# SCENE GRAPH TO IMAGE GENERATION WITH CONTEXTUALIZED OBJECT LAYOUT REFINEMENT

*Maor Ivgi*, Yaniv Benny*, Avichai Ben-David, Jonathan Berant, Lior Wolf*

Blavatnik School of Computer Science, Tel-Aviv University

## ABSTRACT

Generating images from scene graphs is a challenging task that attracted substantial interest recently. Prior works have approached this task by generating an intermediate layout description of the target image. However, the representation of each object in the layout was generated independently, which resulted in high overlap, low coverage, and an overall blurry layout. We propose a novel method that alleviates these issues by generating the entire layout description gradually to improve inter-object dependency. We empirically show on the COCO-STUFF dataset that our approach improves the quality of both the intermediate layout and the final image. Our approach improves the layout coverage by almost 20 points, and drops object overlap to negligible amounts.

***Index Terms***— Image Synthesis, Scene Graph, GAN

## 1. INTRODUCTION

Synthesizing images from natural language descriptions has received substantial attention recently [1, 2, 3, 4], as it has wide applicability for content generation. However, it has been shown that models that accept textual descriptions as their input fail to produce images with multiple detailed objects with complex relations [1, 2, 4]. Thus, *scene graphs* [5], i.e. graphs where nodes correspond to entities and edges describe relations between them, were proposed [6] as an intermediate representation of the desired image. This approach has been widely adopted [7, 8, 9] for this task.

When generating images from scene graphs (SG), there are three main desiderata: (i) *Photo-realism*: the image should look natural with salient objects, (ii) *Correctness*: the image should contain the objects and relations specified in the SG, and (iii) *Diversity*: because an SG is an underspecified representation compatible with many output images, a model should reflect that in its output distribution. Current models for SG-to-image generation invariably combine a supervised learning objective at training time. Specifically, given an SG and an image they predict for each object separately the *exact* location and shape from the gold semantic layout to produce the ground truth image. Although this can achieve correctness for simple geometric relations, it inevitably results in

_____
* equal contribution

poor quality image-layout due to the under-specificity of the SG. In particular, many distinct images can be represented by the same SG, thus maximum likelihood based techniques result in a blurry average of object shapes and positions across possible images. Such generations are likely to exhibit low resolution, low coverage, and high inter-object overlap. Moreover, due to the strict specification of the prediction task, true diversity is inherently impossible.

In this work we propose two main technical contributions to solve these issues: (i) To address the diversity issue, we reduce the dependence on supervised losses and shift towards adversarial ones. In particular, rather than predicting the box and mask of each object according to the target image, we use an adversarial network as a discriminator. It ensures that the generated object layout is truthful to the required object class in both position and shape and that the relation between every pair of objects is sensible and obeys the constraints dictated by the SG. (ii) To address the quality issue, we introduce a novel method to perform high-resolution layout generation. It incorporates the ability of Graph Convolution Networks (GCNs) [10] to work on variable-shaped structured graphs and contextualizes the state of all objects with CNN-based generators. Using this layout refinement network, we fuse predicted object layouts such that each remains true to its class and respects its dictated relations, while maintaining high coverage and few overlaps. We stack multiple copies of this block and present *Contextualized Objects Layout Refiner* (COLoR): the first model to generate layouts directly from SGs without any intermediate steps such as boxes and masks.
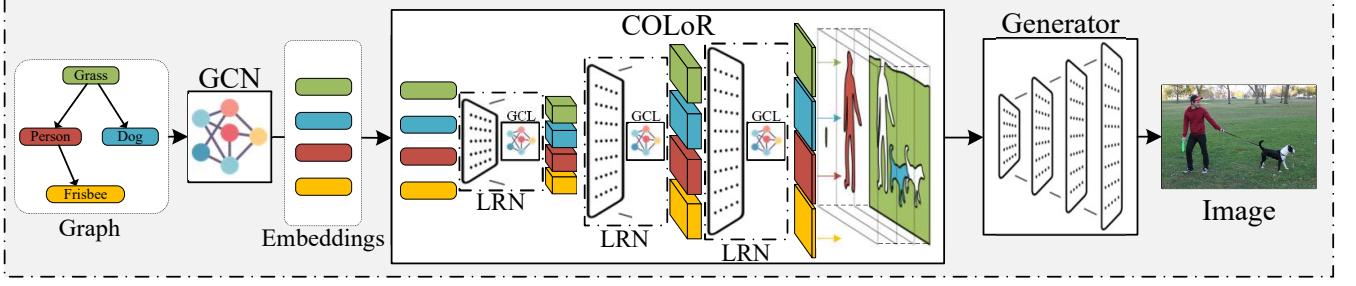
## 2. BACKGROUND

We now describe the architecture of prior work SG2Im [6], which we build upon, and formally define the task.

**Problem setup** Our goal is to train a model $p_\theta(I \mid G, s)$ that takes as input an SG and a random seed $s$ and outputs an image $I$. Given a vocabulary $\mathcal{C}$ of object categories and a set $\mathcal{R}$ of possible relations, an SG is a directed graph $G = (\mathcal{O}, \mathcal{E})$, where each node $o \in \mathcal{O}$ is an object associated with a class by $C : \mathcal{O} \to \mathcal{C}$, and edges $\mathcal{E} \subseteq \mathcal{O} \times \mathcal{O}$ represent directed relations between objects, associated with a type through $R : \mathcal{E} \to \mathcal{R}$.

During training, the available information for every Image $I$ is a segmentation mask, identifying each pixel in the image

**Fig. 1**. Our Scene Graph to Image architecture. A GCN encodes the SG nodes into individual embeddings. The COLoR module upsampled the embeddings with intermediate GC layers into the scene layout. The SPADE generator then produces the image.

to a unique object and its class. This can be used to generate multiple SGs for each image by computing geometric relations between objects and randomly sampling from all possible edges in the complete graph. In addition, this segmentation mask is used to infer the layouts, masks, and bounding boxes for all objects in the image, as defined below.

**Layout generation** An *Image Layout* is a mapping of each pixel in the image to a specific object. Given an object in the layout, we define an *object layout* $\mathbf{l} \in [0,1]^{H \times W}$ as a mapping over the image, indicating pixels that belong to the object. Higher values signify stronger presence of the object. Ideally (as is the case in the annotated layouts), $\mathbf{l}$ is binary. We then define the *Object Box* $\mathbf{b} \in [0,1]^4$ as the minimal axis-aligned bounding box in relative coordinates of all active pixels in $\mathbf{l}$. Finally, cropping $\mathbf{l}$ using $\mathbf{b}$ and projecting it into $[0,1]^{W_m \times W_m}$ where $W_m < \min(H, W)$, we get the *Object Mask* $\mathbf{m} \in [0,1]^{W_m \times W_m}$, which describes its shape, with higher-value pixels corresponding to the existence of the object in said pixel. We note that though $\mathbf{b}$ and $\mathbf{m}$ are derived uniquely from $\mathbf{l}$, it is possible to approximate $\mathbf{l}$ by performing the inverse projection of $\mathbf{m}$ into $[0,1]^{H \times W}$ according to $\mathbf{b}$.

**Scene Graph to Image** We build on SG2Im [6], which includes the following steps: (a) The SG is augmented with an additional dummy node which is connected to all other nodes through an outgoing dummy relation to ensure graph connectivity. (b) Every node and edge in the SG is replaced by a learned embedding $v \in \mathbb{R}^d$ based on its class. (c) The graph is fed to a GCN, which produces a new embedding $\tilde{v}$ for each object (node) in the SG. (d) The embeddings $\tilde{v}_1, \ldots, \tilde{v}_n$ are fed into the layout predictor consisting of two separate decoders, one predicts a bounding box location $\hat{b}_i$, and another predicts a mask $\hat{m}_i$. Those are used to compute $\hat{l}_i$ as explained above. The embedding $\tilde{v}_i$ is multiplied element-wise with $\hat{l}_i$ producing $\hat{\ell}_i \in [0,1]^{H \times W \times d}$. (e) The extended layouts $\hat{\ell}_1, \ldots, \hat{\ell}_n$ are summed element-wise to produce a coarse image layout $\hat{\mathbf{l}} \in [0,1]^{H \times W \times d}$. (f) $\hat{\mathbf{l}}$ is fed along with $z$ random noise channels into a Cascaded Refinement Network (CRN) [11], predicting the final image $\hat{I}$.

In [6], the model is trained with six loss functions: three adversarial loss functions that evaluate object realism, the ability to correctly classify objects, and image similarity to real images. The other three use strong supervision and force

the model to predict boxes and masks which are similar to those in the ground truth image $\mathbf{I}$. Those are:

$$\mathcal{L}_{\text{box}} = \sum_{i=1}^n \|\mathbf{b}_i - \hat{b}_i\|_2, \quad \mathcal{L}_{\text{pix}} = \|\mathbf{I} - \hat{I}\|_1$$
$$\mathcal{L}_{\text{mask}} = \sum_{i=1}^n \sum_{h,w}^{W_m, W_m} \text{BCE}(\mathbf{m}_{i,h,w}, \hat{m}_{i,h,w}) \quad (1)$$
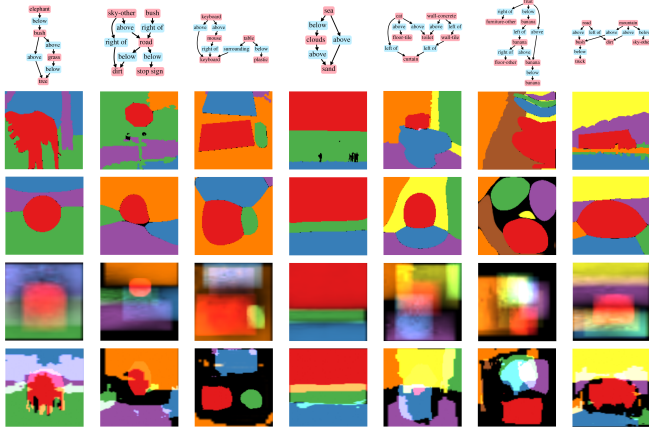
where $\hat{b}_i, \hat{m}_i, \hat{I}$ are the predicted boxes, masks and Image and $\text{BCE}(p, \hat{p}) = -p \log(\hat{p}) - (1-p) \log(1-\hat{p})$.

## 3. METHOD

One major drawback of using the aforementioned supervised loss functions is the underlying assumption that for every SG there exists (in the dataset) **at most one** corresponding image layout. However, in COCO-STUFF [12] which is commonly used for this task, this is far from true, as 73% of the images contain a (multi) set of objects that is shared with many other images in the data and may result in identical SGs. Further, over 25% of the SGs match multiple different layouts, and almost 10% of the SGs describe over 10 different layouts. Thus, a model that maximizes the likelihood of a layout given an SG will be pushed towards predicting the mean of the bounding boxes in the layouts that occur in the training data, and similarly, the average mask. Because the location and shapes of layout substantially vary across images, the model eventually will ignore the context and predict a general location for each object with no distinct shape as can be seen in Figure. 2.

To overcome this difficulty, we remove most of the loss terms that are applied with respect to the **exact** ground truth layout used (§2) and reduce the weight of the rest. Instead, we add adversarial loss functions that encourages the model to generate photo-realistic images that respect the original SG, without forcing it to learn a single SG2Im mapping. The proposed discriminator are applied on the predicted object layouts $\hat{l}_i$. We find that some strong supervision is beneficial to cope with issues that are linked to cold-start.

**Pairwise Layout Discriminator** The main source of training signal in our method is an adversarial network which teaches the generator to be spatially aware, creating objects without overlap that respect the relations set by the SGs. It follows the AC-GAN [13] adversarial loss pattern. Given a pair of neighboring objects in the SG, the discriminator

**Fig. 2**. Comparison of layout generation. Layout values are in $[0, 1]$, depicted as color opacity, thus making overlaps visible. Unassigned pixels are black. Rows from top to bottom: Input Scene graph, true layout, COLoR, SG2Im [6], Grid2Im [7].

accepts their object layouts $\mathbf{l}_i, \mathbf{l}_j \in [0, 1]^{H \times W}$ and their class labels $c_i, c_j \in \mathcal{C}$. It predicts whether this pair comes from a real or a generated layout and classify the relation between the two. Since low-quality layouts will be easily recognized as fake, it also improves the quality of all object layouts individually. In particular, The discriminator $D_l$ performs a mapping $D_l : ([0, 1]^{H \times W}, \{0, 1\}^{|\mathcal{C}|})^2 \rightarrow [0, 1] \times [0, 1]^{|\mathcal{R}|}$. Let $(\hat{y}, \hat{r}) = D_r((\mathbf{l}_i, c_i), (\mathbf{l}_j, c_j))$ be the prediction of the discriminator (real vs. fake and relation prediction) on its input. Let $r \in \{0, 1\}^{|\mathcal{R}|}$ be the true relation and $y = \mathbb{1}_{real}$.

$$\mathcal{L}_{D_l}^d = \text{BCE}(y, \hat{y}) + \text{CE}(r, \hat{r})$$
$$\mathcal{L}_{D_l}^g = \text{BCE}(1, \hat{y}) + \text{CE}(r, \hat{r}) \quad (2)$$

where the discriminator trains on real and fake pairs, and the generator minimizes the loss over generated pairs only.

**Losses** To complement our discriminator and encourage the model to assign objects to every pixel of the image and refrain from overlap, we introduce the *Layout Coverage Regularization* $\mathcal{L}_{\text{reg}} = \mathcal{L}_{\text{coverage}} + \lambda \cdot \mathcal{L}_{\text{overlap}}$. Given $\hat{l}_1, \dots, \hat{l}_n \in [0, 1]^{H \times W}$ we define the summed image layout $\hat{L} = \sum_i^n \hat{l}_i \in [0, n]^{H \times W}$ which gives the following definitions:

$$\mathcal{L}_{\text{coverage}} = \sum_{h,w}^{H,W} \mathbb{1}[\hat{L}_{h,w} \leq 1] \cdot \left(1 - \hat{L}_{h,n}\right) \quad (3)$$

$$\mathcal{L}_{\text{overlap}} = \sum_{h,w}^{H,W} \mathbb{1}[\hat{L}_{h,w} > 1] \cdot \left(\hat{L}_{h,n} - 1\right) \quad (4)$$

The loss reaches 0 if the layouts weights in every pixel sum to exactly 1, and grows as the coverage drops or overlap increases. Since $\mathcal{L}_{\text{overlap}}$ is unbounded and often grows larger than $\mathcal{L}_{\text{coverage}}$, we suppress its contribution by setting $\lambda = 0.4$ which we found achieves the best tradeoff between the terms.

In addition, we found that when tasked to generate layouts with only the losses in equations (2) and (4), the generator fails to learn how to create coherent layouts, and the discriminator falls back to classify fake layouts based on spurious artifacts in the layouts. We attribute this issue to cold-start

problem, and mitigate it by adding small weight to an *Object Layout* loss defined on each predicted object layout $\hat{l}$ and the corresponding ground-truth layout $\mathbf{l}$: $\mathcal{L}_{\text{layout}} = \sum_i^n \|\mathbf{l}_i - \hat{l}_i\|_1$.

**Mapping embeddings to layouts directly** In prior work, boxes $b_i$ and masks $m_i$ were decoded from object embeddings $\tilde{v}_i$ in parallel. Hence, the embeddings $\tilde{v}_i$ computed by the GCN must encode all the information about the location and shape of each object, including avoiding inter-object overlap and maintaining high-coverage of the layout. Furthermore, since the dimension of the mask is $W_m \times W_m$, which is then warped into $H \times W$, and $W_m \ll \min(H, W)$, we can expect a drop in resolution (i.e. very coarse shapes). To mitigate these issues while remaining agnostic to the SG size and structure, we propose an adaptation of the GCN technique to improve the object layouts generation process by contextualizing object layouts on each other. We name this module *Layout Refinement Network* (LRN). Formally, given some intermediate object representations $\hat{v}_1^t, \dots, \hat{v}_n^t \in \mathbb{R}^{C_t \times H_t \times W_t}$, we describe a model that predicts the next representations down the line $\hat{v}_1^{t+1}, \dots, \hat{v}_n^{t+1} \in \mathbb{R}^{C_{t+1} \times H_{t+1} \times W_{t+1}}$, with $H_{t+1} = 2 \cdot H_t$, $W_{t+1} = 2 \cdot W_t$, $C_{t+1} \leq C_t$.

$$\hat{v}_1^{t+1}, \dots, \hat{v}_n^{t+1} = LRN^t(\hat{v}_1^t, \dots, \hat{v}_n^t) \quad (5)$$

First, each representation $\hat{v}_{n'}^t, n' \in [1, n]$ is passed through a decoder $U$ that applies a transposed convolution layer to upsample the representation by a factor of two, followed by a batch normalization layer and a ReLU activation; $q_{n'}^t = U(\hat{l}_{n'}^t)$. Each pair of upsampled representations ($\{q_i^t, q_j^t\}|i \neq j \in [1, n]$) is then passed through a graph convolution layer. Due to the dimensionality of the representations ($C_t \times H_t \times W_t$), the traditional dense layers of the graph convolution are replaced with 2D-convolutional layers: $q_{i,j}^t, q_{j,i}^t = GCL^t(q_i^t, q_j^t)$. Summing the initial and pairwise representation produce a new representation that is contextualized on all objects in the scene: $\hat{v}_{n'}^{t+1} = q_{n'}^t + \frac{1}{n-1} \left(\sum_{i \neq n} q_{n',i}^t + \sum_{i \neq n'} q_{i,n}^t\right)$. In intermediate stages, the residual sum is followed by a ReLU activation. In the final stage, a sigmoid is applied to create the object layout. Stacking $T$ LRN blocks (where $T = \log_2 H - 1$ due to the behavior of transposed convolutions in the first stage), we skip box and mask predictions Entirely. Instead, our model (depicted in Figure. 1), generates the layout directly from object embedding. Samples are shown in Figure. 3. Given embeddings $\tilde{v}_i^1$ of size $1 \times 1$ and depth $K = 2^T$, we stack $T$ layers of upsampling which reduces the depth by a factor of two followed by an LRN. The output of the model is a set of object layouts $\hat{l}_1, \dots, \hat{l}_n \in [0, 1]^{H \times W}$. In each stage, the LRN is used to pass information between the layouts which results in a coherent layout exhibiting extremely high coverage and negligible overlaps. We name our model *Contextualized Objects Layout Refiner* (COLoR). To reduce computational constraints and allow for diversity in the generation, we sample a random subset of all possible layout pairs in each layer.
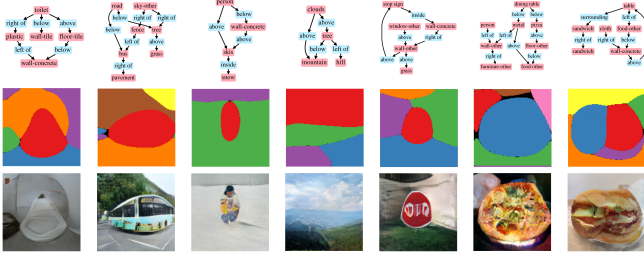
**Fig. 3**. Examples of COLoR generated images from SGs.

| | | Layout | | | | Image | |
|---|---|---|---|---|---|---|---|
| | Model | COV↑ | OVL↓ | DEC↑ | GRS↑ | FID↓ | DIV↑ |
| **Baselines** | SG2Im[6] | 0.76 | 0.18 | 0.83 | 0.62 | 124.3 | 0.0 |
| | SG2Im+SPADE | | unchanged | | | 97.7 | 0.05 |
| | Grid2Im[7] | 0.81 | 0.17 | 0.98 | 0.94 | 96.4 | 0.0 |
| | Grid2Im+SPADE | | unchanged | | | 106.6 | 0.0 |
| **Ours** | COLoR | **0.99** | **0.0** | **1.0** | **0.97** | **95.8** | 0.13 |
| | - $\mathcal{L}_{D_l}$ | **0.99** | **0.0** | **1.0** | 0.55 | 102.9 | 0.09 |
| | - $\mathcal{L}_{\text{layout}}$ | 0.82 | **0.0** | **1.0** | 0.60 | 122.7 | **0.50** |

**Table 1**. Layout quality evaluation with Coverage, Overlap, Decisiveness, and Geometric-Relation-Score. Image quality evaluation with FID and Perceptual Diversity.

## 4. EXPERIMENTS

We train our model to generate 128x128 layouts and use a pretrained SPADE [14] model to predict images from them. We show that it creates high-quality layouts, respects the SG's constraints and results in overall higher quality images compared to prior work. We follow the setup in [6, 7] using the COCO-STUFF dataset [12]. This dataset contains a subset of the images in COCO [15] with additional 91 *stuff* categories. We compare our model against a pretrained model of [7], and a model of [6] trained to produce images of the same resolution. We evaluate layout quality, diversity and adherence to SGs and the predicted image's quality.

**Layout generation**    To evaluate the quality of the generated layouts, we measure the average coverage, overlap, and decisiveness of the layouts. Coverage ranges between 0 to 1 where higher values are better as it means the layout does not contain empty spots. Overlap measures if multiple objects occupy the same pixel which we wish to avoid. The decisiveness measure evaluates how decisive the generator is in deciding the pixels' pertinence. Formally, given predicted layouts $\hat{l}_1, \ldots, \hat{l}_n \in [0,1]^{H \times W}$ we threshold the layouts $\hat{l}_{i,h,w}^t = \mathbb{1}[\hat{l}_{i,h,w} \geq t]$ setting $t = 0.5$ to get $\hat{l}_1^t, \ldots, \hat{l}_n^t \in \{0,1\}^{H \times W}$. We then define the coverage as $\frac{1}{M} \sum_{h,w} \mathbb{1}[\sum_i^n \hat{l}_{i,h,w}^t \geq 1]$, the overlap as $\frac{1}{M} \sum_{h,w} \mathbb{1}[\sum_i^n \hat{l}_{i,h,w}^t \geq 2]$, and the decisiveness is defined as $\frac{4}{M} \sum_{i,h,w} (0.5 - \hat{l}_{i,h,w})^2$. Where $M = H \times W$.

Finally, to evaluate the compliance with relation constraints, we define the *Geometric-Relation-Score* (GRS). Given a pair of predicted object layouts $\hat{l}_i, \hat{l}_j$, we compute the minimal axis-aligned bounding rectangle that contains all pixels with values above 0.5, and projecting it to $W_m \times W_m$

to get a mask. We then use the same heuristic that was used in the construction of the dataset to infer the relations between objects in the predicted layouts, and define the geometric relation score as the accuracy of these predictions.

**Image generation**    To evaluate the quality of the generated images, we use the common *FID* score. We augment this evaluation with the diversity measure suggested by [7], which relies on the Perceptual Similarity measure [16]. There, multiple images are generated from the same SG, and we measures the average distance between every pair, where large distance is correlated with diversity.

**Results**    As depicted in Table 1, our method outperforms on all *layout generation* benchmarks by large margins. It predicts layouts that have both high coverage and low overlap. In addition, the layouts are decisive and fulfill the geometric relations specified by the SG. The image generation quality of our model is preferable compared to the baselines according to both the FID measure and the diversity score. Our model achieves an FID and diversity score of 95.8 and 0.13 respectively. Our model also shows more diversity when generating multiple images from the same SG than the baselines. It should be noted that one of the ablations ($- \mathcal{L}_{\text{layout}}$) scored very high on the diversity benchmark due to its failure to generate reasonable layouts consistently, which means that diversity on its own is not sufficient, and should always be evaluated in conjunction with an image quality benchmark.

**Ablations**    We study the contributions of $\mathcal{L}_{D_l}$ and $\mathcal{L}_{\text{layout}}$ by training COLoR models without them. It can be seen in Table 1 that removing the pairwise layout discriminator $D_l$ impairs the GRS and removing $\mathcal{L}_{\text{layout}}$ hurts the layout coverage. To show that improvements in the final image quality are due to improved layouts, and not due to SPADE's superiority over prior work's layout-to-image networks, we evaluate the image generation quality of SG2Im [6] and Grid2Im [7] layouts by replacing their layout-to-image modules (CRN [11] and Pix2Pix [17] respectively) with the SPADE [14] model we use. We find that the FID score is heavily dependent on the type of generator and not necessarily on the quality of the layout, which was the main focus of this work. Both SG2Im and Grid2Im scored differently using their own generators compared to using SPADE. However, SG2Im score improves while Grid2Im suffers.

## 5. CONCLUSIONS

We presented a new technique to train a model to directly predict object layouts from an abstract scene description while attending all objects simultaneously. Our method achieves a sizable improvements in the layouts' quality compared to prior works, resulting in accurate and photo-realistic images.

## 6. ACKNOWLEDGEMENTS

## 7. REFERENCES

[1] Han Zhang, Tao Xu, and Hongsheng Li, "Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks," in *IEEE International Conference on Computer Vision, ICCV*. 2017, pp. 5908–5916, IEEE Computer Society.

[2] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris Metaxas, "Stackgan++: Realistic image synthesis with stacked generative adversarial networks," *arXiv preprint arXiv:1710.10916*, 2017.

[3] Wenbo Li, Pengchuan Zhang, Lei Zhang, Qiuyuan Huang, Xiaodong He, Siwei Lyu, and Jianfeng Gao, "Object-driven text-to-image synthesis via adversarial training," in *2019 IEEE Conference on Computer Vision and Pattern Recognition, CVPR*. 2019, pp. 12174–12182, Computer Vision Foundation / IEEE.

[4] Tao Xu, Pengchuan Zhang, Qiuyuan Huang, Han Zhang, Zhe Gan, Xiaolei Huang, and Xiaodong He, "Attngan: Fine-grained text to image generation with attentional generative adversarial networks," in *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR*. 2018, pp. 1316–1324, IEEE Computer Society.

[5] Justin Johnson, Ranjay Krishna, Michael Stark, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Fei-Fei Li, "Image retrieval using scene graphs," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*. 2015, pp. 3668–3678, IEEE Computer Society.

[6] Justin Johnson, Agrim Gupta, and Li Fei-Fei, "Image generation from scene graphs," in *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR*. 2018, pp. 1219–1228, IEEE Computer Society.

[7] Oron Ashual and Lior Wolf, "Specifying object attributes and relations in interactive scene generation," in *2019 IEEE/CVF International Conference on Computer Vision, ICCV*. 2019, pp. 4560–4568, IEEE.

[8] Roei Herzig, Amir Bar, Huijuan Xu, Gal Chechik, Trevor Darrell, and A. Globerson, "Learning canonical representations for scene graph to image generation," in *ECCV*, 2020.

[9] Subarna Tripathi, Anahita Bhiwandiwalla, Alexei Bastidas, and Hanlin Tang, "Using scene graph context to improve image generation," *CoRR*, vol. abs/1901.03762, 2019.

[10] Thomas N. Kipf and Max Welling, "Semi-supervised classification with graph convolutional networks," in *5th International Conference on Learning Representations, ICLR, Conference Track Proceedings*. 2017, OpenReview.net.

[11] Qifeng Chen and Vladlen Koltun, "Photographic image synthesis with cascaded refinement networks," in *IEEE International Conference on Computer Vision, ICCV*. 2017, pp. 1520–1529, IEEE Computer Society.

[12] Holger Caesar, Jasper R. R. Uijlings, and Vittorio Ferrari, "Coco-stuff: Thing and stuff classes in context," in *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR*. 2018, pp. 1209–1218, IEEE Computer Society.

[13] Augustus Odena, Christopher Olah, and Jonathon Shlens, "Conditional image synthesis with auxiliary classifier gans," in *Proceedings of the 34th International Conference on Machine Learning, ICML*, Doina Precup and Yee Whye Teh, Eds. 2017, vol. 70 of *Proceedings of Machine Learning Research*, pp. 2642–2651, PMLR.

[14] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu, "Semantic image synthesis with spatially-adaptive normalization," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*. 2019, pp. 2337–2346, Computer Vision Foundation / IEEE.

[15] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick, "Microsoft coco: Common objects in context," in *European conference on computer vision*. Springer, 2014, pp. 740–755.

[16] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR*. 2018, pp. 586–595, IEEE Computer Society.

[17] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A. Efros, "Image-to-image translation with conditional adversarial networks," in *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR*. 2017, pp. 5967–5976, IEEE Computer Society.