# Specifying Object Attributes and Relations in Interactive Scene Generation

Oron Ashual
Tel Aviv University
oron.ashual@gmail.com

Lior Wolf
Tel Aviv University and Facebook AI Research
wolf@cs.tau.ac.il, wolf@fb.com

## Abstract

*We introduce a method for the generation of images from an input scene graph. The method separates between a layout embedding and an appearance embedding. The dual embedding leads to generated images that better match the scene graph, have higher visual quality, and support more complex scene graphs. In addition, the embedding scheme supports multiple and diverse output images per scene graph, which can be further controlled by the user. We demonstrate two modes of per-object control: (i) importing elements from other images, and (ii) navigation in the object space, by selecting an appearance archetype.*

*Our code is publicly available at* https://www.github.com/ashual/scene_generation.

## 1. Introduction

David Marr has defined vision as the process of discovering from images what is present in the world, and where it is [15]. The combination of *what* and *where* captures the essence of an image at the semantic level and therefore, also plays a crucial role when defining the desired output of image synthesis tools.

In this work, we employ scene graphs with per-object location and appearance attributes as an accessible and easy-to-manipulate way for users to express their intentions, see Fig. 1. The *what* aspect is captured hierarchically: objects are defined as belonging to a certain class (horse, tree, boat, etc.) and as having certain appearance attributes. These attributes can be (i) selected from a predefined set obtained by clustering previously seen attributes, or (ii) copied from a sample image. The *where* aspect, is captured by what is often called a scene graph, i.e., a graph where the scene objects are denoted as nodes, and their relative position, such as "above" or "left of", are represented as edge types.

Our method employs a dual encoding for each object in the image. The first part encodes the object's placement and captures a relative position and other global image features, as they relate to the specific object. It is generated based on the scene graph, by employing a graph convolution net-

work, followed by the concatenation of a random vector $z$. The second part encodes the appearance of the object and can be replaced, e.g., by importing it from the same object as it appears in another image, without directly changing the other objects in the image. This copying of objects between images is done in a semantic way, and not at the pixel level.

In the scene graph that we employ, each node is equipped with three types of information: (i) the type of object, encoded as a vector of a fixed dimension, (ii) the location attributes of the objects, which denote the approximate location in the generated image, using a coarse $5 \times 5$ grid and its size, discretized to ten values, and (iii) the appearance embedding mentioned above. The edges denote relations: "right of", "left of", "above", "below", "surrounding", and "inside". The method is implemented within a convenient user interface, which supports a dynamic placement of objects and the creation of a scene graph. The edge relations are inferred automatically, given the relative position of the objects. This eliminates the need for mostly unnecessary user intervention. Rendering is done in real time, supporting the creation of novel scenes in an interactive way, see Fig. 1.

The neural network that we employ has multiple subparts, as can be seen in Fig. 2: (i) A graph convolutional network that converts the input scene graph to a per-object embedding to their location. (ii) A CNN that converts the location embedding of each object to an object's mask. (iii) A parallel network that converts the location embedding to a bounding box location, where the object mask is placed. (iv) An appearance embedding CNN that converts image information into an embedding vector. This process is done off-line and when creating a new image, the vectors can be imported from other images, or selected from a set of archetypes. (v) A multiplexer that combines the object masks and the appearance embedding information, to create a one multidimensional tensor, where different groups of layers denote different objects. (vi) An encoder-decoder residual network that creates the output image.

Our method is related to the recent work of [9], who create images based on scene graphs. Their method also uses a graph convolutional network to obtain masks, a mul-
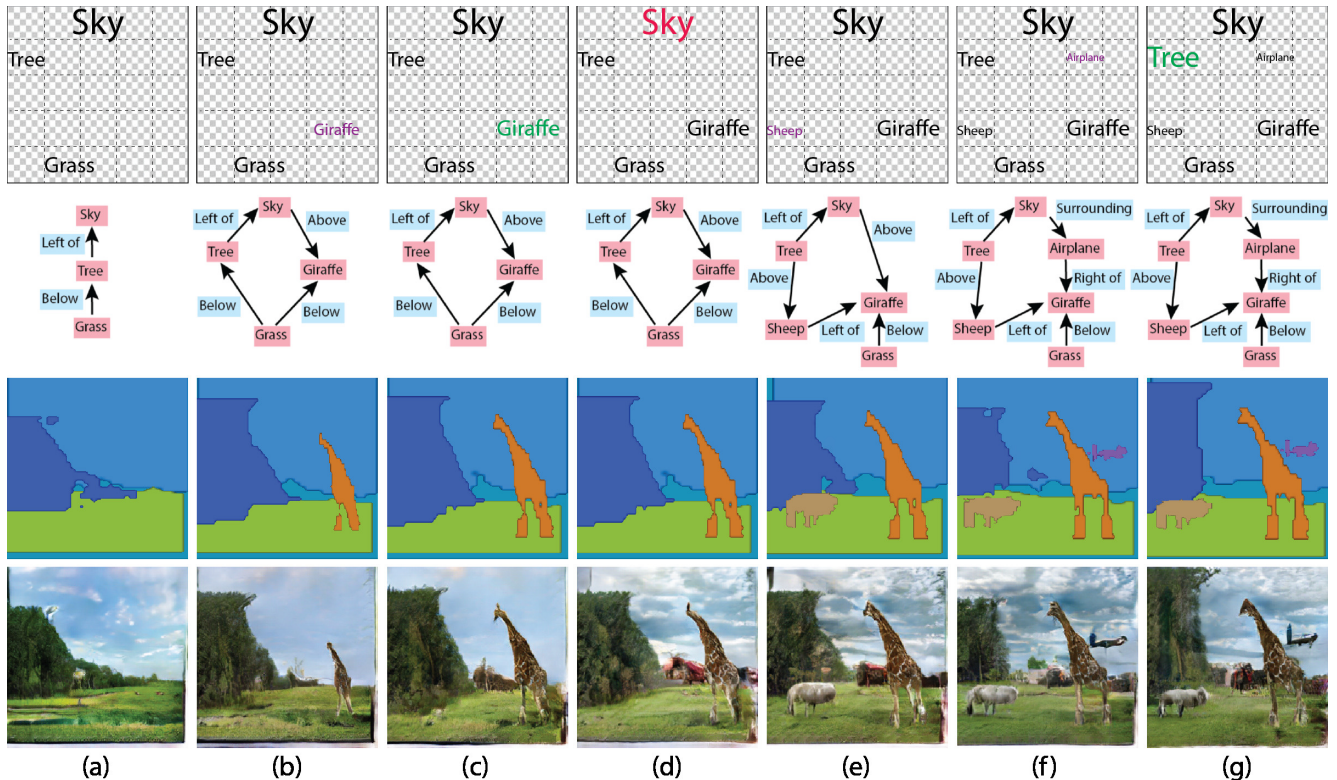
Figure 1. An example of the image creation process. (top row) the schematic illustration panel of the user interface, in which the user arranges the desired objects. (2nd row) the scene graph that is inferred automatically based on this layout. (3rd row) the layout that is created from the scene graph. (bottom row) the generated image. Legend for the GUI colors in the top row: purple – adding an object, green – resizing it, red – replacing its appearance. (a) A simple layout with a sky object, a tree and a grass object. All object appearances are initialized to a random archetype appearance. (b) A giraffe is added. (c) The giraffe is enlarged. (d) The appearance of the sky is changed to a different archetype. (e) A small sheep is added. (f) An airplane is added. (g) The tree is enlarged.

tiplexer that combines the layout information and a subsequent encoder-decoder architecture for obtaining the final image. There are, however, important differences: (i) by separating the layout embedding from the appearance embedding, we allow for much more control and freedom to the object selection mechanism, (ii) by adding the location attributes as input, we allow for an intuitive and more direct user control, (iii) the architecture we employ enables better quality and higher resolution outputs, (iv) by adding stochasticity before the masks are created, we are able to generate multiple results per scene graph, (v) this effect is amplified by the ability of the users to manipulate the resulting image, by changing the properties of each individual object, (vi) we introduce a mask discriminator, which plays a crucial role in generating plausible masks, (vii) another novel discriminator captures the appearance encoding in a counterfactual way, and (viii) we introduce feature matching based on the discriminator network and (ix) a perceptual loss term to better capture the appearance of an object, even if the pose or shape of that object has changed.

## 2. Previous Work

Image generation techniques based on GANs [3] are constantly improving in resolution, visual quality, the diversity of generated images, and the ability to cover the entire visual domain presented during training. In this work, we address conditional image generation, i.e., the creation of images that match a specific input. Earlier work in conditional image generation includes class based image generation [16], which generates an image that matches a given textual description [18, 26]. In many cases, the conditioning signal is a source image, in which case the problem is often referred to as image translation. Pix2pix [7] is a fully supervised method that requires pairs of matching samples from the two domains. The Pix2pixHD architecture that was recently presented by [25] is highly influential and many recent video or image mapping works employ elements of it, including our work.

Image generation based on scene graphs was recently presented in [9]. A scene graph representation is often used for retrieval based on text [10, 17], and a few datasets include this information, e.g., COCO-stuff [2] and the visual
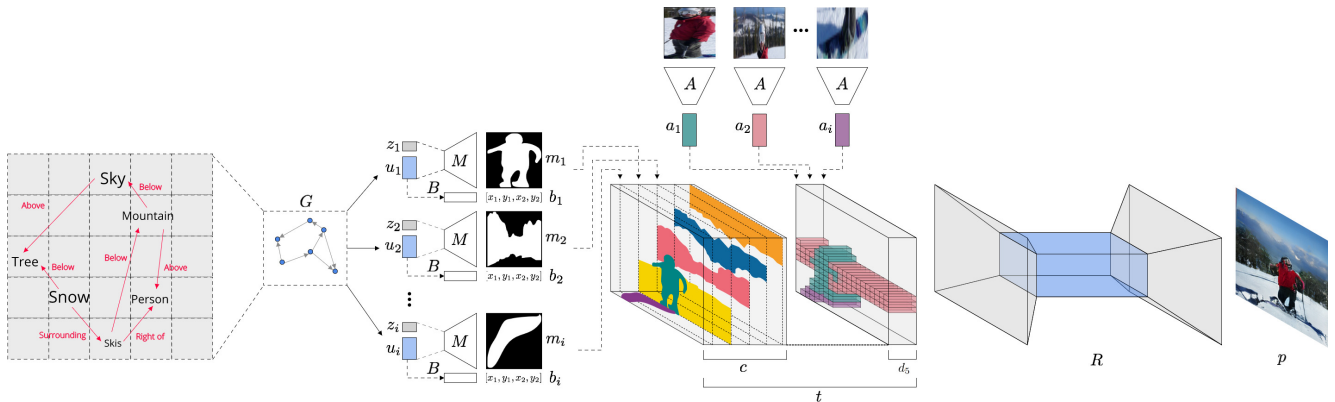
Figure 2. The architecture of our composite network, including the subnetworks $G, M, B, A, R$, and the process of creating the layout tensor $t$. The scene graph is passed to the network $G$ to create the layout embedding $u_i$ of each object. The bounding box $b_i$ is created from this embedding, using network $B$. A random vector $z_i$ is concatenated to $u_i$, and the network $M$ computes the mask $m_i$. The appearance information, as encoded by the network $A$, is then added to create the tensor $t$ with $c + d_5$ channels, $c$ being the number of classes. The autoencoder $R$ generates the final image $p$ from this tensor.

genome [12]. Also related is the synthesis of images from a given input layout of bounding boxes (and not one that is inferred by a network from a scene graph), which was very recently studied by [28] for small 64x64 images. In another line of work, images are generated to match input sentences, without constructing the scene graph as an intermediate representation [18, 26, 6].

Recently, an interactive tool was introduced based on the novel notion of GAN dissection [1]. This tool allows for the manipulation of well-localized neurons that control the occurrence of specific objects across the image using a drawing interface. By either adding or reducing the activations of these neurons, objects can be added, expanded, reduced or removed. The manipulation that we offer here is both more semantic (less related to specific locations and more to spatial relations between the objects), and more precise, in the sense that we provide full control over the exact instance of the object and not just of the desired class.

## 3. Method

Each object $i$ in the input scene graph is associated with a single node $n_i = [o_i, l_i]$, where $o_i \in \mathbb{R}^{d1}$ is a learned encoding of the object class and $l_i \in \{0, 1\}^{d_2+d_3}$ is a location vector. The object class embedding $o_i$ is one of $c$ possible embedding vectors, $c$ being the number of classes, and $o_i$ is set according to the class of object $i$, denoted $c_i$. The embedding size $d_1$ is set arbitrarily to 128. The first $d_2 = 25$ bits of $l_i$ denote a coarse image location using a $5 \times 5$ grid, and the rest denotes the size of the object, using a scale of $d_3 = 10$ values. The edge information $e_{ij} \in \mathbb{R}^{d_1}$ exists for a subset of the possible pairs of nodes, and encodes, using a learned embedding, the relations between the nodes. In other words, the values of $e_{ij}$ are taken from a learned dictionary with six possible values, each associated with one type of pairwise relation.

The location of each generated object is given as a pseudo-binary mask $m_i$ (output of a sigmoid) and a bounding box $b_i = [x_1, y_1, x_2, y_2]^\top \in [0, 1]^4$, which encodes the coordinates of the bounding box as a ratio of the image dimensions. The mask, but not the bounding box, is also determined by a per-object random vector $z_i \sim \mathbb{N}(0, 1)^{d_4}$ to create a variation in the generated masks, where $d_4 = 64$ was set arbitrarily, without testing other values.

The method employs multiple ways of embedding input information. The class identity and every inter-object relation, both taking a discrete value, are captured by embedding vectors of dimension $d_1$, which are learned as part of the end-to-end training. The object appearance $a_i \in \mathbb{R}^{d_5}$ of object $i$ seen during training, is obtained by applying a CNN $A$ to a (ground truth) cropped image $I'_i$ of that object, resized to a fixed resolution of $64 \times 64$. $d_5$ was set arbitrarily to 32 to reflect that it has less information than that of the entire object, which is embedded in $\mathbb{R}^{d_1}$.

The way in which the data flows through the subnetworks, as depicted in Fig. 2, is captured by the equations:

$$u_i = G(\{n_i\}, \{e_{ij}\}) \quad (1) \qquad a_i = A(I'_i) \quad (4)$$
$$m_i = M(u_i, z_i) \quad (2) \qquad t = T(\{c_i, m_i, b_i, a_i\}) \quad (5)$$
$$b_i = B(u_i) \quad (3) \qquad p = R(t) \quad (6)$$

where $G$ is the graph convolutional network [9, 20] that generates the per-object layout embedding, $M$ and $B$ are the networks that generate the object's mask and its bounding box, respectively, $T$ is the fixed (unlearned) function that maps the various object embeddings to a tensor $t$. Finally, $R$ is the encoder-decoder network that outputs an image $p \in \mathbb{R}^{H \times W \times 3}$ based on $t$. The exact architecture of each network is provided in the appendix.

The function $T$ constructs the tensor $t$ as a sum of per-

3

object tensors $t_i \in \mathbb{R}^{H \times W \times (d_5+c)}$, where $c$ is the number of objects. First, the mask $m_i$ is shifted and scaled, according to the bounding box $b_i$, resulting in a mask $m_i^{HW}$ of size $H \times W$. Then, a first tensor $t_i^1$ of size $H \times W \times d_5$ is formed as the tensor product of $m_i^{HW}$ and $a_i$. Similarly, a second tensor $t_i^2 \in \mathbb{R}^{H \times W \times c}$ is formed as the tensor product of $m_i^{HW}$ and the one hot vector of length $c$ encoding class $c_i$. The tensor $t_i$ is a concatenation of the two tensors $t_i^1$ and $t_i^2$ along the third dimension.

For performing adversarial training of the appearance embedding network $A$, we create two other tensors: $t'$ and $t''$. The first one is obtained by employing the ground truth bounding box $b_i'$ of object $i$ and the ground truth segmentation mask $m_i'$. The second one is obtained by incorporating the same ground truth bounding box and mask in a counterfactual way, by replacing $a_i$ with $a_k$, where $a_k$ is an appearance embedding of an object image $I_k'$ of a different object from the same class $c_i$, i.e., $a_k = A(I_k')$, $c_i = c_k$ and

$$t' = T(\{c_i, m_i', b_i', a_i\}) \tag{7}$$
$$t'' = T(\{c_i, m_i', b_i', a_k\}) \tag{8}$$

During training, in half of the training samples, the location and size information vectors $l_i$ are zeroed, in order to allow the network to generate layouts, even when this information is not available.

## 3.1. Training Loss Terms

The loss used to optimize the networks contains multiple terms, which is not surprising, given the need to train five networks (not including the adversarial discriminators mentioned below) and two vector embeddings ($o_i$ and $e_{ij}$).

$$L = \mathcal{L}_{\text{Rec}} + \lambda_1 \mathcal{L}_{\text{box}} + \lambda_2 \mathcal{L}_{\text{perceptual}} + \lambda_3 \mathcal{L}_{\text{D-mask}} + \lambda_4 \mathcal{L}_{\text{D-image}}$$
$$+ \lambda_5 \mathcal{L}_{\text{D-object}} + \lambda_6 \mathcal{L}_{\text{FM-mask}} + \lambda_7 \mathcal{L}_{\text{FM-image}} \tag{9}$$

where in our experiments we set $\lambda_1 = \lambda_2 = \lambda_6 = \lambda_7 = 10$, $\lambda_3 = \lambda_4 = 1$, $\lambda_5 = 0.1$.

The reconstruction loss $\mathcal{L}_{\text{Rec}}$ is the L1 difference between the reconstructed image $p$ and the ground truth training image. The box loss $\mathcal{L}_{\text{box}}$ is the MSE between the computed $b_i$ (summed over all objects) and the ground truth bounding box $b_i'$. Note that unlike [9], we do not employ a mask loss, since our mask contains a stochastic element (Eq. 2). The perceptual loss $\mathcal{L}_{\text{perceptual}} = \sum_{u \in U} \frac{1}{u} ||F^u(p) - F^u(p')||_1$ [8] compares the generated image with the ground truth training image $p'$, using the activations $F^u$ of the *VGG* network [21] at layer $u$ in a set of predefined layers $U$.

Our method employs three discriminators $D_{\text{mask}}$, $D_{\text{object}}$, and $D_{\text{image}}$. The mask discriminator employs a Least Squares GAN (*LS-GAN* [14]) and is conditioned on the object's class $c_i$. Recall that $m_i'$ is the real mask of object $i$

and $m_i$, the generated mask, which depends on a random variable $z_i$. The GAN loss associated with the mask discriminator is given by

$$\mathcal{L}_{D-mask} = [\log D_{\text{mask}}(m_i', c_i)] +$$
$$\mathbb{E}_{z \sim \mathcal{N}(0,1)^{64}} [log(1 - D_{\text{mask}}(M(u_i, z), c_i)] \tag{10}$$

For the purpose of training $D_{\text{mask}}$, we minimize $-\mathcal{L}_{D-mask}$. The second discriminator $D_{image}$, is used for training in an adversarial manner three networks $R$, $M$, and $A$. The loss of $\mathcal{L}_{\text{D-image}}$ is a compound loss that is given as

$$\mathcal{L}_{\text{D-image}} = \mathcal{L}_{\text{real}} - \mathcal{L}_{\text{fake-image}} - \mathcal{L}_{\text{fake-layout}} + \mathcal{L}_{\text{alt-appearance}}$$

where

$$\mathcal{L}_{\text{real}} = \log D_{\text{image}}(t', p') \tag{11}$$
$$\mathcal{L}_{\text{fake-image}} = \log(1 - D_{\text{image}}(t', p)) \tag{12}$$
$$\mathcal{L}_{\text{fake-layout}} = \log(1 - D_{\text{image}}(t, p')) \tag{13}$$
$$\mathcal{L}_{\text{alt-appearance}} = \log(1 - D_{\text{image}}(t'', p')) \tag{14}$$

The goal of the compound loss is to make sure that the generated image $p$, given a ground truth layout tensor $t'$ is indistinguishable from the real image $p'$, and that this is true, even if the layout tensor $t$ is based on estimated bounding boxes and masks (unlike $t'$). In addition, we would like the ground truth image to be a poor match for a counterfactual appearance vector, as given in $t''$.

Following [25] we use a multi-scale *LS-GAN* with two scales. In other words, $\mathcal{L}_{\text{D-image}}$ is computed at the full scale and at half scale (using two different discriminators), and both terms are summed up to obtain the actual $\mathcal{L}_{\text{D-image}}$.

The third discriminator, $D_{\text{object}}$, guarantees that the generated objects, one by one, look real. For this purpose, we crop $p$ using the bounding boxes $b_i$ to create object images $I_i$. Recall that $I_i'$ are ground truth crops of images, obtained from the ground truth image $p'$, using the ground truth bounding boxes $b'$.

$$\mathcal{L}_{\text{D-object}} = \sum_{i=1}^{n} \log D_{\text{object}}(I_i') - \log D_{\text{object}}(I_i) \tag{15}$$

$D_{\text{object}}$ maximizes this loss during training.

The mask feature matching loss $\mathcal{L}_{\text{FM-mask}}$ and the image feature matching loss $\mathcal{L}_{\text{FM-image}}$ are similar to the perceptual loss, i.e., they are based on the L1 difference in the activation. However, instead of the *VGG* loss, the discriminators are used, as in [19]. In these losses, all layers are used. $\mathcal{L}_{\text{FM-mask}}$ compares the activations of the generated mask $m_i$ and the real mask $m_i'$ (the discriminator $D_{\text{mask}}$ also takes the class $c_i$ as input). The other feature matching loss $\mathcal{L}_{\text{FM-image}}$ compares the activations of $D_{\text{image}}(t, p)$ with those of the ground truth layout tensor and image $D_{\text{image}}(t', p')$.
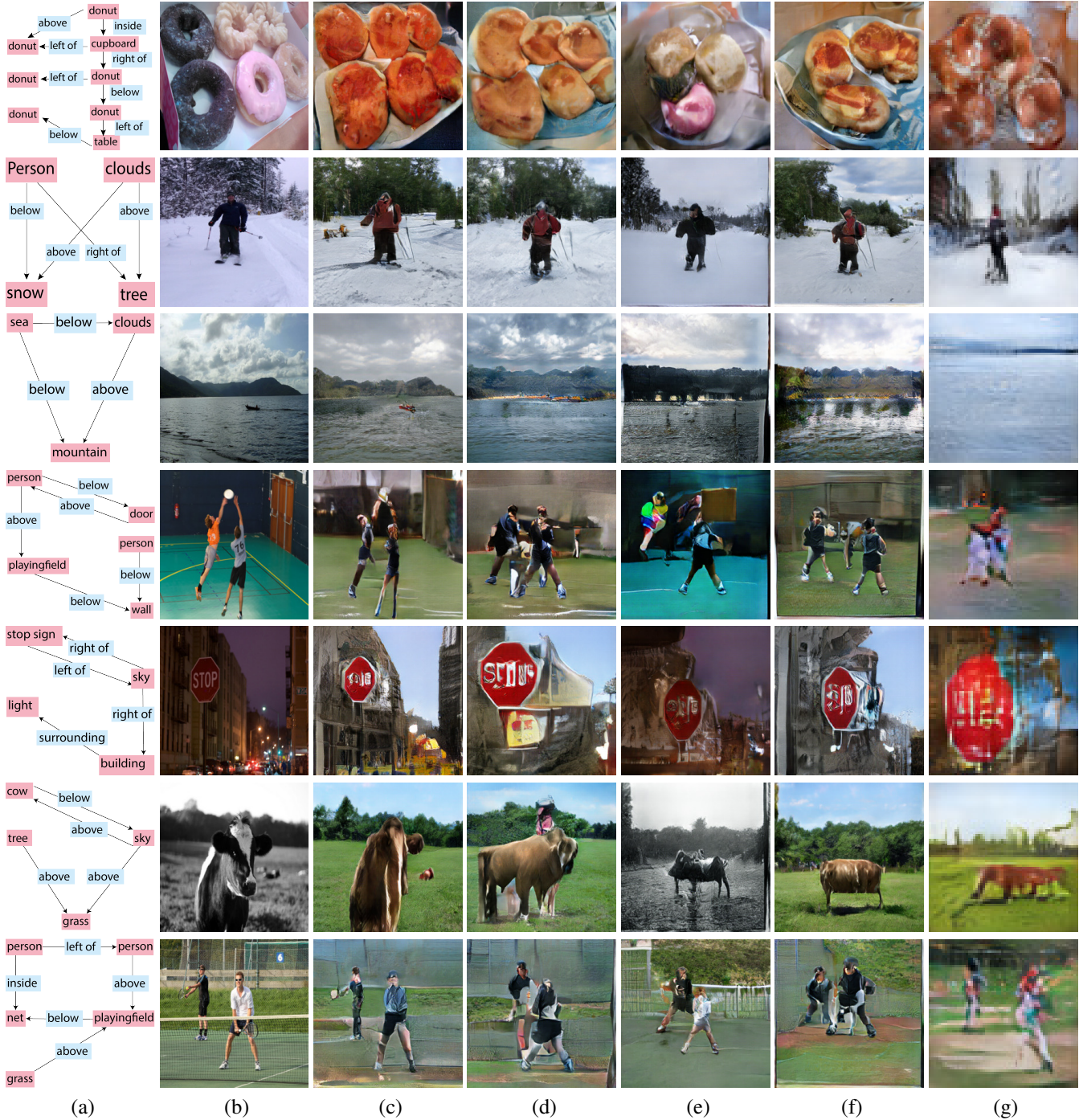
Figure 3. Image generation based on a given scene graph. Each row is a different example. (a) the scene graph, (b) the ground truth image, from which the layout was extracted, (c) our results when we used the ground truth layout of the image, similar to [28], (d) our method's results, where the appearance attributes present a random archetype and the location attributes coarsely describe the ground truth bounding box, (e) our results when we use the ground truth image to generate the appearance attributes, and the location attributes are zeroed $l_i = 0$, (f) our results where $l_i = 0$, and the appearance attributes are sampled from the archetypes, and (g) the results of [9].

## 3.2. Generating the archetypes

The GUI enables the user to select from preexisting object appearances, as well as copying the appearance vector $a_i$ from another image. The existing object appearances are given as 100 archetypes per object class. These are obtained, by applying the learned network $A$ to all objects in a
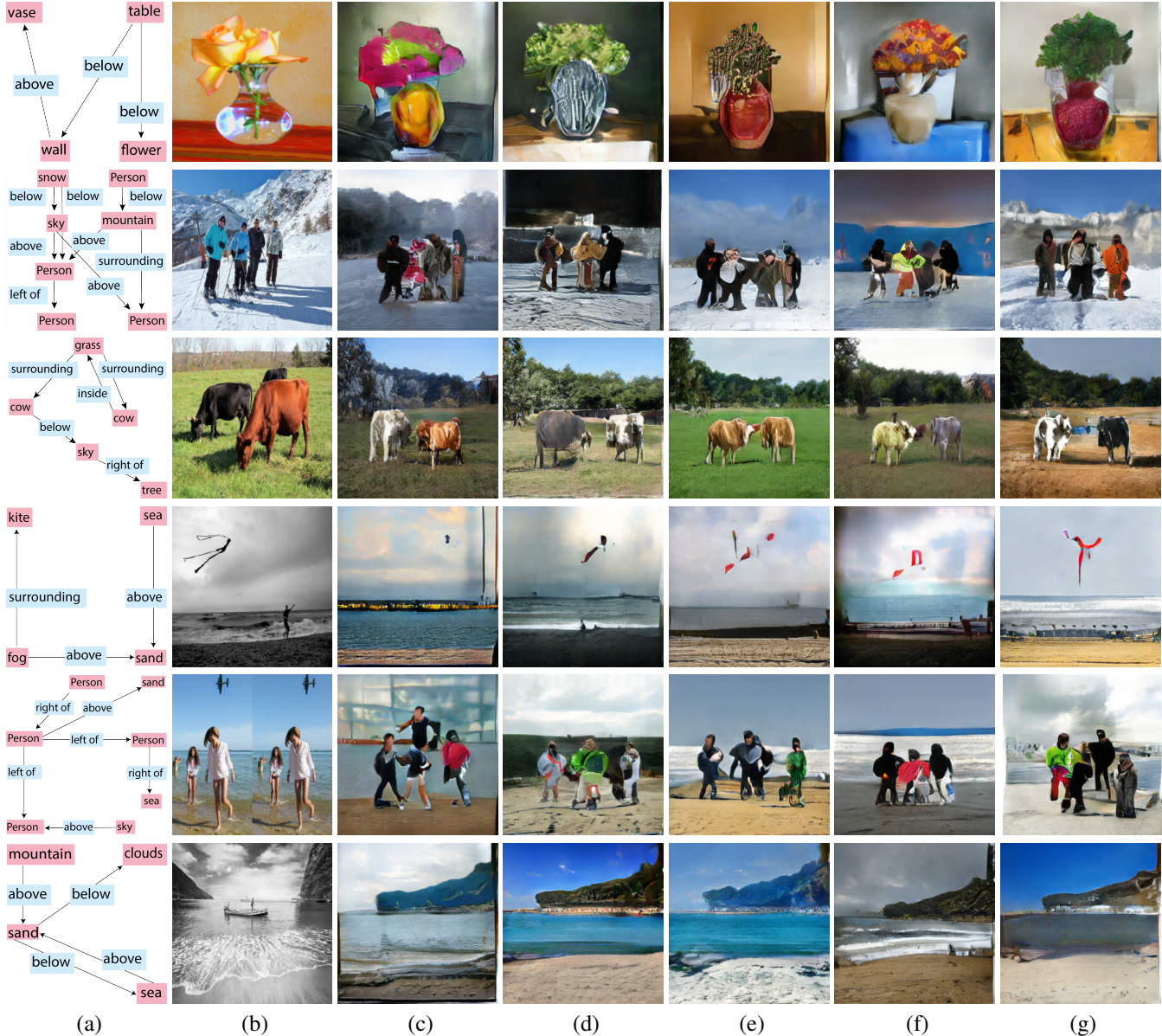
Figure 4. The diversity obtained when keeping the location attributes $l_i$ fixed at zero and sampling different appearance archetypes. (a) the scene graph, (b) the ground truth image, from which the layout was extracted, (c–g) generated images.

given class in the training set and employing k-means clustering, in order to obtain 100 class means.

In the GUI, the archetypes are presented linearly along a slider. The order along the slider is obtained by applying a 1-D t-SNE [24] embedding to the 100 archetypes.

### 3.3. Inferring the scene graph from the layout panel

The GUI lets the users place objects on a schematic layout, see Fig. 1. Each object is depicted as a string in one of ten different font sizes, in order to capture the size element of $l_i$. The location in the layout determines the $5 \times 5$ grid placement, which is encoded in the other part of $l_i$.

Note, however, that the locations and sizes are provided as indications of the structure of the graph layout and not as absolute locations (or scene layout). The generating network maintains freedom in the object placements to match the semantic properties of the objects in the scene.

The coarse placement by the user is more intuitive and less laborious than specifying a scene graph. To avoid adding unwanted work for the users, the edge labels are inferred, based on the relative position and size of the objects. An object which is directly to the left of another object, for example, is labeled "left of". The order in which objects are inserted to the layout determines the reference directing. If
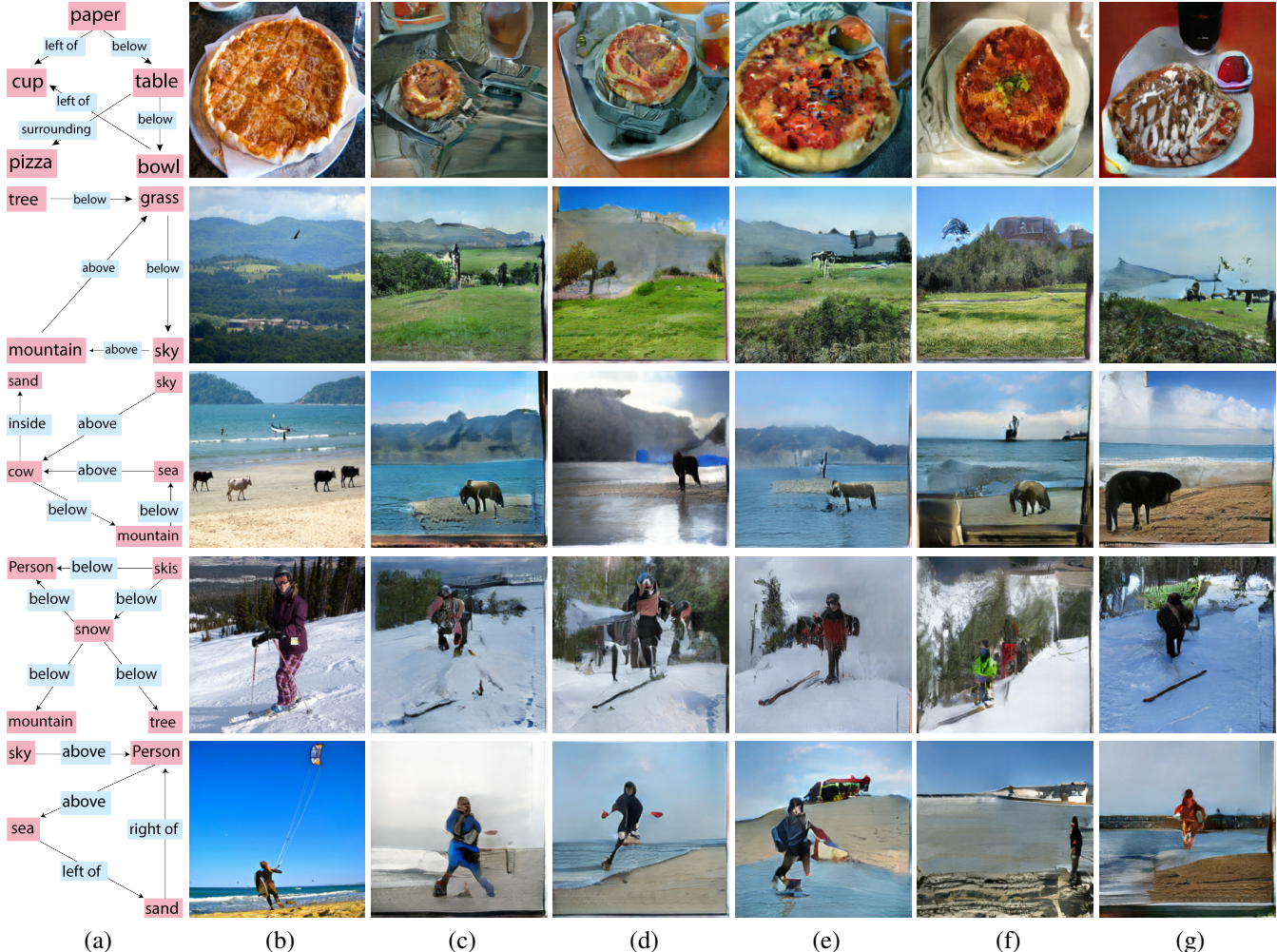
Figure 5. The diversity obtained when keeping the appearance vectors fixed and sampling from the location distribution. (a) the scene graph, (b) the ground truth image from which the layout was extracted, (c–g) generated images.

object $i$ is inserted before object $j$, then "$i$ is to the left of $j$" and not "$j$ is to the right of $i$". The inside and surrounding relations are determined similarly, by considering objects of different sizes, whose centers are nearby.

### 3.4. Training details

All networks are trained using *ADAM* [11] solver with $beta1 = 0.5$ for 1 million iterations. The learning rate was set to $1e^{-4}$ for all components except $\mathcal{L}_{\text{D-mask}}$, where we set it to a smaller learning rate of $1e^{-5}$. The different learning rates help us to stabilize the mask network. We use batch sizes of 32, 16, 4 in our $64 \times 64$, $128 \times 128$, $256 \times 256$ resolutions respectively. Notice that since each image contains up to 8 objects, each batch contains up to $8 \times 32 = 256$ different objects.

## 4. Experiments

We compare our results with the state of the art methods of [9] and [28], using various metrics from the literature. In addition, we perform an ablation analysis to study the relative contribution of various aspects of our method. Our experiments are conducted on the COCO-Stuff dataset [2], which, using the same split as the previous works, contains approximately 25,000 train images, 1000 validation, and 2000 test images.

We employ two modes of experiments: either using the ground truth (GT) layout or the inferred layout. The first mode is the only one suitable for the method of [29]. Whenever possible, we report the statistics reported in the previous work. Some of the statistics reported for [9] are computed by us, based on the published model. We report results for three resolutions $64^2$, $128^2$, and $256^2$. The literature reports numerical results only for the first resolution. While [9] presents visual results for 128x128, our attempts
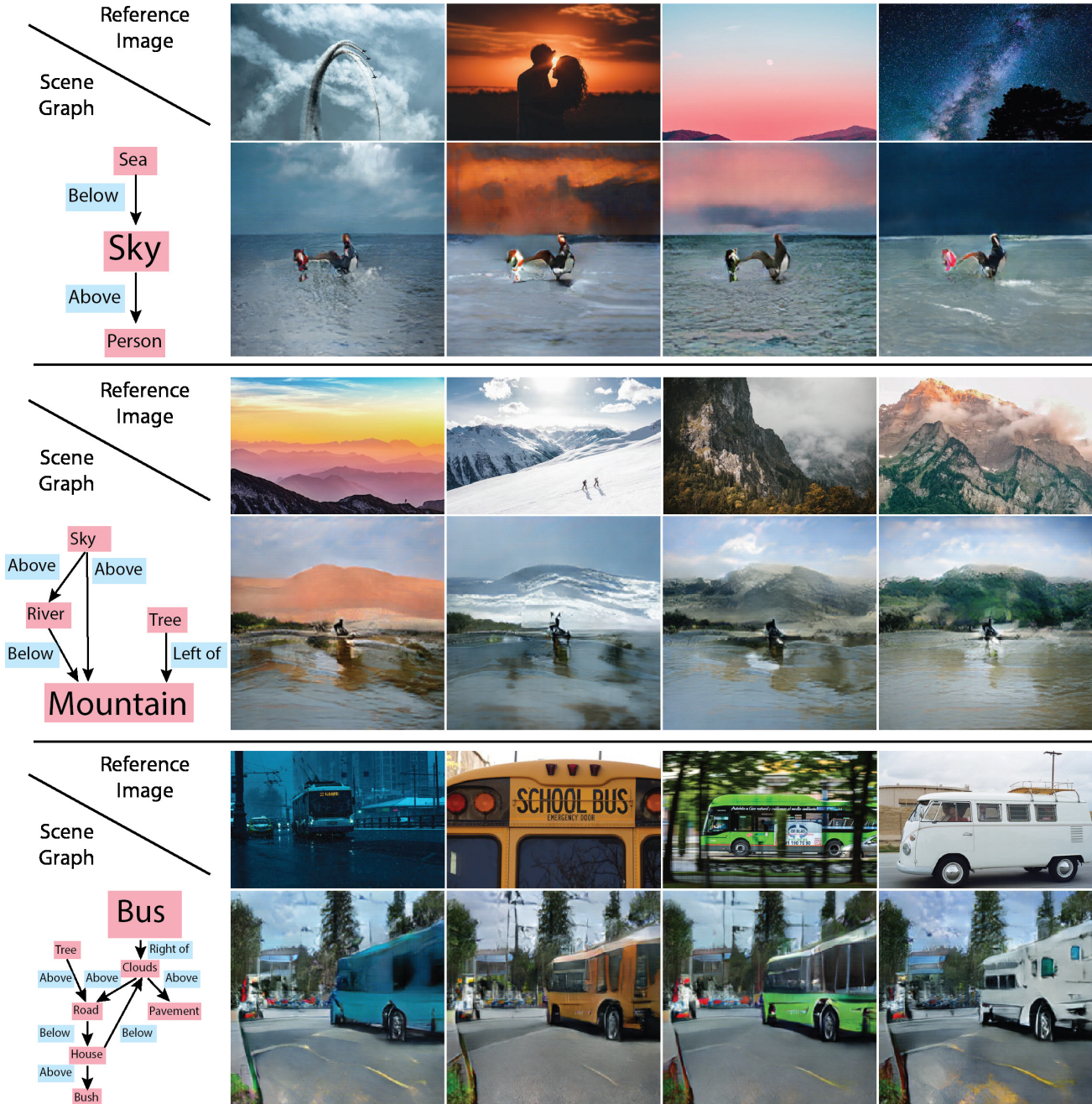
Figure 6. Duplicating an object's appearance in the generated image. Images are created based on the scene graph, such that the appearance is taken from one of five unrelated images. In this example, the sky's appearance is generated from the reference image, while all other objects use the same random appearance archetype.

to train their method using the published code on this resolution, resulted in sub-par performance, despite some effort. We, therefore, prefer not to provide these non-competitive baseline numbers. The code of [29] is not yet available.

We employ multiple acceptable literature evaluation metrics for evaluating the generated images. The *inception score* [19] measures both the quality of the generated images and their diversity. As has been done in previous works, a pre-trained inception network [22] is employed in order to obtain the network activations used to compute the score. Larger inception scores are better. The *FID* [5] measures the distance between the distribution of the generated images and that of the real test images, both modelled as a multivariate Gaussian. Lower *FID* scores are better.

8

| Reso-lution[b] | Method | Inception[a] | FID | Accu-racy |
|---|---|---|---|---|
| 64x64 | Real Images | $16.3 \pm 0.4$ | 0 | 54.5 |
| | [9] GT Layout | $7.3 \pm 0.1$ | 86.5 | 33.9 |
| | [28] GT Layout | $9.1 \pm 0.1$ | [c] | [d] |
| | Ours GT Layout | $10.3 \pm 0.1$ | 48.7 | 46.1 |
| | [9] | $6.7 \pm 0.1$ | 103.4 | 28.8 |
| | Ours | $7.9 \pm 0.2$ | 65.3 | 43.3 |
| 128x128 | Real Images | $24.2 \pm 0.9$ | 0 | 59.3 |
| | Ours GT Layout | $12.5 \pm 0.3$ | 59.5 | 44.6 |
| | Ours | $10.4 \pm 0.4$ | 75.4 | 42.8 |
| 256x256 | Real Images | $30.7 \pm 1.2$ | 0 | 62.4 |
| | Ours GT Layout | $16.4 \pm 0.7$ | 65.2 | 45.3 |
| | Ours | $14.5 \pm 0.7$ | 81.0 | 42.2 |

[a]The inception score of [9] for the complete pipeline is taken from their paper. The other scores are not reported there. The inception score for [28] is the one reported by the authors.

[b][9] and [28] report numerical results only for a resolution of 64x64.

[c]Not reported and cannot be computed due to lack of code/results.

[d]The accuracy reported by [28] is incompatible (different classifiers).

Table 1. A quantitative comparison using various image generation scores. In order to support a fair comparison, our model does not use location attributes and employs random appearance attributes.

Less common, but relevant to our task, is the *classification accuracy* score, used by [29]. A ResNet-101 model [4] is trained to classify the 171 objects available in the training datasets, after cropping and resizing them to a fixed size of 224x224 pixels. On the test image, we report the accuracy of this classifier applied to the object images that are generated, using the bounding box of the image's layout. A higher accuracy means that the method creates more realistic, or at least identifiable, objects.

We also report a *diversity score* [27], which is based on the perceptual similarity [10] between two images. This score is used to measure the distance between pairs of images that are generated given the same input. Ideally, the user would be able to obtain multiple, diverse, alternative outputs to choose from. Specifically, the activations of AlexNet [13] are used together with the LPIPS visual similarity metric [27]. A higher diversity score is better.

In addition, we also report three scores for evaluating the quality of the bounding boxes. The *IoU* score is the ratio between the area of the ground truth bounding box that is also covered by the generated bounding box (the intersection), and the area covered by either box (the union). We also report recall scores at two different thresholds. *R@0.5* measures the ratio of object bounding boxes with an IoU of at least 0.5, and similarly for *R@0.3*.

Tab. 1 compares our method with the baselines and the real test images using the inception, FID, and classification accuracy scores. We make sure not to use information that the baseline method of [9] is not using and use zero location attributes and appearance attributes that are randomly sam-

| Res | Method | Diversity |
|---|---|---|
| 64x64 | Johnson *et al*. [9] | $0.15 \pm 0.08$ |
| | Zhao *et al*. [28] GT layout | $0.15 \pm 0.06$ |
| | Ours fixed appearance attributes and zeroed location attributes | $0.23 \pm 0.01$ |
| | Ours zeroed location attributes | $0.35 \pm 0.01$ |
| | Ours fixed appearance attributes | $0.37 \pm 0.01$ |
| | Our full method | $0.43 \pm 0.07$ |
| 256x256 | Ours fixed appearance attributes and zeroed location attributes | $0.48 \pm 0.09$ |
| | Ours zeroed location attributes | $0.61 \pm 0.07$ |
| | Ours fixed appearance attributes | $0.62 \pm 0.05$ |
| | Our full method | $0.67 \pm 0.05$ |

Table 2. The diversity score of [27]. The results of [9] are computed by us and are considerably higher than those reported for the same method by [28]. The results of [28] are from their paper.

| | IoU | R@0.5 | R@0.3 |
|---|---|---|---|
| Johnson *et al*. [9][a] | 0.37 | 0.32 | 0.52 |
| Ours (w/o location attributes) | 0.41 | 0.37 | 0.62 |
| Ours (w/ location attributes) | 0.61 | 0.66 | 0.86 |

[a]Taken from the paper itself

Table 3. Comparison of predicted bounding boxes

| User Study | [9] | Ours |
|---|---|---|
| More realistic output | 16.7% | 83.3% |
| Better adherence to scene graph | 19.3% | 80.7% |
| Ratio of observed objects among all COCO objects | 27.31% | 45.38% |
| Ratio of observed objects among all COCO stuff | 46.49% | 65.23% |

Table 4. User study results

| Model | Inception | FID |
|---|---|---|
| Full method | $10.4 \pm 0.4$ | 75.4 |
| No $\mathcal{L}_{\text{perceptual}}$ | $6.2 \pm 0.1$ | 125.1 |
| No $\mathcal{L}_{\text{D-mask}}$ | $5.2 \pm 0.1$ | 183.6 |
| No $\mathcal{L}_{\text{D-image}}$ | $7.4 \pm 0.2$ | 122.5 |
| No $\mathcal{L}_{\text{D-object}}$ | $8.7 \pm 0.1$ | 94.5 |
| Using $D_{\text{image}}$ of [9] | $8.1 \pm 0.3$ | 114.2 |

Table 5. Ablation Study

pled (see 3.2). [29] employs bounding boxes and not masks. However, we follow the same comparison (to masked based methods) given in their paper.

As can be seen, our method obtains a significant lead in all these scores over the baseline methods, whenever such a comparison can be made. This is true both when the ground truth layout is used and when the layout is generated. As expected, the ground truth layout obtains better scores.

Sample results of our 256x256 model are shown in Fig. 3, using test images from the COCO-stuff datasets. Each row presents the scene layout, the ground truth image

from which the layout was extracted, our method's results, where the object attributes present a random archetype and the location attributes are zeroed ($l_i = 0$), our results when using the ground truth layout of the image (including masks and bounding boxes), our results where the appearance attributes of each object are copied from the ground truth image and the location vectors are zero, and our results where the location attributes coarsely describe the objects' locations and the appearance attributes are randomly selected from the archetypes. In addition, we present the result of the baseline method of [9] at the 64x64 resolution for which a model was published.

As can be seen, our model produces realistic results across all settings, which are more pleasing than the baseline method. Using ground truth location and appearance attributes, the resulting image better matches the test image.

Tab. 2 reports the diversity of our method in comparison to the two baseline methods. The source of stochasticity we employ (the random vector $z_i$ used in Eq. 2) produces a higher diversity than the two baseline methods (which also include a random element), even when not changing the location vector $l_1$ or appearance attributes $a_i$. Varying either one of these factors adds a sizable amount of diversity. In the experiments of the table, the location attributes, when varied, are sampled using per-class Gaussian distribution that fit to the location vectors of the training set images.

Fig. 4 presents samples obtained when sampling the appearance attributes. In each case, for all $i$, $l_i = 0$ and the object's appearance embedding $a_i$ is sampled uniformly between the archetypes. This results in a considerable visual diversity. Fig. 5 presents results in which the appearance is fixed to the mean appearance vector for all objects of that class and the location attribute vectors $l_i$ are sampled from the Gaussian distributions mentioned above. In almost all cases, the generated images are visually pleasing. In some cases, the location attributes sampled are not compatible with a realistic image. Note, however, that in our method, the default value for $l_i$ is zero and not a random vector.

Tab. 3 presents a comparison with the method of [9], regarding the placement accuracy of the bounding boxes. Even when not using the location attribute vectors $l_i$, our bounding box placement better matches the test images. As expected, adding the location vectors improves the results.

The ability of our method to copy the appearance of an existing image object is demonstrated in Fig. 6. In this example, we generate the same test scene graph, while varying a single object in accordance with five different options extracted from images unseen during training. Despite the variability of the appearance that is presented in the five sources, the generated images mostly maintain their visual quality. These results are presented at a resolution of 256x256, which is the default resolution for our GUI. At this resolution, the system processes a graph in 16.3ms.

**User study** Following [9], we perform a user study to compare with the baseline method the realism of the generated image, the adherence to the scene graph, as well as to verify that the objects in the scene graph appear in the output image. The user study involved $n = 20$ computer graphics and computer vision students. Each student was shown the output images for 30 random test scene-graphs from the COCO-stuff dataset and was asked to select the preferable method, according to two criteria: "which image is more realistic" and "which image better reflects the scene graph". In addition, the list of objects in the scene graph was presented, and the users were asked to count the number of objects that appear in each of the images. The two images, one for the method of [9] and one for our method, were presented in a random order. To allow for a fair comparison, the appearance archetypes were selected at random, the location vectors were set to zero for all objects, and we have used images from our $64 \times 64$ resolution model. The results, listed in Tab. 4, show that our method significantly outperforms the baseline method in all aspects tested.

**Ablation analysis** The relative importance of the various losses is approximated, by removing it from the method and training the 128x128 model. For this study, we use both the inception and the FID scores. The results are reported in Tab. 5. As can be seen, removing each of the losses results in a noticeable degradation. Removing the perceptual loss is extremely detrimental. Out of the three discriminators, removing the mask discriminator is the most damaging, since, due to the random component $z_i$, we do not have a direct loss on the mask. Finally, replacing our image discriminator with the one in [9], results in some loss of accuracy.

## 5. Conclusion

We present an image generation tool in which the input consists of a scene graph with the potential addition of location information. Each object is associated both with a location embedding and with an appearance embedding. The latter can be extracted from another image, allowing for a duplication of existing objects to a new image, where their layout is drastically changed. In addition to the dual encoding, our method presents both a new architecture and new loss terms, which leads to an improved performance over the existing baselines.

## Acknowledgement

# References

[1] David Bau, Jun-Yan Zhu, Hendrik Strobelt, Zhou Bolei, Joshua B. Tenenbaum, William T. Freeman, and Antonio Torralba. Gan dissection: Visualizing and understanding generative adversarial networks. *arXiv preprint arXiv:1811.10597*, 2018. 3

[2] Holger Caesar, Jasper R. R. Uijlings, and Vittorio Ferrari. Coco-stuff: Thing and stuff classes in context. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2, 7

[3] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 2672–2680. Curran Associates, Inc., 2014. 2

[4] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 9

[5] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 6626–6637. Curran Associates, Inc., 2017. 8

[6] Seunghoon Hong, Dingdong Yang, Jongwook Choi, and Honglak Lee. Inferring semantic layout for hierarchical text-to-image synthesis. *CoRR*, abs/1801.05091, 2018. 3

[7] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2

[8] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European Conference on Computer Vision*, 2016. 4

[9] Justin Johnson, Agrim Gupta, and Li Fei-Fei. Image generation from scene graphs. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 1, 2, 3, 4, 5, 7, 9, 10, 12

[10] Justin Johnson, Ranjay Krishna, Michael Stark, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Li Fei-Fei. Image retrieval using scene graphs. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3668–3678, 2015. 2, 9

[11] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2016. 7

[12] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Li Fei-Fei. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International Journal of Computer Vision*, 123:32–73, 2016. 3

[13] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural net-works. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates, Inc., 2012. 9

[14] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley. Least squares generative adversarial networks. In *IEEE International Conference on Computer Vision (ICCV)*, 2017. 4

[15] David Marr. *Vision: A Computational Investigation into the Human Representation and Processing of Visual Information*. Henry Holt and Co., Inc., New York, NY, USA, 1982. 1

[16] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014. 2

[17] Max H. Quinn, Erik Conser, Jordan M. Witte, and Melanie Mitchell. Semantic image retrieval via active grounding of visual situations. *2018 IEEE 12th International Conference on Semantic Computing (ICSC)*, pages 172–179, 2018. 2

[18] Scott Reed, Zeynep Akata, Xinchen Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. Generative adversarial text to image synthesis. In *ICML*, 2016. 2, 3

[19] Tim Salimans, Ian J. Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. In *Advances in Neural Information Processing Systems*, 2016. 4, 8

[20] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. The graph neural network model. *IEEE Transactions on Neural Networks*, 20:61–80, 2009. 3

[21] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 4

[22] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015. 8

[23] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016. 12

[24] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9:2579–2605, 2008. 6

[25] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 2, 4

[26] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaolei Huang, Xiaogang Wang, and Dimitris N. Metaxas. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. *IEEE International Conference on Computer Vision (ICCV)*, pages 5908–5916, 2017. 2, 3

[27] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 9

[28] Bo Zhao, Lili Meng, Weidong Yin, and Leonid Sigal. Image generation from layout. *CoRR*, abs/1811.11389, 2018. 3, 5, 7, 9

[29] Han Zhao, Shanghang Zhang, Guanhang Wu, Jo ao P. Costeira, José M. F. Moura, and Geoffrey J. Gordon. Multiple source domain adaptation with adversarial learning. In *ICLR workshop*, 2018. 7, 8, 9

$B$   $L_{128}, L_{512}, L_4$

$A$   $CB_{64}, CB_{128}, CB_{256}, GA, L_{192}, L_{64}, L_{32}$

$R$   $C_{c+32}, D_{128}, D_{256}, D_{512}, D_{1024}, V_{1024}, V_{1024}, V_{1024},$ $V_{1024}, V_{1024}, V_{1024}, V_{1024}, V_{1024}, V_{1024}, U_{512}, U_{256},$ $U_{128}, U_{64}, C_3$

$D_{\text{mask}}$   $C_{1-64-2}, LR, C^*_{(c+64)-128-1}, IN, LR, C_{128-1-1},$ $AP$

$D_{\text{image}}$   $C_{(c+32)-2}, LR, D_{64-2}, D_{128-2}, D_{256-1}, C_{512-1}$

$D_{\text{object}}$   $CB_{64}, CB_{128}, C4_{256}, GA, L_{1024}, L_1$

($^*$Concatenating the $c_i$ data in $D_{\text{image}}$ is done at the third layer)

## A. Network architecture

The graph convolutional network used is the same as the one used in [9], which was modified in order to support the added node information. We concatenate the embedding of the objects to the location attributes creating vectors in $R^{128+35}$. These vectors, together with the relation embedding is feed-forward into a fully-connected layer which results in a vector embedding in $\mathbb{R}^{128}$ for each of the objects and each of the relations. The network then follows the architecture of [9], using a shared symmetric function to calculate the object vector from all the relations it participate in. Our graph convolution has overall five layers.

To describe the rest of the networks, we follow a semi-conventional shorthand notation for defining architectures. Let $C_k$ denote a Convolution layer of $k$ filters with a kernel size of $7x7$ and a stride of 1, followed by instance normalization [23] and a ReLU activation function. Similarly, we use $D_k$ to a layer which uses a stride of 2, reflection padding, and $k$ filters. In addition, we use $B_k$ to denote a $3x3$ upsample-convolution-BatchNormalization-ReLU layer with $k$ filters and a stride and padding of 1. We use $V_k$ to define Residual blocks with two $3x3$ convolutional layers, both with $k$ filters. Moreover, $U_k$ denotes a layer with $k$ filters of size $3x3$ and a fractional stride of 0.5, followed by instance normalization. $CB_k$ denotes a stride-2 k-filter, $4x4$ convolution followed by batch normalization. $GA$ denotes a global average pooling layer. Fully connected layers with $k$ hidden units followed by a ReLU activation are denoted by $L_k$. The ReLU is not applied to the $L_k$ layer, if it is the top layer.

The discriminators call for an even more elaborate terminology. Let $C_{i-k-o}$ denote a 3x3 Convolution layer with $i$ input channels and $k$ output filters, a stride of $o$ and a padding of 1. In addition, $LR$ denotes Leaky-ReLU with negative slope of 0.2, $IN$ denotes Instance Normalization, and $AP_k$ denotes a $3x3$ Average Pool stride 2 and padding of 1. Also, let $C_{k-s}$ denote a Convolution layer with $k$ filters, stride of $s$, kernel size of $4x4$, and a padding of size 2. $D_{k-s}$ denotes a $4x4$ Convolution-InstanceNorm-LeakyReLU layer with $k$ filters, a stride of $s$ padding of 2 and a LeakyReLU with a negative slope of 0.2.

The different components of the networks can be described as:

$M$   $B_{192}, B_{192}, B_{192}, B_{192}, B_{192}, B_{192}, C_1$, Sigmoid activation