

VARSCENE: A Deep Generative Model for Realistic Scene Graph Synthesis

Anonymous Authors¹

Abstract

Scene graphs are powerful abstractions that capture relationships between objects in images by modeling objects as nodes and relationships as edges. Generation of realistic novel scene graphs has applications like scene synthesis and data augmentation for supervised learning. Existing graph generative models are predominantly targeted toward molecular graphs, leveraging the limited vocabulary of atoms and bonds and also the well-defined semantics of chemical compounds. In contrast, scene graphs have much larger object and relation vocabularies, and their semantics are latent. To address this challenge, we propose VARSCENE, a variational autoencoder for scene graphs, which is optimized for the maximum mean discrepancy (MMD) between the ground truth scene graph distribution and distribution of the generated scene graphs. VARSCENE views a scene graph as a collection of star graphs and encodes it into a latent representation of the underlying stars. The decoder generates scene graphs by learning to sample the component stars and edges between them. Our experiments show that our method is able to mimic the underlying scene graph generative process more accurately than several state-of-the-art baselines.

1. Introduction

Scene graphs (Johnson et al., 2015; 2018; Chang et al., 2021) have emerged as a powerful data structure to represent images. The nodes in these graphs represent objects and their attributes, and edges capture relationships between them. They are a succinct and human-consumable summary of the content within the image. Their popularity is partially derived from the availability of rich datasets (Lu et al., 2016; Krishna et al., 2016) that contain ground-truth scene graphs associated with images. Extracting scene graphs from images continues to be a very active field of research

(Tang et al., 2020b; Yang et al., 2018; Gu et al., 2019). The success of scene graphs can be seen in the range of visual understanding tasks that utilize them: retrieval (Schroeder & Tripathi, 2020; Wang et al., 2020; Maheshwari et al., 2021), editing (Dhamo et al., 2020), question answering (Teney et al., 2017) and captioning (Milewski et al., 2020).

While there has been progress in the field of image generation, the end goal of achieving rich, diverse and complex scenes is far from satisfactory (Casanova et al., 2020; Li et al., 2019). In the current paper, we propose the novel task of *scene graph generation* as a stepping stone to scene synthesis. Here, our goal is to learn to synthesize scene graphs based on a set of training examples of scene graphs. Our work is relevant to the theme of generating new and complex scenes, where scene graphs are often used as a means for conditioning image synthesis (Johnson et al., 2018). We decouple the generation of scene semantics (sets of objects and how they relate to each other) from its visual manifestation, and focus specifically on the first part in the current paper. Where required, we will use existing methods to convert our synthesized scene graphs to images (Johnson et al., 2018; Tseng et al., 2020).

Driven by the success at text and image synthesis, there is much recent interest in deep generative models for graphs, e.g., GraphVAE (Simonovsky & Komodakis, 2018), GraphRNN (You et al., 2018), NeVAE (Samanta et al., 2019), MolGAN (De Cao & Kipf, 2018), GRAN (Liao et al., 2019), etc. However, most of them are tailored to molecule discovery and cannot be effectively used to generate scene graphs, as we demonstrate in our experiments. Very recently, Garg et al. (2021) proposed an unconditional scene graph generator, which, however, is not optimized to capture the underlying data distribution.

1.1. Present Work

We propose VARSCENE, a novel variational autoencoder, specifically designed for scene graph generation. (We use graph ‘generation’ and ‘synthesis’ interchangeably.) In contrast to prior VAE-based graph generative models which maximize Evidence Lower Bound (ELBO) for model training, VARSCENE is optimized to generate graphs having minimum distribution discrepancy with the ground truth (training) graph distribution, while using its generative capacity to introduce plausible variations. We thus obtain

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

a model from which we can sample graphs representing realistic scenes. Moreover, due to the encoder-decoder architecture, our approach facilitates conditional scene graph generation, which would otherwise not be possible (Garg et al., 2021).

Generating scene graphs with meaningful semantics is a challenging task. Designing simplistic masking functions, as for molecular graphs, does not work in this context, since the notion of underlying semantics is far more complex than in molecules. Responding to this challenge, VARSCENE views a scene graph as a collection of *stars*, where a star is comprised of a ‘hub’ node with its incident edges (‘spokes’). Our encoder embeds the graph into a collection of latent representation vectors based on its component stars. At the other end, unlike most existing generative models (You et al., 2018; Samanta et al., 2019; De Cao & Kipf, 2018), our decoder does not expand the graph by generating a node or an edge at a time, but instead by sampling stars. Specifically, it first samples stars from a trainable distribution in a sequential manner, and then connects these stars to generate semantically meaningful scene graphs. Such a model allows us to generate novel scene graphs while retaining the semantic correlations between the different components of a scene. In contrast to generative models for molecular graphs, which work with few node and edge labels, VARSCENE is able to generate (and be trained on) scene graphs consisting of a large number of node and edge labels.

In theory, the variational autoencoder described above should be able to learn the underlying data distribution, given sufficient training data. However, due to the use of an approximate posterior and lower bound of the true objective, it may not show optimal performance in practice. In response, we further re-train the decoder to directly minimize maximum mean discrepancy (MMD) between the distributions of generated and ground truth scene graphs. Such a construction allows our method to effectively trade off between the underlying MMD and the fidelity to the prior decoder model previously obtained using ELBO maximization.

We evaluate VARSCENE on three real world datasets. Our experiments show that VARSCENE is able to mimic the underlying distribution of scene graphs more accurately than several baselines.

2. Preliminaries

2.1. Notation and Definitions

Scene Graphs, Objects, Relations. Given a collection of directed scene graphs $\{G = (V, E)\}$ along with the set of objects T and the set of relations R , we denote t_u as the *object* represented by the node $u \in V$ and r_e as the *relation* represented by the edge $e \in E$. For example, if the node pair (u, v) in the scene graph contains the semantic

content: “car \rightarrow moving on \rightarrow road”, then we use $t_u = \text{car}$, $t_v = \text{road}$ and $r_e = \text{moving on}$. We denote t_u and r_e as the feature vectors for node u and edge e respectively. Such feature vectors may be obtained by various means, depending on the application. Here we use BERT embeddings of the names of the object t_u and the relation r_e . Features may also be derived of object color and texture. Finally, we define $\text{nbr}(u)$ as the neighbors of node u .

Star and Neighbor-Stars. Given a graph G and a node u , the *star* s rooted at u is represented by a pair consisting of the node feature vector t_u and the multiset of features of edges incident on u . I.e., the representation of s is $s := \langle t_u, \{r_{(u,v)} \mid v \in \text{nbr}(u)\} \rangle$. We denote u as $\text{root}(s)$. Here, the star consists of only one node, i.e., the central node, characterized by the underlying object and the associated open-ended edges characterized by the relations. Thus, the identity of u , or the identity or type of its neighbors are not included in the representation of s . Only the object type of u and the types of incident edges matter. Therefore, it is quite possible that two stars having different root node IDs have identical representation, i.e., $s = s'$ but $\text{root}(s) \neq \text{root}(s')$. Given a dataset of graphs D , we denote the set of stars in its “star vocabulary” as $S(D) = \{s \mid s \subset G, G \in D\}$.

Given a star s in graph G , we define its neighbor-stars $\mathcal{N}(s)$ as those stars whose roots are connected to the root of s via an edge. We also define $\gamma(s, s')$ as the relation type r_e of the edge e that connects $\text{root}(s)$ and $\text{root}(s')$. Formally,

$$\mathcal{N}(s) = \{(s', \gamma(s, s')) \mid (\text{root}(s), \text{root}(s')) \in E(G) \text{ for some } G \in D\}.$$

Thus, $\mathcal{N}(s)$ contains pairs of stars s and relation r_e , where s and s' are connected via the edge e . When clear from context, we will drop the edge type and write “ $s \in \mathcal{N}(s)$ ” instead of “ $s, \gamma \in \mathcal{N}(s)$ ”.

Maximum Mean Discrepancy (MMD). We will use MMD to compute the discrepancy between the distribution of the generated graphs and the corresponding true graph distribution (Goyal et al., 2020; You et al., 2018). We define MMD between distributions P and Q as:

$$\text{MMD}^2(P, Q) = \mathbb{E}_{G_1, G_2 \sim P} k(G_1, G_2) + \mathbb{E}_{G'_1, G'_2 \sim Q} k(G'_1, G'_2) - 2\mathbb{E}_{G \sim P, G' \sim Q} k(G, G') \quad (1)$$

Here, $k(G, G')$ is the kernel induced by a suitable RKHS. We have samples D and D' drawn from p_1 and p_2 respectively. Thus, an unbiased estimator of MMD, denoted as $\widehat{\text{MMD}}(D, D')$, is given by

$$\left[\frac{1}{\binom{|D|}{2}} \sum_{G_1, G_2 \in D} k(G_1, G_2) + \frac{1}{\binom{|D'|}{2}} \sum_{G'_1, G'_2 \in D'} k(G'_1, G'_2) - \frac{2}{|D||D'|} \sum_{G \in D, G' \in D'} k(G, G') \right]^{1/2} \quad (2)$$

In practice, we use the following RBF kernel on some rep-

resentation vectors $\nu(G)$ and $\nu(G')$ of the graphs, *i.e.*,

$$k(G, G') = \exp(-\|\nu(G) - \nu(G')\|_2^2/2). \quad (3)$$

2.2. Overview of Graph Variational Autoencoder

Graph VAEs are trainable generative models involving an observable graph G and its latent representations¹ Z , together with three distributions: (i) a prior distribution $p_0(Z)$ over the latent code Z , (ii) the encoder $q(Z|G)$ which embeds the observed graphs G into (a distribution over) the latent code z , and (iii) the decoder $p(G|Z)$ which samples observable graphs G based on the latent representation Z . Computation of the log-likelihood of observables under a VAE model requires marginalization over the distribution of Z , which is usually intractable. Therefore, we typically maximize the evidence lower bound (ELBO) on the log-likelihood:

$$\sum_{G \in D} \left[\mathbb{E}_{Z \sim q(\cdot|G)} [\log p(G|Z)] + KL(q(Z|G)||p_0(Z)) \right] \quad (4)$$

The encoder $q(Z|G)$ is intended to approximate the true posterior of Z . Therefore, the quality of the encoder, and thus, the tightness of the lower bound, depends on the expressivity of q . Typically, it is modeled using a universal distribution approximator, parameterized on G .

3. VARSCENE Model

Here we develop VARSCENE in full detail, formulating the underlying MMD optimization problem and presenting its different components and design decisions.

3.1. Overview

Our goal is to design a variational autoencoder that is able to synthesize realistic scene graphs in both unconditional settings (given a corpus of real scene graphs), and conditional settings (given a corpus and also a specific scene graph as a starting point). While we want the synthetic graphs to exhibit stochastic variation beyond training graphs, we want the distribution of various graph properties to remain close to those measured on the training sample. The cornerstone of our approach is to build the decoder that minimizes the MMD between these two distributions, rather than only maximizing ELBO (4).

Given a hypothesis class \mathcal{P} , we wish to learn a generative distribution $p^{\text{MMD}}(G) \in \mathcal{P}$ which minimizes MMD between the distribution of real scene graphs and the distribution of the generated graphs. However, training such a generative model without any inductive bias is unlikely to succeed. To tackle this problem, we first train a base variational autoencoder (p_0, q, p) on the training scene graphs. While the trained decoder p is not optimized for minimum MMD, it provides an approximate generative model to start with. We

¹ Z is usually a set of latent codes associated with embedding vectors of nodes, edges or substructures.

next use this trained decoder p to guide the training of our MMD optimized generative model p^{MMD} . To this end, we design an optimization problem around p^{MMD} which keeps it close to the pre-trained decoder p , while also minimizing MMD with respect to the training data. Finally, we develop a gradient based algorithm to search for p^{MMD} .

3.2. Design of MMD-Optimized Decoder p^{MMD}

Learning p^{MMD} Given the True Generator p_{true} . If we are given the true generator p_{true} , then, in principle, we can directly estimate a generative distribution $p^{\text{MMD}} \in \mathcal{P}$ over scene graphs, which minimizes MMD between p^{MMD} and the true distribution:

$$\underset{p^{\text{MMD}} \in \mathcal{P}}{\text{minimize}} \text{MMD}(p_{\text{true}}, p^{\text{MMD}}). \quad (5)$$

The above setup faces two bottlenecks: (i) It does not have access to the true distribution p_{true} . It can only use the observed scene graphs present in the training dataset. (ii) Specifying \mathcal{P} without any prior knowledge about the generative process is difficult in practice.

Learning p^{MMD} Given Observed Graphs. Responding to the above challenges, we approximate (5) through several intermediate steps. First, we use the observed scene graphs to estimate $\widehat{\text{MMD}}$, as we cannot access the true graph distribution p_{true} . Specifically, given D , a set of observed scene graphs, we replace the objective in the optimization (5) with its sample estimate defined in Eq. (2) as follows:

$$\underset{p^{\text{MMD}} \in \mathcal{P}}{\text{minimize}} \mathbb{E}_{D' \sim_{\text{id}} p^{\text{MMD}}(\cdot)} \widehat{\text{MMD}}(D, D') \quad (6)$$

However, absence of any prior knowledge of the underlying generative process \mathcal{P} makes it difficult to search for p^{MMD} . Now, suppose that we have a trained VAE model $(p_0, q_{\hat{\phi}}, p_{\hat{\theta}})$ modeled by neural networks with estimated parameters $\hat{\phi}$ and $\hat{\theta}$, where the encoder $q_{\hat{\phi}}(\cdot|G)$ embeds the graph into a latent representation Z , the decoder $p_{\hat{\theta}}(\cdot|Z)$ generates a graph G conditioned on Z and p_0 is a prior distribution on Z . We can use this trained model to guide the training of p^{MMD} , where we use the same neural network $p_{\hat{\theta}}$ to parameterize $p^{\text{MMD}} \approx p_{\theta}^{\text{MMD}}$ with a new parameter vector θ . To this end, given a trained VAE model $(p_0, q_{\hat{\phi}}, p_{\hat{\theta}})$ we learn p_{θ} which minimizes the required MMD, while penalizing its KL divergence with respect to the trained decoder $p_{\hat{\theta}}$, *i.e.*, we solve:

$$\underset{\theta \in \Theta}{\text{minimize}} \mathbb{E}_Z \left[\mathbb{E}_{D(Z) \sim_{\text{id}} p_{\theta}^{\text{MMD}}(\cdot|Z)} \widehat{\text{MMD}}(D, D(Z)) + \rho \text{KL}(p_{\theta}^{\text{MMD}}(\cdot|Z) || p_{\hat{\theta}}(\cdot|Z)) \right], \quad (7)$$

where $D(Z)$ is a set of generated graphs by p_{θ}^{MMD} conditioned on the latent representations Z . Here, $p_{\hat{\theta}}$ and p_{θ}^{MMD} share the same parameterized class of generative models $\mathcal{P}(\Theta)$. In the outer expectation, Z can be drawn from the prior distribution p_0 for unconditional graph generation or from the trained encoder $q_{\hat{\phi}}$ for conditional graph genera-

tion. Moreover, $\rho > 0$ is a tunable coefficient for the KL regularizer which allows us to trade off between minimizing MMD and the fidelity of p^{MMD} to the original decoder p . Such a KL divergence based regularization allows p_{θ}^{MMD} to generate plausible scene graphs with meaningful semantics.

Apart from VAEs, such regularization using KL divergence with respect to a prior distribution has also been used in reinforcement learning (Todorov, 2009), social welfare based control (Tabibian et al., 2020), property oriented molecule generation (Samanta et al., 2019), etc. However, our key goal here is MMD optimization for scene graph generation, where we must steer the property of the entire set of generated graphs, whereas the above tasks are concerned with property of one instance, e.g., generating molecule with a desired property, re-ranking tweets for misinformation mitigation, etc.

Computation of $\widehat{\text{MMD}}(D, \{G\})$. As suggested by Eq. (2) and (3), computation of $\widehat{\text{MMD}}(D, \{G\})$ requires us to featurize G into $\nu(G)$. Graph kernels (Vishwanathan et al., 2010) provide some guidance to our choice of ν , but judging whether a synthetic graph is “close to natural” requires the comparison of a diverse bouquet of graph properties, covering complex interactions between topology, node and edge labels. Because our unit of graph generation is a star, we will focus, for the most part, on a definition of $\nu(G)$ in terms of the number of different stars present in the graph. (For evaluation, though, we will include other important properties.) We gather the universe of stars $S(D)$ present in the dataset D , defined as

$$\nu_s(G) = \# \text{ occurrences of } s \text{ in } G \quad (8)$$

and then compute $\nu(G) = [\nu_s(G) : s \in S(D)]$, which is then used to compute the required MMD using Eq. (3).

Gradient Computation for Solving (7). We solve optimization (7) using a gradient descent algorithm on the finite sample estimate of the underlying objective, i.e.,

$$C_{\theta}(\{G\}) = \mathbb{E}_Z \left[\widehat{\text{MMD}}(D, D(Z)) + \frac{\rho}{|D(Z)|} \sum_{G \in D(Z)} \log \frac{p_{\theta}^{\text{MMD}}(G | Z)}{p_{\hat{\theta}}(G | Z)} \right]. \quad (9)$$

Computing the gradient of such a quantity may seem difficult, since the samples $\{G\}$ are drawn from the distribution itself which is being trained. To circumvent this problem, we use the log-derivative trick (Williams, 1992) to estimate the gradient $\nabla_{\theta} C_{\theta}(\{G\})$ as follows:

$$\nabla_{\theta} \left[\sum_{G \in D(Z)} \widehat{\text{MMD}}(D, D(Z)) \log p_{\theta}^{\text{MMD}}(G | Z) + \frac{\rho}{|D(Z)|} \sum_{G \in D(Z)} \left(\frac{1}{2} \log^2 \frac{p_{\theta}^{\text{MMD}}(G | Z)}{p_{\hat{\theta}}(G | Z)} + \log p_{\theta}^{\text{MMD}}(G | Z) \right) \right], \quad (10)$$

The term within the above gradient can be viewed as a

pseudo loss function which can be easily coded in any ML library to solve the optimization in Eq. (7).

3.3. Neural Architecture of VARSCENE

In this section, we present the neural architecture of the various components of VARSCENE, beginning with a brief outline.

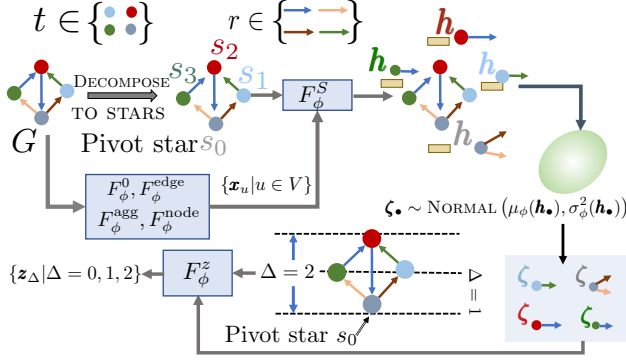
Outline. VARSCENE consists of a base variational autoencoder $(p_0, q_{\theta}, p_{\theta})$ and the final generative distribution p_{θ}^{MMD} . The encoder $q_{\theta}(\cdot)$ views a graph as a set of stars $\{s_0, s_1, \dots\}$, where we call the star s_0 to be the pivot star. Then, we consider various stars $\{s_i \mid \text{DISTANCE}(\text{root}(s_0), \text{root}(s_i)) = \Delta\}$ lying at a given number of hops Δ from the pivot star s_0 , and obtain an aggregated representation vector z_{Δ} over these stars. Thus, our encoder q_{θ} encodes the stars into the corresponding latent representations $Z = \{z_0, \dots, z_{\Delta_{\max}}\}$, where Δ_{\max} is the maximum distance of any star from the pivot s_0 . The corresponding decoder $p_{\hat{\theta}}(\cdot)$ generates stars in a sequential manner by sampling from a softmax distribution conditioned on the latent representations Z . Given a set of training graphs along with their object and relation types, i.e., $D = \{G\}$ along with the object types $\{T\}$ and the relation types $\{R\}$, VARSCENE first learns the encoder $q_{\theta}(\cdot)$ and decoder $p_{\theta}(\cdot)$ by maximizing ELBO (evidence lower bound) of the likelihood function. More specifically, VARSCENE consists of four components:

1. Prior: $p_0(Z)$ with $\Delta_{\max} \sim \text{Poisson}(\lambda)$
2. Encoder: $q_{\theta}(Z \mid G, T = \{t_u\}, R = \{r_e\})$
3. Base decoder: $p_{\theta}(\{s_i\}, \{\gamma_{ij}\} \mid Z)$
4. MMD-optimized decoder: $p_{\theta}^{\text{MMD}}(\{s_i\}, \{\gamma_{ij}\} \mid Z)$

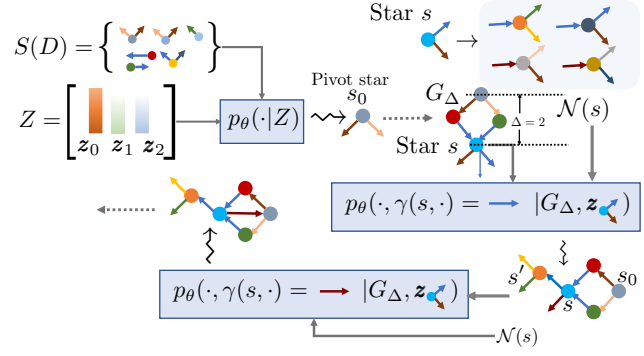
where, $Z = \{z_0, \dots, z_{\Delta_{\max}}\}$ is the set of star representations at different distances and $\gamma_{ij} = \gamma(s_i, s_j)$. Figure 1 gives an overview of our encoder decoder architecture.

Encoder q_{θ} . Given a graph $G = (V, E)$ along with the object types $T = \{t_u\}$ associated with the nodes and the relation types $R = \{r_e\}$, the encoder q_{θ} aims to characterize the scene graphs as a collection of star graphs $\{s_i\}$ and represent the graph as a collection of embeddings of these stars, i.e., $Z = \{z_{\Delta} \mid 0 \leq \Delta \leq \Delta_{\max}\}$. Such an approach is likely to preserve the spatial semantics of the visual content in a scene graph more effectively than existing graph generative models (Samanta et al., 2019; De Cao & Kipf, 2018; Li et al., 2018b) which represent a graph as a collection of node embeddings. In VARSCENE, we compute the graph representation Z in two steps. First, we compute the node and edge embeddings using a graph neural network (GNN) and then combine these to compute star embeddings.

Node and Edge Embeddings. Given a graph G and the object types $T = \{t_u\}$ and $R = \{r_e\}$, we compute the representations of the nodes x_u and edges x_e using a GNN



(a) Neural architecture of VARSCENE encoder



(b) Neural architecture of VARSCENE decoder

Figure 1. Encoder (Panel (a)): Node colors represent object types; edge colors represent relation labels. Stars are defined by hub node type and relation labels but not neighbor node types. Given a graph $G = (V, E)$ along with its object types $\{t\}$ and relation types $\{r\}$, we first decompose it into the set of stars $\{s_0, s_1, s_2, s_3\}$ where s_0 is a pivot star. Next, we feed the graph into a graph neural network ($F_\phi^0, F_\phi^{\text{edge}}, F_\phi^{\text{node}}, F_\phi^{\text{agg}}$) to derive node embeddings x_u . These node embeddings along with the stars are fed into a neural module F_ϕ^S to compute star embeddings $\{h_s\}$, which are used to generate the latent random codes $\{\zeta_s\}$. Finally, they are aggregated based on the distance from pivot star s_0 , providing embeddings $\{z_\Delta\}$. **Decoder (Panel (b)):** The base decoder p_θ and p_θ^{MMD} follow same neural architecture. Given latent codes $Z = \{z_\Delta | \Delta = 0, 1, \dots\}$, the decoder first samples the pivot star s_0 . Then it keeps adding stars at increasing distance Δ from s_0 . Given a snapshot G_Δ , the decoder p_θ draws a new star s' which is going to be connected with an existing star s at distance Δ . In this case, s' can be drawn either from the neighbor stars $\mathcal{N}(s)$ (first term in Eq. (19); the star s' in this figure, connected on the edge \rightarrow) or an already existing star (second term in Eq. (19); the star s_0 in this figure, connected on the edge \rightarrow).

proposed by Gilmer et al. (2017). More specifically, we start with the initial node representation $x_u(0)$ which transforms the feature vector t_u of the object type associated with node u using a neural network F_ϕ^0 .

$$x_u(0) = F_\phi^0(t_u) \quad (11)$$

Then, given a depth limit K , we aggregate structural information from $k = 1, \dots, K$ hops from each node into an embedding vector in a recursive manner. Formally:

$$x_{(u,v)}(k-1) = F_\phi^{\text{edge}}(x_u(k-1), x_v(k-1), r_{(u,v)}) \quad (12)$$

$$\bar{x}_u(k-1) = F_\phi^{\text{agg}}(\{x_{(u,v)}(k-1) | v \in \text{nbr}(u)\}) \quad (13)$$

$$x_u(k) = F_\phi^{\text{node}}(x_u(k-1), \bar{x}_u(k-1)). \quad (14)$$

Here, F_ϕ^{edge} and F_ϕ^{node} are modeled as universal approximator neural networks; and, F_ϕ^{agg} is a symmetric aggregator.

Star Embedding. Given S as the set of stars in G , we represent each star $s \in S$ with depth- K embeddings:

$$h_s = F_\phi^S(x_{\text{root}(s)}, \{x_{(u,v)} | v \in \text{nbr}(\text{root}(s))\}) \quad (15)$$

Next, we generate the latent representations ζ_s from a normal distribution, which parameterizes the mean and the variance using neural networks μ_ϕ and σ_ϕ , i.e.,

$$\zeta_s \sim \text{NORMAL}(\mu_\phi(h_s), \sigma_\phi^2(h_s)). \quad (16)$$

In theory, ζ_s encodes the information about the substructures up to K hops away from s and thus the collection of the representations $\{\zeta_s\}$ should be able to accurately characterize the graph. However, this may not be true in practice, as it does not explicitly specify the relative positions of the stars. To ameliorate this limitation, we first randomly draw the pivot star s_0 from set of all stars S and compute the dis-

tances $\{\Delta_s\}$ of all stars $s \in S$ from s_0 . Next, we aggregate $\{\zeta_s | \Delta_s = \Delta\}$ into a latent vector z_Δ as follows:

$$z_\Delta = F_\phi^z(\{\zeta_s | \Delta_s = \Delta\}) \text{ for } 0 \leq \Delta \leq \Delta_{\max} \quad (17)$$

where Δ_{\max} is the maximum possible distance from s_0 to any star $s \in S$. Finally, we represent the graph as $Z = \{z_0, \dots, z_{\Delta_{\max}}\}$.

Decoders p_θ and p_θ^{MMD} . The decoders p_θ and p_θ^{MMD} share the same neural network structure. They take the latent representations $Z = \{z_0, \dots, z_{\Delta_{\max}}\}$ as input and generate a new graph G . To do so, they incrementally draw new stars s_j from a database of stars $S(D)$ and connect them with previously sampled stars $s_i, i \leq j-1$. The decoders first sample a star s_0 and set $\Delta = 0$. Then, at the step $\Delta \geq 1$, they select the stars $\{s_i | \text{DISTANCE}(\text{root}(s_0), \text{root}(s_i)) = \Delta - 1\}$ having distance $\Delta - 1$ from s_0 and connect one of them with another star s_j using one of the open edges e , i.e., $\gamma(s_i, s_j) = r_e$. In this context, s_j is either drawn from the database of stars $S(D)$ or chosen from one of the existing stars in the graph.

To connect a new star s_j at an open edge e of a star s_i lying at a distance Δ from s_0 , our decoders first featurize all the potential stars $s := \langle t_u, \{r_{(u,v)} | v \in \text{nbr}(u)\} \rangle$ as $s = [t_u, \sum_v r_{(u,v)}]$; then represent the pair $[s, z_\Delta]$ using a logit; and finally, feeds the logit into a softmax distribution. This softmax distribution is used to sample the new star s' . Formally, given the snapshot of the scene graph G_Δ at step

Dataset	$ D $	Folds	$\mathbb{E}[V]$	$\mathbb{E}[E]$	$ T $	$ R $
Visual Genome	110000	81:9:9	3.23	2.31	16943	8411
Small-sized Visual Genome	124854	55:12:32	5.31	4.87	150	50
Visual Relationship Detection	7721	74:12:12	4.29	3.78	100	70

Table 1. Dataset statistics where Folds = $|D_{tr}| : |D_{dev}| : |D_{test}|$.

Δ , we have:

$$p_{\theta}(G|Z) = p_{\theta}(\{s_i\}, \{\gamma(s_i, s_j)\} | Z) \\ = \prod_{\Delta=1}^{\Delta_{\max}} \prod_{\gamma(s_i, s_j): \Delta_{s_i}=\Delta} p_{\theta}(s_j, \gamma(s_i, s_j) | G_{\Delta}, Z) \quad (18)$$

where, $\Delta_{s_i} = \text{DISTANCE}(\text{root}(s_0), \text{root}(s_i))$. Moreover, we model $p_{\theta}(s_j, \gamma(s_i, s_j) = r | G_{\Delta}, Z)$ as,

$$\frac{\text{MASK}(s_i, s_j, r) \mathbb{I}[s_j \notin G_{\Delta}] \exp(M_{\theta}(s_j, z_{\Delta}))}{\sum_{s \in S(D) \setminus G_{\Delta}} \text{MASK}(s_i, s, r) \exp(M_{\theta}(s, z_{\Delta}))} \\ + \frac{\text{MASK}(s_i, s_j, r) \mathbb{I}[s_j \in G_{\Delta}]}{|\{s : \text{MASK}(s_i, s, r) = 1, s \in G_{\Delta}\}|} \quad (19)$$

Here, M_{θ} is a neural network and the masking function MASK aims to ensure that generated scene graphs are semantically meaningful. Specifically, $\text{MASK}(s, s', r) = 1$ encodes that s' is connected to s via the relation r in at least one scene graph in the dataset D , i.e., $\text{MASK}(s, s', r) = \mathbb{I}[(s', r) \in \mathcal{N}(s)]$. The first term computes the probability when s_j is not present in G_{Δ} and is drawn from $S(D)$. The second term computes the probability when s_j is already present in G_{Δ} .

Prior Distribution p_0 . $p_0(z_0, \dots, z_{\Delta_{\max}})$ is modeled using a Normal distribution, i.e., $p_0(z_0, \dots, z_{\Delta_{\max}}) = \text{NORMAL}(0, \mathbb{I}_{\Delta_{\max}})$. Finally, we maximize $\text{ELBO}(q_{\phi}, p_{\theta})$ defined in Eq. (4) with respect to ϕ and θ to obtain $\hat{\phi}$ and $\hat{\theta}$.

4. Experiments

In this section, we provide a comprehensive evaluation of VARSCENE, addressing the following research questions: **RQ1**: How does VARSCENE compare against state-of-the-art graph generators in terms of its ability to capture the true graph distribution? **RQ2**: Does the MMD-optimized decoder p_{θ}^{MMD} provide any performance gain as compared to the the base decoder p_{θ} ? **RQ3**: How does the quality of the images created from scene graphs generated by VARSCENE compare against the quality of images created via state-of-the-art graph generators? **RQ4**: What do the generated scene graphs look like? Does the variational code space ensure smooth transitions between scene graphs? Appendix F contains additional experiments.

4.1. Experimental setup

Datasets. We use three datasets: (i) Visual Genome (VG) (Krishna et al., 2016), (ii) Small-sized Visual Genome (SVG) (Xu et al., 2017) and (iii) Visual Relationship Detection (VRD) (Lu et al., 2016). Visual Genome has a total

of 16493 object categories and 8411 relationship categories. VARSCENE can deal with this, but some recent systems report results on SVG, derived from VG by reducing the number of object and relationship categories to 150 and 50 respectively, by combining related labels and removing poor-quality labels (Xu et al., 2017; Garg et al., 2021). Table 1 summarizes details about these datasets. Appendix E provides additional details.

State-of-the-art Competitors. We compare VARSCENE against several competitive graph generative methods: two generative models which were developed for molecule synthesis, viz., (1) DeepGMG (Li et al., 2018b), (2) MolGAN (De Cao & Kipf, 2018); two domain-agnostic graph generative models, viz., (3) GraphRNN (You et al., 2018), (4) GraphGen (Goyal et al., 2020), and one graph generative model specifically aimed at synthesizing scene graphs, similar to our setup, i.e., (5) SceneGen (Garg et al., 2021).

Training, Validation and Testing. Given a dataset of scene graphs D , we split them in training (D_{tr}), validation (D_{dev}) and test (D_{test}) folds, where the exact size of these folds $|D_{\bullet}|$ varies across datasets due to their varying sizes (see Table 1). We first train the encoder q_{ϕ} and the base decoder p_{θ} using the training set D_{tr} and then train the MMD optimized decoder p_{θ}^{MMD} using the validation set D_{dev} . Finally, p_{θ}^{MMD} is evaluated using D_{test} .

Evaluation Metrics. We evaluate the quality of a generative model by measuring the similarity between the graphs generated by the trained model and the graphs in the test set. However, unlike normed-distance in Euclidean feature space, a single similarity measure cannot capture the similarity between two graphs in terms of all the properties. Hence, we resort to several similarity metrics.

Star-Sim: The *star distribution similarity* is measured in terms of the cosine similarity between the distributions of the stars present in the test and generated graphs, i.e., $\cos(\mathbb{E}_{G' \sim p_{\theta}^{\text{MMD}}} \nu(G'), \mathbb{E}_{G \sim D_{test}} \nu(G))$ with $\nu(G) = [\nu_s(G)]$ where $\nu_s(G)$ is the number of occurrences of star s in G , as defined in Eq. (8).

Edge-Sim: It is measured in terms of the cosine similarity between the distributions of the edge bigrams of the form of the tuple $e_b = \langle r_{(\bullet, u)}, t_u, r_{(u, \bullet)} \rangle$ across the test and generated graphs, i.e., $\cos(\mathbb{E}_{G' \sim p_{\theta}^{\text{MMD}}} \nu(G'), \mathbb{E}_{G \sim D_{test}} \nu(G))$ where $\nu(G) = [\nu_{e_b}(G)]$, $\nu_{e_b}(G) = \#$ occurrences of edge bigram e_b in G .

Node-Sim: It is measured in terms of the cosine similarity between the distributions of the node bigrams of the form of the tuple $n_b = \langle t_u, r_{(u, v)}, t_v \rangle$ across the test and generated graphs, i.e., $\cos(\mathbb{E}_{G' \sim p_{\theta}^{\text{MMD}}} \nu(G'), \mathbb{E}_{G \sim D_{test}} \nu(G))$ where $\nu(G) = [\nu_{n_b}(G)]$, $\nu_{n_b}(G)$ being the number of occurrences of node bigram n_b in G .

SP-K: Shortest path kernel (Vishwanathan et al., 2010).

WL-K: Weisfeiler-Lehman kernel (Vishwanathan et al.,

Model	Star-Sim	Edge-sim	Node-sim	SP-K	WL-K	NSPD-K
Visual Genome (VG)						
DeepGMG	0.69	0.46	0.15	0.01	0.09	0.01
MolGAN	0.00	0.00	0.00	0.00	0.04	0.01
GraphGen	0.66	0.37	0.11	0.00	0.03	0.01
GraphRNN	0.63	0.00	0.03	0.00	0.03	0.01
SceneGen	0.73	0.50	0.32	0.02	0.08	0.01
VARSCENE ^{unc}	0.59	0.45	0.40	0.22	0.11	0.01
VARSCENE ^{cond}	0.86	0.52	0.62	0.08	0.07	0.01
Small-sized Visual Genome (SVG)						
DeepGMG	0.80	0.64	0.49	0.10	0.35	0.05
MolGAN	0.00	0.00	0.00	0.04	0.40	0.08
GraphGen	0.52	0.64	0.37	0.05	0.21	0.04
GraphRNN	0.25	0.07	0.09	0.38	0.59	0.07
SceneGen	0.86	0.88	0.93	0.68	0.54	0.06
VARSCENE ^{unc}	0.92	0.70	0.83	1.00	0.65	0.06
VARSCENE ^{cond}	0.91	0.69	0.81	0.96	0.64	0.06
Visual Relationship Detection (VRD)						
DeepGMG	0.74	0.73	0.60	0.99	1.41	0.20
MolGAN	0.00	0.00	0.00	0.01	0.97	0.21
GraphGen	0.64	0.75	0.64	0.31	0.79	0.17
GraphRNN	0.54	0.29	0.71	0.21	0.76	0.18
SceneGen	0.81	0.94	0.95	0.60	1.12	0.21
VARSCENE ^{unc}	0.91	0.93	0.94	1.03	1.56	0.23
VARSCENE ^{cond}	0.91	0.93	0.93	1.45	1.92	0.22

Table 2. Performance of different graph generative models which include two variants of our model, *i.e.*, VARSCENE^{unc} (unconditional generation, $Z \sim p_0(\cdot)$), VARSCENE^{cond} (conditional generation, $Z \sim q_\phi(\cdot | G)$), and the baselines, *viz.*, DeepGMG (Li et al., 2018b), MolGAN (De Cao & Kipf, 2018), GraphGen (Goyal et al., 2020), GraphRNN (You et al., 2018), SceneGen (Garg et al., 2021), measured in terms of cosine similarity between stars (Star-Sim), edge bigrams (Edge-Sim), node bigram (Node-Sim), Shortest path kernel (SP-K) (Vishwanathan et al., 2010), Weisfeiler Lehman Kernel (WL-K) (Vishwanathan et al., 2010), Neighborhood Subgraph Pairwise Distance Kernel (NSPD-K) (Costa & De Grave, 2010). In all cases, the encoder q_ϕ and the base decoder p_θ were trained using the training set D_{tr} and the MMD optimized decoder was obtained using the validation set D_{dev} . In all cases, the weight of the KL divergence term in Eq. (7) was set as $\rho = 1000$. Numbers in **bold** (underline) indicate the (second) best performer.

2010).

NSPD-K: Neighborhood pairwise distance graph kernel (Costa & De Grave, 2010).

In all kernel-based metrics, we compute the similarity as $\mathbb{E}_{G' \sim p_\theta^{MMD}, G \sim D_{test}} K(G', G)$. Note that most of these measure the similarity in terms of structures as well as object and relation types. In particular, SP-Kernel and WL-Kernel take object types into account and ignore relationship types, whereas NSPD-Kernel incorporates both object and relationship types.

4.2. Results

Comparison with Baseline Graph Generators. We first address research question **RQ1** by comparing VARSCENE against baselines, in terms of the metrics defined above. We present results for two variants of VARSCENE, *viz.*, unconditional generative model VARSCENE^{unc} where $Z \sim p_0(\cdot)$

	VG		SVG		VRD	
	Edge	Node	Edge	Node	Edge	Node
VARSCENE ^{cond}	0.5268	0.6226	0.6964	0.8111	0.9372	0.9377
SceneGen, NULL $\in R$	0.0890	0.1121	0.2661	0.8207	0.4598	0.8873

Table 3. Performance measured in terms of Edge-Sim (‘Edge’) and Node-Sim (‘Node’) for VARSCENE^{cond} and the variant of SceneGen where we introduce a new relation type NULL between an isolated node and any other node. Numbers in bold indicate the best performer. With this modification, SceneGen shows rather poor performance compared to Table 2, where it is not penalized for isolated nodes. Real graphs in our data sets do not have isolated nodes.

and conditional generative model VARSCENE^{cond} where $Z \sim q_\phi(\cdot | G)$ with $G \in D_{tr}$, Table 2 summarizes the results. We make the following observations.

- (1) VARSCENE outperforms all the competitors in terms of Star-Sim, SP-Kernel and WL-Kernel with a significant boost. There is no consistent winner between VARSCENE^{unc} and VARSCENE^{cond}. Since VARSCENE^{cond} generates graphs which are structurally similar to the existing training graphs, one may expect that it is likely to mimic the underlying distribution more closely than VARSCENE^{unc}. However, this need not always be the case, if there is a large empirical distribution shift between finite training and test folds. VARSCENE^{unc} often captures such drift, as it does not look into the training graphs during new graph generation.
- (2) SceneGen outperforms the other baselines by a substantial margin in a majority of the cases. Since it is specifically aimed at scene graph generation, it is able to capture the underlying scene graph distribution more effectively than the baselines which are predominantly domain agnostic or designed for molecule synthesis.
- (3) Among the other baselines, DeepGMG shows good performance in terms of WL-Kernel or NSPD-Kernel across Visual Genome and Visual Relationship Detection datasets. Although it was primarily applied for molecule generation, its design choice is quite domain agnostic. The performance of MolGAN is extremely poor due to its design choices that are specifically aimed at molecule generation.
- (4) While VARSCENE outperforms SceneGen in terms of edge and node bigram similarity measures for the Visual Genome dataset, SceneGen appears to be the best performer for the two other datasets. However, this apparent superior performance of SceneGen may be misleading — careful investigation revealed that SceneGen generates many isolated nodes whose divergence from the test distribution cannot be captured by bigram similarity.

To shed light on the observation noted in item (4), we next introduce a dummy relation type NULL and assign this relation between any isolated node and every other node in a graph. Thus, the set of edge bigrams now include $\langle r, t, r' \rangle$ where r and r' may also assume the NULL type in addition to the existing legitimate relation types R . Similarly, the set

of node bigrams now include $\langle t, r, t' \rangle$ where r may assume the NULL type. Table 3 summarizes the effect of this augmented set of relation types, which shows that SceneGen actually performs not as well as VARSCENE^{cond}.

	VG		SVG		VRD	
	Star	Edge	Star	Edge	Star	Edge
p_{θ}^{MMD}	0.8660	0.5268	0.9182	0.6964	0.9140	0.9372
p_{θ}	0.5867	0.2588	0.7120	0.4195	0.8988	0.9339

Table 4. Performance measured in terms of Star-Sim (‘Star’) and Edge-Sim (‘Edge’) for the scene graphs generated by the MMD optimized decoder p_{θ}^{MMD} and the base decoder p_{θ} . Numbers in bold indicate the best performer. We observe that p_{θ}^{MMD} outperforms p_{θ} . In all cases, we used conditional graph generation.

Model	FID (\downarrow)	IS (\uparrow)	Precision (\uparrow)	Recall (\uparrow)
DeepGMG	9.8344	4.3566	0.9891	0.9800
MolGAN	240.3760	1.1707	0.0137	0.0986
GraphGen	13.2806	4.2802	0.9780	0.9607
GraphRNN	18.8156	4.6561	0.9432	0.9707
SceneGen	19.2634	4.0283	0.9230	0.9513
VARSCENE ^{unc}	<u>6.8591</u>	5.1505	0.9894	<u>0.9872</u>
VARSCENE ^{cond}	6.0195	<u>5.0211</u>	0.9894	0.9874

Table 5. Evaluation of various generative models by assessing the quality of the the images obtained from the corresponding scene graphs. Performance is measured in terms of Fréchet Inception Distance (Heusel et al., 2017), Inception Score (Salimans et al., 2016) and Precision & Recall (Sajjadi et al., 2018). Numbers in bold (underline) show the best (second best) performer. \downarrow means smaller is better; \uparrow means larger is better.

Effect of MMD-Optimized Decoder. Next, we address research question RQ2. Specifically, we compare the quality of the scene graphs generated by our MMD optimized decoder p_{θ}^{MMD} with the base decoder p_{θ} . Table 4 summarizes the results, which shows that our MMD-optimized decoder is able to mimic the true distribution of stars, as well as edge bi-grams, more accurately than the base decoder. Other metrics are reported in Appendix F.

Image Quality. Thus far, we have assessed the quality of the generated graphs, rather than the images that might be generated from them. Here, we address research question RQ3 by evaluating the quality of the images corresponding to the synthesized scene graphs. Specifically, we use sg2im², a scene graph to image generation system (Johnson et al., 2018) that comes pre-trained on the Small-sized Visual Genome dataset. Similar to Garg et al. (2021), we evaluate the quality of the images using Fréchet Inception Distance (Heusel et al., 2017), Inception Score (Salimans et al., 2016), Precision and Recall (Sajjadi et al., 2018). Both Inception Score (IS) and Fréchet Inception Distance (FID) aim to provide proxies for qualitative human evaluation. They assess the images from two perspectives: (a) the quality of the images via the output of an image classification model like InceptionV3 and (b) diversity across

²<https://github.com/google/sg2im>

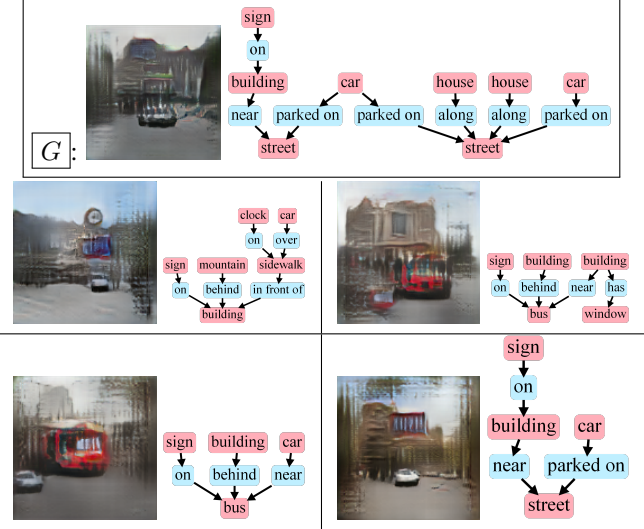


Figure 2. Four scene graphs and the corresponding images, generated using $G \sim p_{\theta}^{\text{MMD}}(\bullet | Z)$, where $Z \sim q_{\phi}(\bullet | \boxed{G})$. Here, \boxed{G} is the graph used for conditioning, which is chosen from Small-sized Visual Genome dataset. The images corresponding to the scene graphs G' are close to the image corresponding to \boxed{G} .

the set of the images. We formally define FID, IS, Precision and Recall in Appendix E. Table 5 summarizes the results for SVG, which shows that, (1) the images corresponding to the scene graphs provided by the variants of VARSCENE have substantially better quality than all other baselines; (2) DeepGMG outperforms the other baselines in majority of the cases.

Visualization of Generated Scene Graphs. Finally, we address research question RQ4, where we visualize the graphs G' generated by conditioning on an existing graph G . Specifically, we draw $G' \sim p_{\theta}^{\text{MMD}}(\bullet | Z)$, where $Z \sim q_{\phi}(\bullet | G)$. Figure 2 provides some samples, which show that VARSCENE is able to generate graphs which provides similar images to the image corresponding to scene graph \boxed{G} used for conditioning.

5. Conclusion

We presented VARSCENE, a variational autoencoder tailored to synthesize scene graphs from a seed set of ‘gold’ graphs. VARSCENE uses a novel sampling vocabulary of stars annotated with node and edge types from large vocabularies typical in image collections. It directly minimizes distributional discrepancies for features observed in the gold and generated graphs. These two strengths enable VARSCENE to outperform recent graph generators on three data sets. Additional experiments provide further insight into VARSCENE and prior systems. It would be interesting to use third-party knowledge (e.g., knowledge graphs) to enhance to semantic structure of the generated scenes.

References

- Ashual, O. and Wolf, L. Specifying object attributes and relations in interactive scene generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4561–4569, 2019.
- Blum, L. C. and Reymond, J.-L. 970 million druglike small molecules for virtual screening in the chemical universe database GDB-13. *J. Am. Chem. Soc.*, 131:8732, 2009.
- Caesar, H., Uijlings, J., and Ferrari, V. Coco-stuff: Thing and stuff classes in context. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1209–1218, 2018.
- Casanova, A., Drozdal, M., and Romero-Soriano, A. Generating unseen complex scenes: are we there yet? *arXiv preprint arXiv:2012.04027*, 2020.
- Chang, X., Ren, P., Xu, P., Li, Z., Chen, X., and Hauptmann, A. Scene graphs: A survey of generations and applications. *arXiv preprint arXiv:2104.01111*, 2021.
- Costa, F. and De Grave, K. Fast neighborhood subgraph pairwise distance kernel. In *ICML*, 2010.
- Dai, B., Zhang, Y., and Lin, D. Detecting visual relationships with deep relational networks. In *Proceedings of the IEEE conference on computer vision and Pattern recognition*, pp. 3076–3086, 2017.
- De Cao, N. and Kipf, T. Molgan: An implicit generative model for small molecular graphs. *arXiv preprint arXiv:1805.11973*, 2018.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Dhamo, H., Farshad, A., Laina, I., Navab, N., Hager, G. D., Tombari, F., and Rupperecht, C. Semantic image manipulation using scene graphs. In *CVPR*, pp. 5213–5222, 2020.
- Garg, S., Dhamo, H., Farshad, A., Musatian, S., Navab, N., and Tombari, F. Unconditional scene graph generation. *ICCV*, 2021. URL <https://arxiv.org/abs/2108.05884>.
- Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O., and Dahl, G. E. Neural message passing for quantum chemistry. In *International conference on machine learning*, pp. 1263–1272. PMLR, 2017.
- Goyal, N., Jain, H. V., and Ranu, S. Graphgen: A scalable approach to domain-agnostic labeled graph generation. *Proceedings of The Web Conference*, Apr 2020. doi: 10.1145/3366423.3380201. URL <http://dx.doi.org/10.1145/3366423.3380201>.
- Gu, J., Zhao, H., Lin, Z., Li, S., Cai, J., and Ling, M. Scene graph generation with external knowledge and image reconstruction. In *CVPR*, June 2019.
- Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and Hochreiter, S. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *arXiv preprint arXiv:1706.08500*, 2017.
- Johnson, J., Krishna, R., Stark, M., Li, L.-J., Shamma, D. A., Bernstein, M. S., and Fei-Fei, L. Image retrieval using scene graphs. *CVPR*, pp. 3668–3678, 2015. URL <http://hci.stanford.edu/publications/2015/scenegraphs/JohnsonCVPR2015.pdf>.
- Johnson, J., Gupta, A., and Fei-Fei, L. Image generation from scene graphs. In *CVPR*, pp. 1219–1228, 2018.
- Klawonn, M. and Heim, E. Generating triples with adversarial networks for scene graph construction. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- Krishna, R., Zhu, Y., Groth, O., Johnson, J., Hata, K., Kravitz, J., Chen, S., Kalantidis, Y., Li, L.-J., Shamma, D. A., Bernstein, M., and Fei-Fei, L. Visual genome: Connecting language and vision using crowdsourced dense image annotations. 2016. URL <https://arxiv.org/abs/1602.07332>.
- Li, W., Zhang, P., Zhang, L., Huang, Q., He, X., Lyu, S., and Gao, J. Object-driven text-to-image synthesis via adversarial training. In *CVPR*, June 2019.
- Li, Y., Ouyang, W., Zhou, B., Wang, K., and Wang, X. Scene graph generation from objects, phrases and region captions. In *Proceedings of the IEEE international conference on computer vision*, pp. 1261–1270, 2017.
- Li, Y., Ouyang, W., Zhou, B., Shi, J., Zhang, C., and Wang, X. Factorizable net: an efficient subgraph-based framework for scene graph generation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 335–351, 2018a.
- Li, Y., Vinyals, O., Dyer, C., Pascanu, R., and Battaglia, P. Learning deep generative models of graphs. *arXiv preprint arXiv:1803.03324*, 2018b.
- Liao, R., Li, Y., Song, Y., Wang, S., Nash, C., Hamilton, W. L., Duvenaud, D., Urtasun, R., and Zemel, R. S. Efficient graph generation with graph recurrent attention networks. *CoRR*, abs/1910.00760, 2019. URL <http://arxiv.org/abs/1910.00760>.

- Lu, C., Krishna, R., Bernstein, M., and Fei-Fei, L. Visual relationship detection with language priors. In *ECCV*, 2016.
- Maheshwari, P., Chaudhry, R., and Vinay, V. Scene graph embeddings using relative similarity supervision. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 2328–2336, 2021.
- Milewski, V., Moens, M.-F., and Calixto, I. Are scene graphs good enough to improve image captioning? 2020. URL <https://arxiv.org/abs/2009.12313>.
- Rupp, M., Tkatchenko, A., Müller, K.-R., and von Lilienfeld, O. A. Fast and accurate modeling of molecular atomization energies with machine learning. *Physical Review Letters*, 108:058301, 2012.
- Sajjadi, M. S., Bachem, O., Lucic, M., Bousquet, O., and Gelly, S. Assessing generative models via precision and recall. *Advances in Neural Information Processing Systems*, 2018. URL <https://arxiv.org/abs/1806.00035>.
- Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A., and Chen, X. Improved techniques for training gans. *Advances in Neural Information Processing Systems*, 29:2234–2242, 2016.
- Samanta, B., DE, A., Jana, G., Chattaraj, P. K., Ganguly, N., and Rodriguez, M. G. Nevae: A deep generative model for molecular graphs. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):1110–1117, 2019. doi: 10.1609/aaai.v33i01.33011110.
- Schroeder, B. and Tripathi, S. Structured query-based image retrieval using scene graphs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2020.
- Simonovsky, M. and Komodakis, N. Graphvae: Towards generation of small graphs using variational autoencoders. *CoRR*, abs/1802.03480, 2018. URL <http://arxiv.org/abs/1802.03480>.
- Sterling, T. and Irwin, J. J. Zinc 15 – ligand discovery for everyone. *Journal of Chemical Information and Modeling*, 55(11):2324–2337, 2015. doi: 10.1021/acs.jcim.5b00559. URL <https://doi.org/10.1021/acs.jcim.5b00559>. PMID: 26479676.
- Sun, W. and Wu, T. Learning layout and style reconfigurable gans for controllable image synthesis. *IEEE transactions on pattern analysis and machine intelligence*.
- Tabibian, B., Gomez, V., De, A., Schölkopf, B., and Rodriguez, M. G. On the design of consequential ranking algorithms. In *Conference on Uncertainty in Artificial Intelligence*, pp. 171–180. PMLR, 2020.
- Tan, F., Feng, S., and Ordonez, V. Text2scene: Generating compositional scenes from textual descriptions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6710–6719, 2019.
- Tang, H., Xu, D., Yan, Y., Torr, P. H., and Sebe, N. Local class-specific and global image-level generative adversarial networks for semantic-guided scene generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7870–7879, 2020a.
- Tang, K., Niu, Y., Huang, J., Shi, J., and Zhang, H. Unbiased scene graph generation from biased training. In *CVPR*, pp. 3716–3725, 2020b.
- Teney, D., Liu, L., and van Den Hengel, A. Graph-structured representations for visual question answering. In *CVPR*, pp. 1–9, 2017.
- Todorov, E. Efficient computation of optimal actions. *Proceedings of the national academy of sciences*, 106(28): 11478–11483, 2009.
- Tseng, H.-Y., Lee, H.-Y., Jiang, L., Yang, M.-H., and Yang, W. Retrievegan: Image synthesis via differentiable patch retrieval. In *ECCV*, pp. 242–257. Springer, 2020.
- Vishwanathan, S. V. N., Schraudolph, N. N., Kondor, R., and Borgwardt, K. M. Graph kernels. *Journal of Machine Learning Research*, 11:1201–1242, 2010.
- Wang, S., Wang, R., Yao, Z., Shan, S., and Chen, X. Cross-modal scene graph matching for relationship-aware image-text retrieval. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 1508–1517, 2020.
- Williams, R. J. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256, 1992.
- Xu, D., Zhu, Y., Choy, C. B., and Fei-Fei, L. Scene graph generation by iterative message passing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- Xu, T., Zhang, P., Huang, Q., Zhang, H., Gan, Z., Huang, X., and He, X. Attngan: Fine-grained text to image generation with attentional generative adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1316–1324, 2018.
- Yang, J., Lu, J., Lee, S., Batra, D., and Parikh, D. Graph rnn for scene graph generation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018.

- You, J., Ying, R., Ren, X., Hamilton, W. L., and Leskovec, J. Graphrnn: A deep generative model for graphs. *CoRR*, abs/1802.08773, 2018. URL <http://arxiv.org/abs/1802.08773>.
- Zareian, A., Karaman, S., and Chang, S.-F. Bridging knowledge graphs to generate scene graphs. In *ECCV*, 2020.
- Zhang, J., Elhoseiny, M., Cohen, S., Chang, W., and Elgammal, A. Relationship proposal networks. In *CVPR*, pp. 5678–5686, 2017.

VARSCENE: A Deep Generative Model for Realistic Scene Graph Synthesis (Supplementary Material)

A. Broader Impact

Generating realistic and complex scenes is a key task in image retrieval, image editing, question-answering, etc. In this work, we show that it is possible to sample such scene graphs from a distribution estimated from a set of real scene graphs, without explicitly inspecting training images. Any bias (e.g., gender, race) in the training scene graphs may propagate to the synthetic images produced by our system. Our model has sufficient flexibility to attempt to reduce or mitigate such bias. Specifically, one can replace the first term quantifying MMD in our objective function in Eq. (7) with a suitable alternative which can mitigate the biases in the generated graphs.

B. Limitations

There are a few potential limitations of VARSCENE, which we leave to be addressed as ongoing and future work.

1. The MMD computation depends heavily on the choice of graph kernel. It is likely that the decoder p_{θ}^{MMD} optimized for one specific choice of graph kernel may not show good performance for other types of kernels. Multitask objectives may be of interest here.
2. Currently, our model is trained in two stages. This increases the time and complexity of training. It would be of interest to train a single stage variational autoencoder, which would be also optimized for a given MMD measure (or combination of MMD measures).
3. Our generative mechanism based on star agglomeration allows us to incorporate salient semantic structures from training scenes alone. Knowledge graphs representing facts (such as [WikiData](#)) or [common sense](#) may provide valuable additional “realism prior” or “world bias” for better training with less supervision.

C. Related Work

Our work is related to graph generative models and scene graph extraction from images.

C.1. Graph Generative Models

Driven by the success of text and image generative models, there is much recent interest in deep generative models for graphs, primarily designed for molecule discovery, *e.g.*, GraphVAE (Simonovsky & Komodakis, 2018), GraphRNN (You et al., 2018), NeVAE (Samanta et al., 2019), MolGAN (De Cao & Kipf, 2018), GRAN (Liao et al., 2019), etc. However, the edge labels in molecular graphs are restricted by the types of bonds present in chemical compounds, which does not exceed 10. On the other hand, the node labels correspond to the unique elements, which is typically <15 in the ZINC (Sterling & Irwin, 2015) and QM9 (Blum & Raymond, 2009; Rupp et al., 2012) datasets. In contrast, the distinct number of node and edge labels in scene graph datasets are significantly higher than those in molecular graphs. *E.g.*, Visual Genome (Krishna et al., 2016) has 75,729 unique node labels and 40,480 unique edge labels. Even the smaller Visual Relationship Detection dataset (Lu et al., 2016) has 100 unique node labels and 70 unique edge labels. Apart from the smaller number of node and edge labels, molecules and protein compounds follow well-defined graph semantics, *e.g.*, valency rules of atoms, non-existence of strained bridges and rings, and occurrence of specific motifs in molecules. These place constraints on the set of plausible graphs, which can be exploited by generative models for molecular graphs. Codifying knowledge about the visual world to aid the training process may not be straightforward in the context of scene graphs.

C.2. Scene Graph Extraction from Images

A graph where the nodes represent objects with edges being relationships between them encapsulates the semantic content of an image in a succinct manner. We refer to the task of constructing such a graph from the visual modality as scene graph extraction (although it is sometimes called ‘generation’, which can be confusing). Literature associated with this (essentially inference) task is vast (Xu et al., 2017; Yang et al., 2018; Li et al., 2018a; Tang et al., 2020b), with alternatives characterized by the modeling approaches — *e.g.*, do they use conditional random fields or recurrent/convolutional network variations, are the models trained in an adversarial manner, and so on. Most of these models’ inputs are derived from the current image alone — these would be visual features computed from regions of interest that help detect the object type or identify the relationship between previously identified objects (Li et al., 2017; Dai et al., 2017; Zhang et al., 2017; Klawonn & Heim, 2018). More recently, inference methods aided by external knowledge have been developed. This knowledge can take the form of word embeddings (Lu et al., 2016) and knowledge graphs (Zareian et al., 2020). The task considered in the current paper is that of scene graph *synthesis*, where we do not have access to the image modality. In the absence of this visual information, toward synthesizing realistic and plausible scene graphs, we mirror the existing practice of utilizing information from the linguistic domain in the form of word embeddings of the objects and relationship names.

C.3. Image Generation from Scene Graphs

While there has been progress in the field of image generation, the end goal of achieving rich, diverse and complex scenes of high quality still needs work (Caesar et al., 2018). In the current paper, we propose that an existing collection of scene graphs can be utilized to train a generative model of such graphs. Sampling from such a model provides us with the semantic content for a novel scene. The sampled graph can then be used as an input into a conditional image generation model (Johnson et al., 2018; Ashual & Wolf, 2019; Tseng et al., 2020). Conditional image generation models that operate on other input modalities also exist. For example, textual phrases (Xu et al., 2018; Tan et al., 2019), semantic segmentation masks (Tang et al., 2020a), and object positions and layouts (Sun & Wu).

D. Additional Details of VARSCENE

We provide more details about VARSCENE as follows.

t and r . Given an object type t and relation type r , we use BERT embedding (Devlin et al., 2018) to obtain the corresponding representations t and r . Here, we used the code from <https://huggingface.co/sentence-transformers/paraphrase-MiniLM-L6-v2> and thus have $\text{dimension}(t) = \text{dimension}(r) = 384$.

Specifications of the Encoder. We summarize the components of the encoder as follows:

1. The GNN module of q_ϕ contains four networks: F_ϕ^0 , F_ϕ^{edge} , F_ϕ^{agg} and F_ϕ^{node} , which follow similar architecture as in (Gilmer et al., 2017). We summarize their architecture in details as follows:
 - F_ϕ^0 : It consists of a single layer neural network with linear activation function which outputs a 64 dimensional node features.
 - F_ϕ^{edge} : It is a single layer neural network with linear activation function. It takes the pair of node features and the corresponding relation type of the edge between them as input and then outputs a 64 dimensional edge embedding vector.
 - F_ϕ^{agg} : It is a sum-pool aggregator.
 - F_ϕ^{node} : It is a GRU which sequentially takes $x_u(k-1)$ and $\bar{x}_u(k-1)$ as input and outputs a 64 dimensional node embedding.
2. F_ϕ^S is kept the same as F_ϕ^{node} . It is a GRU which is supplied only one element, the concatenation of $x_{\text{root}(s)}$ and $\text{sum}(\{x_{(\text{root}(s),v)} \mid v \in \text{nbr}(\text{root}(s))\})$. It outputs a 64 dimensional vector.
3. μ_ϕ in Eq. (16) is a two layer neural network which consists of a 128 dimensional hidden layer, ReLU activations and outputs a 64 dimensional mean vector.
4. σ_ϕ in Eq. (16) is a two layer neural network which consists of a 128 dimensional hidden layer, ReLU activations and outputs a 64x64 matrix.
5. F_ϕ^z in Eq. (17) is a sum aggregator.

Specification of the Decoder. Here, M_θ in Eq. (19) is a two layer neural network which consists of a 32 dimensional hidden layer, ReLU activations and outputs a 16 dimensional mean vector.

E. Details About Experiments

E.1. Pre-processing of Datasets

For the SVG dataset, we start with pre-processed graphs that were provided with the code of (Garg et al., 2021), that we obtained from the authors. For VG and VRD datasets, we construct graphs directly from the json formatted scene graphs. After obtaining these graphs, we split them into weakly connected components, i.e., components in which each node is reachable from all other nodes considering edges to be undirected. This step is done for all datasets. For the VG dataset, the total number of graphs after splitting exceeded $\sim 800K$. Due to the large number of graphs, we randomly sampled 90K train graphs, 10K validation graphs and 10K test graphs for the VG dataset. For the SVG dataset, train-test split was provided in (Garg et al., 2021), hence we used the same set of test graphs and randomly sampled graphs from the train set to create validation graphs. For the VRD dataset, train-test graphs were provided, but the number of graphs was quite small. Due to this, we randomly sampled graphs from the train and test sets to create the validation graphs.

E.2. Implementation Details of Baselines

DeepGMG. We used the version released with the official PyTorch implementation of Goyal et al. (2020)³.

MolGAN. We used an open-source PyTorch implementation (De Cao & Kipf, 2018)⁴.

GraphRNN. We used the label adapted version released with the official PyTorch implementation of Goyal et al. (2020)⁵.

GraphGen. We used the official PyTorch implementation (Goyal et al., 2020)⁶.

SceneGen. We obtained the code from the authors (Garg et al., 2021).

We summarize the total number of parameters in our model and the baselines in Table 6, which shows that our model uses minimum number of parameters and still outperforms all the baselines. Moreover, note that, our number of parameters is independent of the dataset size, whereas the complexity of baseline models widely varies across most of the datasets.

	Parameter Count		
	VG	SVG	VRD
DeepGMG	9,096,436	653,771	641,836
MolGAN	623,102,433	18,895,885	11,670,574
GraphRNN	19,561,966	529,421	510,150
GraphGen	25,056,650	2,865,855	2,781,760
SceneGen	277,012	222,781	228,019
VARSCENE	221,249	221,249	221,249

Table 6. Number of parameters in each trainable model.

E.3. Hyperparameters

VARSCENE. For the base decoder (p_θ), we trained the model for 1000 epochs. For the optimized decoder (p_θ^{MMD}), we again trained for 1000 epochs. The batch size $|\mathcal{B}| = 1024$. For training p_θ , the learning rate was kept as 10^{-4} . For p_θ^{MMD} , the learning rate was kept as 10^{-3} . Moreover, we used the Adam optimizer with a 10^{-5} weight decay factor.

SceneGen. We trained their model for 300 epochs. The batch size was $|\mathcal{B}| = 64$. A step learning rate scheduler was used for both nodes and edges, with initial learning rate 10^{-3} , final learning rate 10^{-4} and decay factor 0.95. The Adam optimizer was used for updates. These hyper-parameters were the default values.

DeepGMG, GraphRNN, GraphGen. We trained the model for 10000 epochs. Batch size $|\mathcal{B}| = 2048$. Initial learning rate was 0.003 and decayed by a factor of 0.3 after epochs 100, 200, 400 and 800. The Adam optimizer was used for parameter updates. These hyper-parameters were the default values.

MolGAN. For VRD and SVG datasets, we trained the model for 12500 epochs. For the VG dataset, we trained the model for 800 epochs (since the model took a lot of time for training). Batch size $|\mathcal{B}| = 16$. Learning rate was kept as 10^{-4} . The Adam optimizer was used for parameter updates with $\beta_1 = 0.5$, $\beta_2 = 0.999$. These hyper-parameters were the default values.

³<https://github.com/idea-iitd/graphgen>

⁴<https://github.com/yongqyu/MolGAN-pytorch>

⁵<https://github.com/idea-iitd/graphgen>

⁶<https://github.com/idea-iitd/graphgen>

E.4. Infrastructure details

We performed all experiments on an Intel Xeon server with 1 TB RAM running Ubuntu 18.04.5 LTS. Our code ran on any of a TITAN RTX GPU with 24 GB RAM and two TITAN X (Pascal) GPUs with 12 GB RAM each.

E.5. Evaluation metrics for Image Synthesis

Let \mathcal{I}_r be a set of real images, and \mathcal{I}_g be a set of images that were output by an image generation model. Metrics that compute the quality of the generated set typically use a feature representation \mathbf{x}_r and \mathbf{x}_g respectively for elements of these sets. Typically, the activation values of the penultimate layer of a pre-trained image classification model, most often the Inception v3, is used for this purpose. Let \mathbf{X}_g be the feature representation matrix for the set of generated images, and the corresponding matrix for real images is given by \mathbf{X}_r . The metrics below utilize these feature representations in different ways, each of which assumes that \mathcal{I}_g and \mathcal{I}_r are of the same size.

Fréchet Inception Distance. (Heusel et al., 2017): The FID is a distance between two multivariate Gaussians each of which is represented by a mean and covariance matrix: $(\mathbf{m}_r, \mathbf{C}_r)$ and $(\mathbf{m}_g, \mathbf{C}_g)$. The FID is given by $\|\mathbf{m}_r - \mathbf{m}_g\|_2^2 + \text{trace}(\mathbf{C}_r + \mathbf{C}_g - 2(\mathbf{C}_r \times \mathbf{C}_g)^{0.5})$. In the context of evaluating the output of models that generate images, the summary statistics (\mathbf{m}, \mathbf{C}) are computed over the feature representation matrices \mathbf{X} as defined above.

Inception Score. (Salimans et al., 2016): The IS compares the marginal class label distribution over a collection of real images to the conditional label distribution given a set of generated images. It is defined as $\exp(\mathbb{E}_{I \sim \mathcal{I}_g} [KL(p(y|I) || p(y))])$. I here represents an individual output of the image generation model. The vector y is taken to be the set of labels of an image classification model, again typically taken to be the Inception v3 trained on ImageNet. The intuition is that a good image generation output will have $p(y|I)$ with low entropy indicating the presence of recognisable objects. But at a set level, we would like a set of varied images, and this is captured by the KL term.

Precision & Recall. (Sajjadi et al., 2018): The precision is given $P(\mathbf{X}_r, \mathbf{X}_g) = \mathbb{E}_{\mathbf{x} \sim \mathbf{X}_g} f(\mathbf{x}, \mathbf{X}_r)$, and recall $R(\mathbf{X}_r, \mathbf{X}_g) = \mathbb{E}_{\mathbf{x} \sim \mathbf{X}_r} f(\mathbf{x}, \mathbf{X}_g)$. The binary function $f(\mathbf{x}, \mathbf{X})$ returns if the sample \mathbf{x} lies within the volume of the set \mathbf{X} . In the current context, Recall represents the fraction of the real images that were generated. While Precision is the fraction of generated images that are similar to real images.

F. Additional Experiments

Effect of MMD-Optimized Method. Table 4 compared p_{θ}^{MMD} against p_{θ} using Star-Sim and Edge-Sim measurements. Here, in Table 7, we present additional similarity and kernel measurements.

Visual Genome (VG)						
Model	Star-Sim	Edge-Sim	Node-Sim	SP-Kernel	WL-Kernel	NSPD-Kernel
p_{θ}^{MMD}	0.8660	0.5268	0.6226	0.0840	0.0779	0.0135
p_{θ}	0.5867	0.2588	0.2212	0.1183	0.1105	0.0143
Visual Relationship Detection (VRD)						
p_{θ}^{MMD}	0.9140	0.9372	0.9377	1.4588	1.9218	0.2275
p_{θ}	0.8988	0.9339	0.9306	0.9426	1.4560	0.2156
Small-sized Visual Genome (SVG)						
p_{θ}^{MMD}	0.9182	0.6964	0.8111	0.9621	0.6417	0.0633
p_{θ}	0.7120	0.4195	0.5616	1.9098	0.6944	0.0586

Table 7. Performance for the scene graphs generated by the MMD optimized decoder p_{θ}^{MMD} and the base decoder p_{θ} . Numbers in bold indicate the best performer. We observe that p_{θ}^{MMD} outperforms p_{θ} in a majority of cases. In all cases, we used conditional graph generation.

Effect of Optimizing Other Similarity Metrics. In Eq. (7), we focus on minimizing MMD between the generated and test scene graphs in terms of the RBF kernel (3). Here, we replace this MMD metric with negative of the six similarity measures used for evaluation, *i.e.*, Star-Sim, Edge-Sim, Node-Sim, SP-Kernel, WL-Kernel and NSPD-Kernel. Tables 8–10 summarize the results for different datasets, and show that the superiority of our method is consistent across them.

Visual Genome (VG)						
Model	Star-Sim	Edge-Sim	Node-Sim	SP-Kernel	WL-Kernel	NSPD-Kernel
DeepGMG	0.6996	0.4668	0.1579	0.0137	0.0994	0.0144
MolGAN	0.0000	0.0000	0.0000	0.0002	0.0439	0.0131
GraphGen	0.6696	0.3760	0.1145	0.0037	0.0378	0.0126
GraphRNN	0.6352	0.0000	0.0332	0.0003	0.0339	0.0121
SceneGen	0.7378	0.5050	0.3275	0.0276	0.0804	0.0159
VARSCENE ^{unc}	0.5963	0.4576	0.4055	0.2293	0.1176	0.0147
VARSCENE ^{cond}	0.8660	0.5268	0.6226	0.0840	0.0779	0.0135
VARSCENE ^{unc} (Edge)	0.6999	0.3892	0.3640	0.0789	0.0806	0.0132
VARSCENE ^{cond} (Edge)	0.5930	0.1677	0.2789	0.1129	0.0899	0.0132
VARSCENE ^{unc} (Node)	0.7702	0.5497	0.4433	0.0558	0.0680	0.0129
VARSCENE ^{cond} (Node)	0.7393	0.3636	0.3726	0.0622	0.0674	0.0126
VARSCENE ^{unc} (WL)	0.7350	0.4709	0.4039	0.0865	0.0753	0.0133
VARSCENE ^{cond} (WL)	0.6819	0.4088	0.3422	0.0989	0.0784	0.0131
VARSCENE ^{unc} (SP)	0.6134	0.4246	0.3240	0.1334	0.0868	0.0135
VARSCENE ^{cond} (SP)	0.5883	0.3193	0.2960	0.1720	0.0985	0.0136
VARSCENE ^{unc} (NSPD)	0.6704	0.4831	0.3934	0.1399	0.0887	0.0134
VARSCENE ^{cond} (NSPD)	0.5932	0.3046	0.3766	0.3390	0.1121	0.0136
VARSCENE ^{unc} (Node, Edge)	0.7880	0.6319	0.5571	0.0672	0.0743	0.0133
VARSCENE ^{cond} (Node, Edge)	0.7592	0.3620	0.4573	0.0816	0.0745	0.0133
VARSCENE ^{unc} (Star, Node, Edge)	0.8013	0.5094	0.5711	0.1065	0.0887	0.0145
VARSCENE ^{cond} (Star, Node, Edge)	0.6813	0.2477	0.4131	0.1476	0.0957	0.0142
VARSCENE ^{unc} (SP, Node, Edge)	0.8066	0.5606	0.5476	0.0647	0.0746	0.0135
VARSCENE ^{cond} (SP, Node, Edge)	0.7344	0.3062	0.4343	0.0774	0.0749	0.0132

Table 8. Performance of different graph generative models (analogous to Table 2) for different optimization objectives in the second stage of our learning for Visual Genome dataset. The metrics against VARSCENE, in parentheses, imply that in Eqn. (7) we replace $\widehat{\text{MMD}}$ with negative of these similarity metrics. Here VARSCENE (Metric) means that in Eq. 7 we replace $\widehat{\text{MMD}}$ with $-\text{Metric}$. Moreover, for last few rows, we maximize the sum of multiple similarities. For example, VARSCENE (Metric₁, Metric₂) indicates that in Eq. 7, we replace $\widehat{\text{MMD}}$ with $-(\text{Metric}_1 + \text{Metric}_2)$. We write Star for Star-Sim, Edge for Edge-Sim, Node for Node-Sim, SP for SP-Kernel, WL for WL-Kernel and NSPD for NSPD-Kernel. Recall that VARSCENE as-it-is stands for optimizing the discrepancy between the star distributions of the generated and the test set.

Small-sized Visual Genome (SVG)						
Model	Star-Sim	Edge-Sim	Node-Sim	SP-Kernel	WL-Kernel	NSPD-Kernel
DeepGMG	0.8097	0.6401	0.4961	0.1092	0.3589	0.0538
MolGAN	0.0000	0.0000	0.0000	0.0496	0.4084	0.0879
GraphGen	0.5223	0.6461	0.3713	0.0545	0.2117	0.0432
GraphRNN	0.2590	0.0775	0.0920	0.3833	0.5991	0.0705
SceneGen	0.8644	0.8828	0.9311	0.6876	0.5480	0.0663
VARSCENE ^{unc}	0.9235	0.7049	0.8335	1.0046	0.6556	<u>0.0685</u>
VARSCENE ^{cond}	0.9182	0.6964	0.8111	0.9621	0.6417	0.0633
VARSCENE ^{unc} (Edge)	<u>0.9400</u>	0.7451	<u>0.8620</u>	0.7177	0.6100	0.0673
VARSCENE ^{cond} (Edge)	0.9302	<u>0.7538</u>	0.8362	0.7140	0.6078	0.0624
VARSCENE ^{unc} (Node)	0.9351	0.7452	0.8599	0.7193	0.6135	0.0677
VARSCENE ^{cond} (Node)	0.9210	0.7446	0.8287	0.6726	0.6058	0.0623
VARSCENE ^{unc} (WL)	0.9469	0.7268	0.8546	<u>0.9840</u>	0.6342	0.0679
VARSCENE ^{cond} (WL)	0.9300	0.7043	0.8223	0.9395	0.6005	0.0614
VARSCENE ^{unc} (SP)	0.9328	0.7246	0.8576	0.8961	<u>0.6444</u>	<u>0.0691</u>
VARSCENE ^{cond} (SP)	0.9286	0.7171	0.8280	0.8721	0.6170	0.0636
VARSCENE ^{unc} (NSPD)	0.9349	0.7231	0.8498	0.8096	0.6215	0.0669
VARSCENE ^{cond} (NSPD)	0.9334	0.7204	0.8289	0.8038	0.6266	0.0628
VARSCENE ^{unc} (Node, Edge)	0.9254	0.6918	0.8443	0.9661	0.6373	0.0679
VARSCENE ^{cond} (Node, Edge)	0.9160	0.7044	0.8128	0.7688	0.5996	0.0621
VARSCENE ^{unc} (Star, Node, Edge)	0.9308	0.7039	0.8346	0.9234	0.6397	0.0679
VARSCENE ^{cond} (Star, Node, Edge)	0.9204	0.7028	0.8127	0.8861	0.6196	0.0627
VARSCENE ^{unc} (SP, Node, Edge)	0.9394	0.7322	0.8577	0.8346	0.6329	0.0684
VARSCENE ^{cond} (SP, Node, Edge)	0.9332	0.7206	0.8366	0.8087	0.6116	0.0629

Table 9. Performance of different graph generative models (analogous to Table 8) for different optimization objectives in the second stage of our learning for Small-sized Visual Genome dataset.

Visual Relationship Detection (VRD)						
Model	Star-Sim	Edge-Sim	Node-Sim	SP-Kernel	WL-Kernel	NSPD-Kernel
DeepGMG	0.7459	0.7338	0.6035	0.9977	1.4125	0.2038
MolGAN	0.0000	0.0000	0.0001	0.0149	0.9721	0.2121
GraphGen	0.6478	0.7586	0.6431	0.3184	0.7964	0.1770
GraphRNN	0.5455	0.2941	0.7163	0.2125	0.7639	0.1888
SceneGen	0.8112	0.9486	0.9552	0.6069	1.1274	0.2130
VARSCENE ^{unc}	0.9194	0.9374	0.9403	1.0327	1.5619	0.2301
VARSCENE ^{cond}	0.9140	0.9372	0.9377	1.4588	1.9218	0.2275
VARSCENE ^{unc} (Edge)	0.9216	0.9290	0.9363	1.1150	1.6175	0.2349
VARSCENE ^{cond} (Edge)	0.9177	0.9219	0.9343	1.4194	1.9164	0.2290
VARSCENE ^{unc} (Node)	0.9063	0.9296	0.9240	1.0671	1.5865	0.2321
VARSCENE ^{cond} (Node)	0.9062	<u>0.9387</u>	0.9268	1.2975	1.7934	0.2238
VARSCENE ^{unc} (WL)	0.8989	0.9281	0.9217	1.1509	1.6318	0.2345
VARSCENE ^{cond} (WL)	0.8979	0.9249	0.9251	1.4902	1.9244	0.2294
VARSCENE ^{unc} (SP)	0.9122	0.9317	0.9328	1.1432	1.6431	<u>0.2379</u>
VARSCENE ^{cond} (SP)	0.9108	0.9290	0.9354	<u>1.5389</u>	<u>1.9726</u>	0.2334
VARSCENE ^{unc} (NSPD)	0.8893	0.9011	0.9044	1.0247	1.5454	0.2276
VARSCENE ^{cond} (NSPD)	0.8898	0.8858	0.9065	1.2732	1.7864	0.2204
VARSCENE ^{unc} (Node, Edge)	0.8941	0.9254	0.9190	1.1111	1.6225	0.2332
VARSCENE ^{cond} (Node, Edge)	0.9029	0.9298	0.9261	1.4332	1.8953	0.2256
VARSCENE ^{unc} (Star, Node, Edge)	0.9223	0.9324	0.9425	1.1342	1.6316	0.2357
VARSCENE ^{cond} (Star, Node, Edge)	<u>0.9222</u>	0.9267	<u>0.9455</u>	1.4657	1.9168	0.2292
VARSCENE ^{unc} (SP, Node, Edge)	0.9103	0.9248	0.9364	1.2269	1.6996	0.2398
VARSCENE ^{cond} (SP, Node, Edge)	0.9173	0.9264	0.9416	1.5882	2.0036	0.2339

Table 10. Performance of different graph generative models (analogous to Table 8) for different optimization objectives in the second stage of our learning for Visual Relationship Detection dataset.

Visual Genome (VG)						
Model	Star-Sim	Edge-Sim	Node-Sim	SP-Kernel	WL-Kernel	NSPD-Kernel
VARSCENE ^{unc}	0.5963	0.4576	0.4055	0.2293	0.1176	0.0147
VARSCENE ^{cond}	0.8660	0.5268	0.6226	0.0840	0.0779	0.0135
VARSCENE ^{unc} (Star)	0.7779	0.5289	0.4732	0.0863	0.0794	0.0134
VARSCENE ^{cond} (Star)	0.6589	0.2630	0.3659	0.1204	0.0869	0.0133
Small-sized Visual Genome (SVG)						
Model	Star-Sim	Edge-Sim	Node-Sim	SP-Kernel	WL-Kernel	NSPD-Kernel
VARSCENE ^{unc}	0.9235	0.7049	0.8335	1.0046	0.6556	0.0685
VARSCENE ^{cond}	0.9182	0.6964	0.8111	0.9621	0.6417	0.0633
VARSCENE ^{unc} (Star)	0.9339	0.7339	0.8505	0.7679	0.6104	0.0667
VARSCENE ^{cond} (Star)	0.9235	0.7307	0.8322	0.7654	0.6153	0.0622
Visual Relationship Detection (VRD)						
Model	Star-Sim	Edge-Sim	Node-Sim	SP-Kernel	WL-Kernel	NSPD-Kernel
VARSCENE ^{unc}	0.9194	0.9374	0.9403	1.0327	1.5619	0.2301
VARSCENE ^{cond}	0.9140	0.9372	0.9377	1.4588	1.9218	0.2275
VARSCENE ^{unc} (Star)	0.9067	0.9294	0.9275	1.0799	1.5911	0.2306
VARSCENE ^{cond} (Star)	0.9074	0.9210	0.9301	1.4106	1.8871	0.2242

Table 11. Performance provided by two variants of VARSCENE, viz., VARSCENE(RBF) which is the default variant having k as the RBF kernel as defined in Eq. (3) and VARSCENE(Star-Sim) which uses $k(G, G') = \cos(\nu(G), \nu(G'))$. Numbers in bold indicate the best performer.

Effect of Using RBF kernel. Recall that we use RBF kernel in Eq. (3) in our experiments. Here, we compare such an approach against using a simple cosine-similarity based kernel computation, that is, $k(G, G') = \cos(\nu(G), \nu(G'))$. Table 11 summarizes the results, which shows that using RBF kernel provides the performance boost for most similarity metrics.

Efficiency Analysis. Here, we analyze efficiency of different methods. Table 12 summarizes the results, which shows that training the base VAE model of VARSCENE is fastest on Visual Genome which is the largest dataset. For smaller datasets, i.e., Small-sized Visual Genome and Visual Relationship Detection, the training time is comparable with baselines. However, we note that training the MMD optimized decoder is slower in all datasets. The reason for more time in MMD optimization is that we need to sample a set of graphs (we sample 1000 graphs) at each iteration for computing $\widehat{\text{MMD}}$ in the objective as per Eqn. 9. For SVG and VRD datasets specifically, the graph generation process takes a higher amount of time. This is because its decoders sample stars beyond a limit (we set this limit as 50). In such cases, we sample graphs repeatedly until we get a valid graph or until we have repeated the procedure for a certain number of times (we set this limit as 20).

	VG	SVG	VRD
DeepGMG	2.0642	3.6121	2.5645
MolGAN	0.1274	0.1128	0.1414
GraphRNN	1.1582	0.0530	0.0435
GraphGen	0.0635	0.0429	0.0319
SceneGen	0.1945	0.0650	0.0510
VARSCENE (p_θ)	0.0334	0.0650	0.4188
VARSCENE (p_θ^{MMD})	5.3866	22.1826	24.9556

Table 12. Time (in seconds) for one iteration for a fixed batch size $|\mathcal{B}| = 64$.