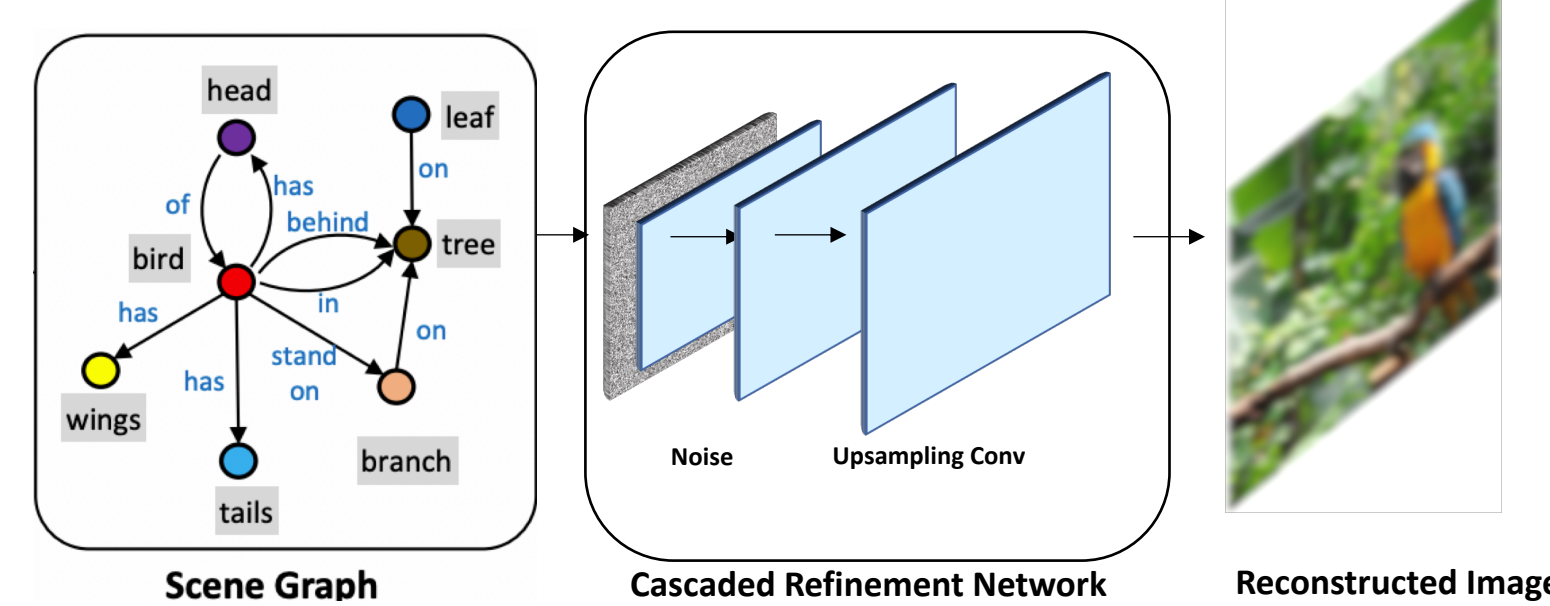
Given an image I , the goal is to decompose the probability distribution of the scene graph $P(G | I)$ into three components, as demonstrated previously by [17]:

$$Pr(G | I) = Pr(B | I)Pr(O | B, I)Pr(R | O, B, I) \quad (1)$$

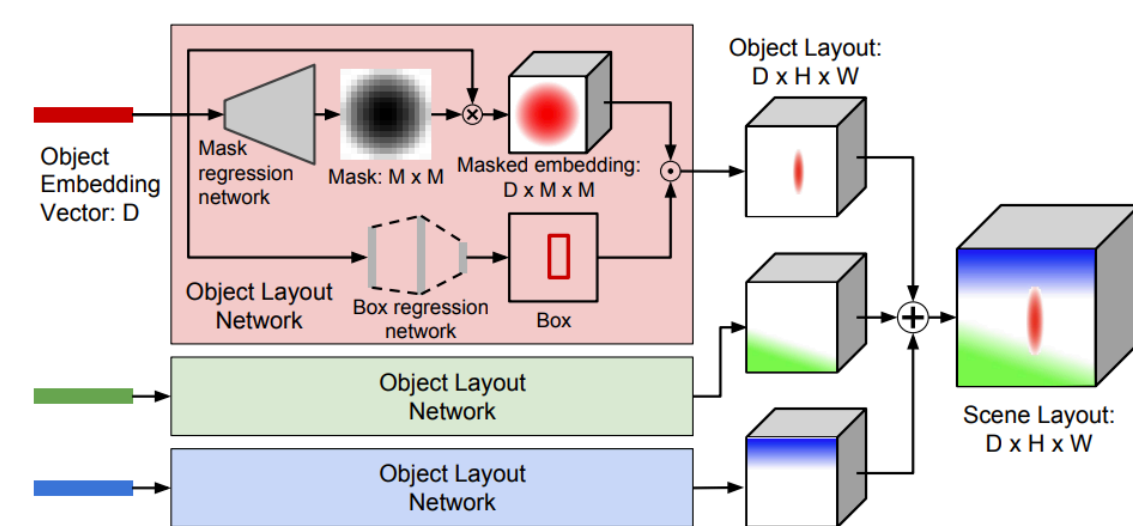
Scene Graph Generation



Conditional Image Synthesis



Scene Layout



Procedure for Image Generation

Function: Gen_{o2i}

Input: Real image I , objects (GT / predicted).

- Object Layout Generation: $O_i^{\text{layout}} \leftarrow \{o_i, r_i\}$
- Scene Layout Generation: $S^{\text{layout}} = \sum_i O_i^{\text{layout}}$
- Image Reconstruction: $\hat{I} = G_i(z, S^{\text{layout}})$
- Update image generator G_i parameters using (17).
- Update image discriminator D_i parameters using (16).

Image-level Supervision

$$\mathcal{L}_{D_i} = \mathbb{E}_{I \sim p_{\text{real}}} [\log D_i(I)]$$

$$\mathcal{L}_{G_i} = \mathbb{E}_{\hat{I} \sim p_G} [\log(1 - D_i(\hat{I})) + \lambda_p \mathcal{L}_{\text{pixel}}]$$

$$\|I - \hat{I}\|_1$$

Experiments

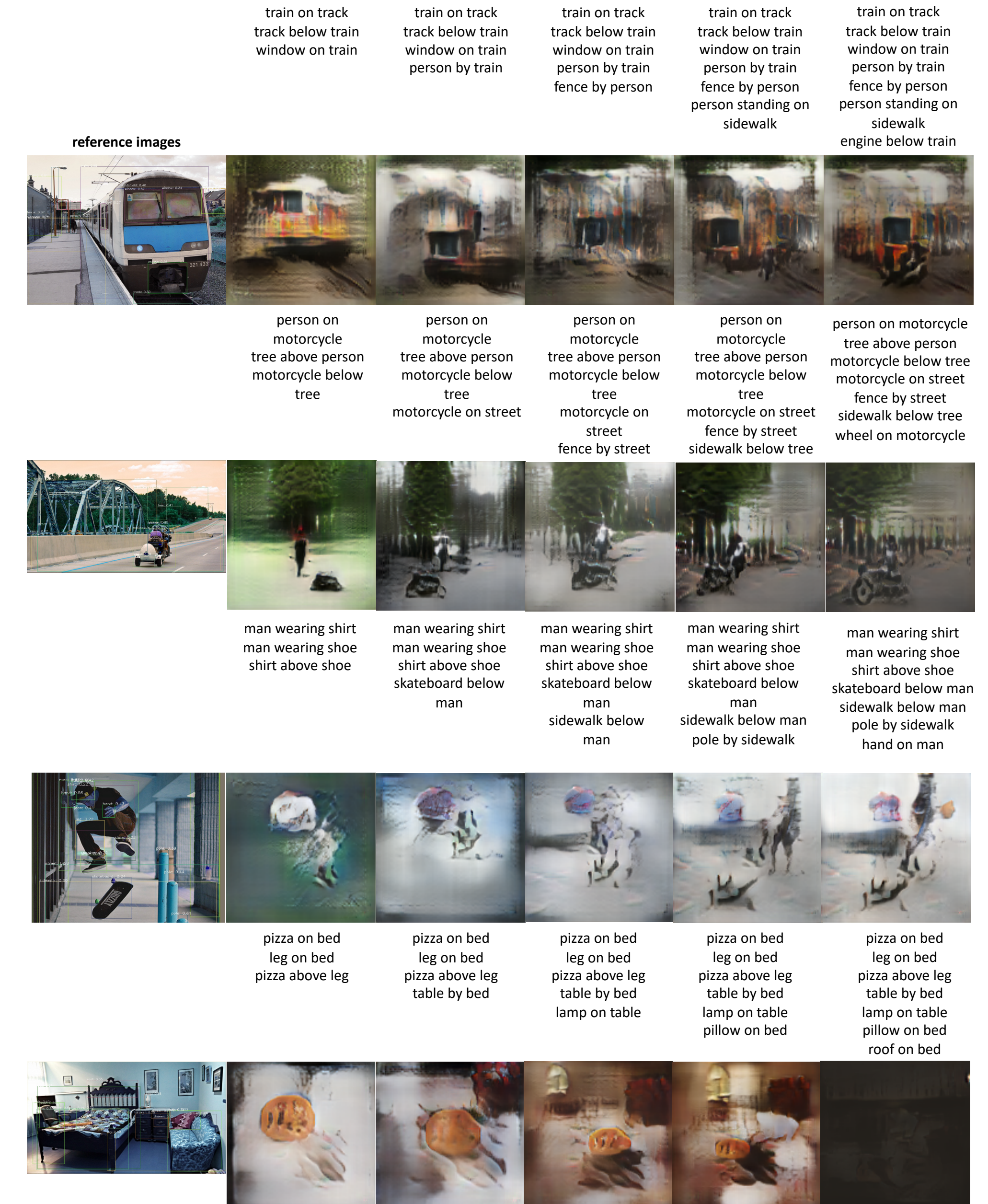


Figure: Qualitative results for conditional image synthesis on Visual Genome (VG). The first column depicts a reference image sampled from VG with the predicted bounding boxes. Each row begins with a simple scene graph and gradually adds more objects and relationships to gain complexity. Generated images demonstrate contextual knowledge and respect object relationships (e.g. pizza on bed).

Method	SGGen+		SGGen		PhrCls		PredCls	
	R@50	R@100	R@50	R@100	R@50	R@100	R@50	R@100
IMP	25.6	27.7	6.4	8.0	20.6	22.4	40.8	45.2
MSDN	25.8	28.2	7.0	9.1	27.6	29.9	53.2	57.9
NM-Freq	26.4	27.8	6.9	9.1	23.8	27.2	41.8	48.8
Graph R-CNN	28.5	35.9	11.4	13.7	29.6	31.6	54.2	59.1

Table: Comparisons on Visual Genome test set. We use Graph-RCNN for the scene graph generator outperforming the baseline of Iterative Message Passing (IMP). MSDN and NM-Freq refer to Multi-level Scene Description Network and Neural Motifs Frequency Prior, respectively.