

Probabilistic Debiasing of Scene Graphs

Bashirul Azam Biswas and Qiang Ji

Rensselaer Polytechnic Institute, Troy, NY-12180

{biswab, jiq}@rpi.edu

Abstract

The quality of scene graphs generated by the state-of-the-art (SOTA) models is compromised due to the long-tail nature of the relationships and their parent object pairs. Training of the scene graphs is dominated by the majority relationships of the majority pairs and, therefore, the object-conditional distributions of relationship in the minority pairs are not preserved after the training is converged. Consequently, the biased model performs well on more frequent relationships in the marginal distribution of relationships such as ‘on’ and ‘wearing’, and performs poorly on the less frequent relationships such as ‘eating’ or ‘hanging from’. In this work, we propose virtual evidence incorporated within-triplet Bayesian Network (BN) to preserve the object-conditional distribution of the relationship label and to eradicate the bias created by the marginal probability of the relationships. The insufficient number of relationships in the minority classes poses a significant problem in learning the within-triplet Bayesian network. We address this insufficiency by embedding-based augmentation of triplets where we borrow samples of the minority triplet classes from its neighborhood triplets in the semantic space. We perform experiments on two different datasets and achieve a significant improvement in the mean recall of the relationships. We also achieve better balance between recall and mean recall performance compared to the SOTA de-biasing techniques of scene graph models. Code is publicly available at <https://github.com/bashirulazam/within-triplet-debias>.

1. Introduction

Any visual relationship can be expressed as a triplet *subject-relationship-object* and all triplets in an image can be represented as a concise graph called Scene Graph (SG) [19] where the nodes represent the objects and the edges represent relationships. This representation has been proven useful for many downstream tasks such as image captioning [37], visual reasoning [24], and image generation [10]. Scene Graph Generation (SGG) has become one of the ma-

jor computer vision research arenas after the introduction of Visual Genome (VG) dataset [11]. The distribution of triplets in VG images has two distinct characteristics: (1) the presence of strong within-triplet prior, and (2) the long-tail distribution of the relationship. As shown in Figure 1 (a), the within-triplet prior dictates that ‘window’ will most likely be ‘on’ the ‘building’ rather than ‘eating’ it. Zeller *et al.* [43] has utilized this within-triplet prior as the conditional probability of relationships given subject and object by proposing a frequency baseline in the SGG task. On the other hand, the distribution of relationship labels suffers from a long-tailed nature and Tang *et al.* [27] addressed this long-tailed issue by considering a causal interpretation of the biased prediction. We argue that these two seemingly different characteristics of the relationship distribution are interrelated. The abundance of the *head* classes of the relationship distribution in Figure 1 (c), such as ‘on’ and ‘wearing’, arises from the abundance of their parent subject and object lying in the *head* region of Figure 1 (b).

As indicated by [5], the long-tailed distribution exists both in relationship and object label. Since relationship labels are dependent on their object pair because of the within-triplet prior, the long-tail distribution of the relationship worsens due to the long-tail nature of the object pairs. Crowd-collection of VG images creates *selection bias* and crowd-annotation of these images create *label-bias* [29] and *co-occurring-bias* [25]. To analyze such biases, we investigate the distribution of the object pair of the triplets in VG database. As shown in Figure 1 (b), ‘window-building’ and ‘man-shirt’ are the most frequently annotated pairs and top 1% object pair covers 33% of all triplets. As a result, the dominant relationships in these *head* pairs, such as ‘on’ and ‘wearing’, dominate the marginal distribution of Figure 1 (c).

In training a deep-learning-based SGG model, samplers will sample more relationships from the *head* pairs. As a result, the Maximum Likelihood Estimation (MLE) of the parameters is biased to predict the relationship classes in the *head* pairs [27] and the object-conditional representation of the relationship in the *tail* pairs will be lost in the training process. Therefore, various deep learning-based models,

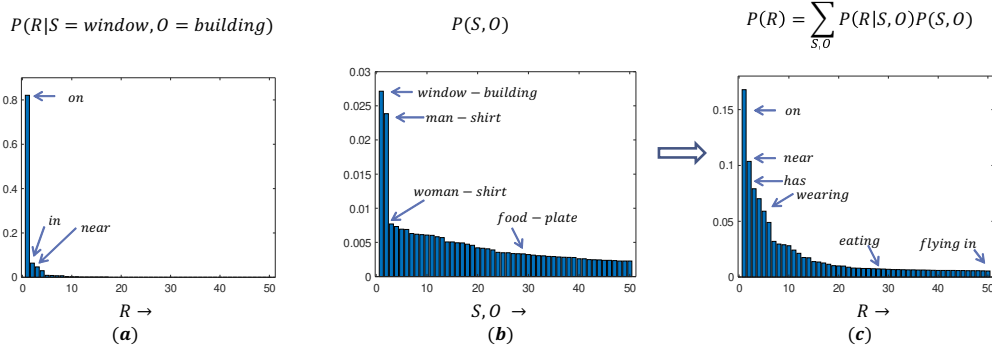


Figure 1. (a) Within-triplet dependency of relationship on its parent object pair; (b) long-tail nature of the pair statistics where 33% pair samples originated from top 1% pairs; (c) long-tail nature of the relationships showing the dominance of ‘on’ and ‘wearing’. The skewness in (c) is an effect of skewness in (b). Since ‘on’, ‘near’ or ‘wearing’ dominates in these top 1% pairs, they become the majority relationships in (c) and many other relationships, such as ‘eating’ or ‘flying in’, which dominate in the tail pairs of (b), are suppressed in the training process.

which attempt to implicitly capture such object-conditional representation [4, 45], fail to preserve the representation in the trained model and perform poorly on the tail region of the relationships.

Previous works attempt to retrieve the tail regions through re-sampling/re-weighting the minority classes in training [3, 5, 7] or through causal intervention in testing [27]. Their success is well-demonstrated by the significant increase of minority-driven evaluation metric *mean recall*. However, these approaches do not consider the strong within-triplet prior of triplets and hurt the performance of majority-driven evaluation metric *recall*. Keeping this gap in mind, we propose an inference-time post-processing methodology that bolsters the minority tail classes as well as hurts the majority head classes less brutally. We propose a within-triplet Bayesian Network (BN) that combines the within-triplet prior with uncertain biased evidence from SOTA models. Posterior inference with this BN simultaneously eradicates the long-tailed bias in the marginal distribution of the relationship and restores the object-conditional within-triplet prior.

Learning such a small within-triplet BN from the training data is a seemingly trivial task where we can perform simple MLE of parameters by counting. However, because of restricting our training samples only belonging to some top- N_r classes based on the marginal probability of relationship, we sacrifice many information revealing triplets in the minority pairs. For example, in the ‘man-pizza’ pair, we see there exist many interesting relationships such as ‘man-biting-pizza’ or ‘man consuming pizza’ which are semantically similar to one of the top- N_r valid triplets ‘man-eating-pizza’. This phenomenon is also a result of *label bias* [29]

where the annotator chooses some labels over another for the same category of objects or relationships. We propose a novel method of borrowing samples from such invalid triplets into learning the distribution of the valid triplets using embedding-based augmentation.

The posterior inference is the most efficient probabilistic tool to combine domain-dependent prior with instance-dependent evidence and, to the best of our knowledge, no prior work in SGG literature formulates the problem of triplet generation as a posterior inference problem. The overview of our approach is illustrated in Figure 2. In summary, our contribution is proposing a posterior inference-based post-processing method where we

- integrate the within-triplet priors with the evidence uncertainties generated by the measurement model and,
- introduce a simple yet novel learning scheme of the within-triplet network where we borrow samples from the semantically similar yet invalid triplet categories.

2. Related Works

Our proposed SGG model combines prior knowledge of triplets with the uncertain evidence of measurement models to address the long-tailed issue of relationships. Moreover, we learn the prior model using the similarity of triplets in the language embedding space. Hence, we divide the related works into four major categories as following

Implicit prior and context incorporation: Global context of an image has been captured either by BiLSTM, Graph Neural Network (GNN), attentional Graph Convolution Network (aGCN), or Conditional Random Field (CRF)

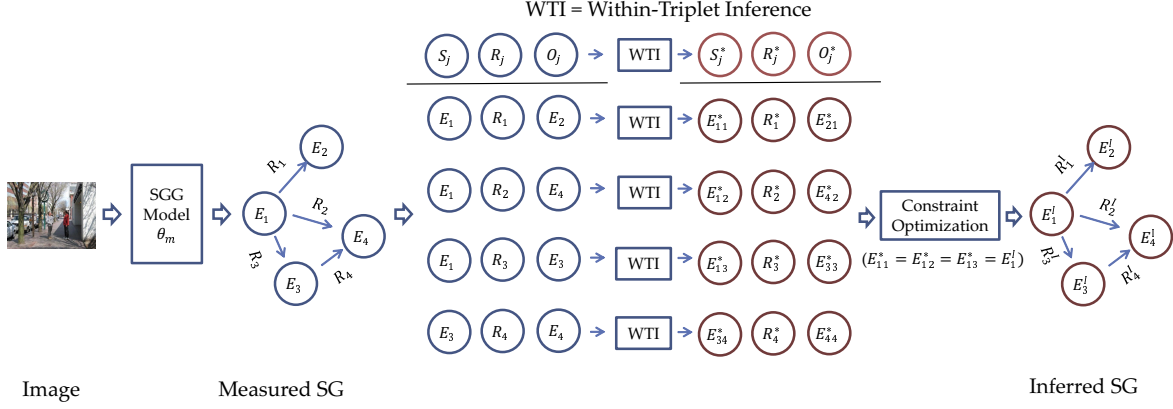


Figure 2. Overview of our proposed approach. For each testing image I , SGG baseline model θ_m generates physically connected triplets with associated uncertainties for the subject, object, and relationship. Our within-triplet inference framework takes all the uncertainties as uncertain evidence and performs posterior inference to infer each triplet separately. Afterward, a constrained optimization procedure is performed to resolve the conflict between object entities.

in [26, 32, 34, 36, 43]. Statistical information of triplets is encoded in [2, 4, 33]. Tang *et al.* [28] composed a dynamic tree structure to capture triplet context. Transformer-based context capturing is adopted in [13, 17, 23, 42]. Within-triplet relations are incorporated in learnable modules as embeddings in [22, 30, 44]. All of these above methods extract the relational information from the training data whereas several other approaches [6, 16, 19, 40, 41] rely on external knowledge-base such as ConceptNet [18].

Language models and ambiguity: Language and vision modules are combined together to guide the training process in [9, 19, 21, 40]. Embeddings from phrasal context [39], word embedding-based external knowledge incorporation [6, 41], and caption database [38] are used for supervision. The ambiguity of scene graph triplets is addressed in [12, 35, 46].

Long-tail distribution: An unbiased metric (*mean recall*) is proposed by [2, 28] to measure the performance of the SGG models in *tail* classes. Several works addressed this long-tailedness in training through modified loss functions [17, 41] and re-weighting/re-sampling training samples [3, 5, 7, 12, 14, 15]. The most relevant to our work is the inference-time causal intervention proposed by Tang *et al.* [27] where they identified the ‘bad bias’ from counterfactual causality and attained significantly higher mean recall than other SOTA models. We, instead of applying causal intervention, resort to a graphical model-based approach that can remove the ‘bad bias’ through uncertain evidence insertion into a Bayesian network while maintaining the ‘good bias’ by within-triplet prior incorporation.

Uncertain evidence: Incorporation of evidence uncertainty into belief networks has been discussed by Judea Pearl in [20]. A thorough discussion for interested readers can be found in [1].

3. Problem formulation

In scene graph generation database, every image I has an annotation of a scene graph $\mathcal{G}_I = (\mathcal{E}, \mathcal{R})$ where $\mathcal{E} = \{E_i, B_i\}_{i=1}^{N_E}$ contains the object classes E_i and their bounding boxes B_i whereas $\mathcal{R} = \{R_j(S_j, O_j)\}_{j=1}^{N_T}$ contains the relationships of a scene graph. Each relationship R_j exists between its subject S_j and object O_j where $S_j, O_j \in \{E_i\}$. Now, in training any SGG model parameterized by θ , we can write the cost function $\mathcal{J}(\theta)$ as

$$\mathcal{J}(\theta) = \mathbb{E}_{\sim p(I, \mathcal{G})} L(I, \mathcal{G}, \theta) \approx \frac{1}{M} \sum_{i=1}^M L(I_i, \mathcal{G}_i, \theta) \quad (1)$$

where $\mathbb{E}(\cdot)$ is the expectation operator, L is the loss function with parameter θ for a sample image I_i with associated scene graph \mathcal{G}_i and M is total number of images. Now, as shown in Figure 1, the distribution of scene graphs is skewed towards certain few categories of object pairs and their relationships. Therefore, while training, the cost function $\mathcal{J}(\theta)$ is driven by the *head* pairs, and the relationships dominating in the *tail* pairs are ignored. As a result, many object-conditional distributions of relationships are not preserved after the training is converged. We propose a test-time post-processing method where the object-conditional distribution of relationships is restored by a within-triplet prior Bayesian network.

4. Posterior inference of Scene Graphs

4.1. Within-triplet Bayesian network

A scene graph is a collection of connected triplets and each triplet has three semantic components - subject (S), object (O), and their relationship (R) and two spatial components - subject bounding box B_s and object bounding box

B_o . We consider the semantic components as random categorical variables and we aim to model their joint distribution $P(S, R, O)$. We assume there exists no latent confounder between the subject and object node whereas the relationship node depends on its parent subject and object. Formally, we assume the following statements hold true for any triplet -

1. Relationship label of a triplet depends on its subject and object $\implies S \rightarrow R \leftarrow O$;
2. Subject and object are independent, not given the relationship $\implies S \perp\!\!\!\perp O \nmid R$;
3. Subject and object becomes dependent, given the relationship $\implies S \not\perp\!\!\!\perp O \mid R$

Based on these assumptions, we can build a Bayesian network of triplet, as shown in Figure 3, which encodes the joint distribution using the chain rule as follows

$$P(S, R, O) = P(S)P(O)P(R|S, O) \quad (2)$$

where $P(S)$ and $P(O)$ represent the marginal distribution of the parents S and O , and $P(R|S, O)$ represents the conditional distribution of the relationship R given its parent subject and object. This Bayesian network encodes the prior joint distribution which resides within a triplet and hence we term it as within-triplet Bayesian network. In the next subsection, we discuss how we can debias the measurement probability of a relationship by incorporating them as uncertain evidence into this Bayesian network to perform posterior inference.

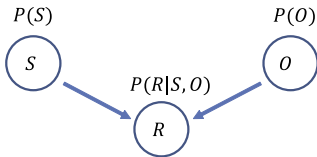


Figure 3. Within-triplet Bayesian network where prior probabilities of subject and object are combined with conditional probability of relationship.

4.2. Uncertain evidence

We denote any trained SGG model with parameter θ_m which generates measurement probabilities of subject, object, and relationship of every triplet for an image I as $P_{\theta_m, I}(S)$, $P_{\theta_m, I}(O)$, and $P_{\theta_m, I}(R)$. We consider these measurements as uncertain evidence of the nodes in the within-triplet BN in Figure 3. Since this evidence is uncertain, we incorporate it as virtual evidence nodes into the BN as shown in Figure 4. Following the virtual evidence method proposed by Judea Pearl in [20], we introduce binary virtual evidence nodes Z_s, Z_o , and Z_r as children of

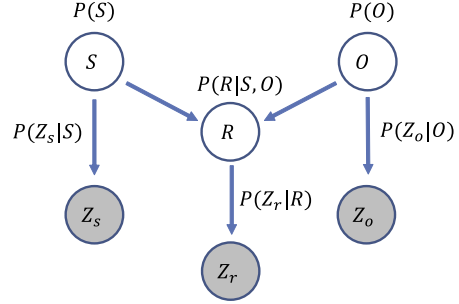


Figure 4. Uncertain evidence of each entity of triplet is incorporated into the Bayesian network as the conditional probability of virtual evidence node denoted as Z_s, Z_o , and Z_r .

their respective parent evidence nodes S, R , and O and instantiate them as *True*. According to Theorem 5 in [1], the conditional distributions of these virtual nodes maintain the following likelihood ratios

$$\begin{aligned} P(Z_s = 1|s_1) : \dots : P(Z_s = 1|s_n) &= \frac{P_{I, \theta_m}(s_1)}{P(s_1)} : \dots : \frac{P_{I, \theta_m}(s_n)}{P(s_n)} \\ P(Z_o = 1|o_1) : \dots : P(Z_o = 1|o_n) &= \frac{P_{I, \theta_m}(o_1)}{P(o_1)} : \dots : \frac{P_{I, \theta_m}(o_n)}{P(o_n)} \\ P(Z_r = 1|r_1) : \dots : P(Z_r = 1|r_n) &= \frac{P_{I, \theta_m}(r_1)}{P(r_1)} : \dots : \frac{P_{I, \theta_m}(r_n)}{P(r_n)} \end{aligned} \quad (3)$$

where $P(S = s), P(O = o)$, and $P(R = r)$ are the marginal probabilities of subject, object, and relationship node and $P_{\theta_m, I}(S = s), P_{\theta_m, I}(O = o)$, and $P_{\theta_m, I}(R = r)$ are their observed measurement probabilities from image I with model θ_m . Now, we have a complete Bayesian network in Figure 4 with well-defined marginal and conditional probabilities. A brief discussion on uncertain evidence and its incorporation into the Bayesian network is discussed in Appendix A of the supplementary material.

4.3. Within-Triplet Inference (WTI) of triplets

After the evidence incorporation as virtual evidence nodes, the posterior joint distribution of triplet nodes becomes

$$\begin{aligned} &P(S, R, O|Z_s = 1, Z_o = 1, Z_r = 1) \\ &\propto P(Z_s = 1|S)P(S)P(Z_o = 1|O)P(O)P(Z_r = 1|R)P(R|S, O) \\ &\propto P_{I, \theta_m}(S)P_{I, \theta_m}(O)\frac{P_{I, \theta_m}(R)}{P(R)}P(R|S, O) \end{aligned} \quad (4)$$

The Maximum a-Posterior (MAP) of this posterior joint distribution becomes

$$\begin{aligned} S^*, R^*, O^* &= \arg \max_{S, R, O} P(S, R, O|Z_s = 1, Z_o = 1, Z_r = 1) \\ &= \arg \max_{S, R, O} P_{I, \theta_m}(S)P_{I, \theta_m}(O)\frac{P_{I, \theta_m}(R)}{P(R)}P(R|S, O) \end{aligned} \quad (5)$$

In Eqn. (5), the within-triplet dependency of relationship is encoded in $P(R|S, O)$ and the measurement probability of relationship $P_{I,\theta_m}(R)$ is debiased by its marginal probability $P(R)$. The subject and object uncertainties are encoded in $P_{I,\theta_m}(S)$ and $P_{I,\theta_m}(O)$. We include some special cases of the MAP Eqn. (5) in Appendix B of supplementary material.

4.4. Constraint optimization in inferred triplets

Any object entity E_i can reside in multiple triplets as subject S or object O and after individual triplet inference, their inferred values can be different in different triplets. However, to form a valid scene graph, their values should be the same after the inference. Formally, if object entity E_i resides in J triplets, their inferred values E_{ij}^* must satisfy the following constraint

$$E_{i1}^* = E_{i2}^* = \dots = E_{ij}^* = \dots = E_{iJ}^* = E_i^I \quad (6)$$

One of the most straightforward ways to satisfy such constraint would be to take the mode of these inferred values as the final value for E_i . However, any object entity E_i^* should be consistent with respect to all of its connected triplets, and hence we formulate a two-step optimization algorithm to infer E_i^* and R_j^* from their connections.

Object Updating: In the first step, we infer each object label E_i^* combining its measurement probability $P_{I,\theta_m}(E_i)$ with the within-triplet probabilities of its connected triplets. We denote T_i^S and T_i^O as the sets of triplets where E_i acts as subject and object respectively

$$\begin{aligned} T_i^S &= \{t_p : t_p(S) = E_i\} \\ T_i^O &= \{t_q : t_q(O) = E_i\} \end{aligned} \quad (7)$$

The updated object probability for object E_i and inference of E_i is derived as

$$\begin{aligned} f(E_i) &= P_{I,\theta_m}(E_i) \left(\sum_{t_p \in T_i^S} P(R = r_{t_p}^I | S = E_i, O = o_{t_p}^I) \right. \\ &\quad \left. + \sum_{t_q \in T_i^O} P(R = r_{t_q}^I | S = s_{t_q}^I, O = E_i) \right) \\ E_i^* &= \arg \max f(E_i) \end{aligned} \quad (8)$$

Intuitively speaking, the updated object probability $f(E_i)$ derived in Eqn. (8) combines the uncertain evidence of an object entity $P_{I,\theta_m}(E_i)$ with the prior probabilities of the within-triplet Bayesian networks of all of its connected triplet.

Relationship Updating: After the first step is completed for each object entity, the conflicts of object entities are resolved with updated entity values. In the second step, we

update the relationship label of each triplet based on the updated subject and object values

$$\begin{aligned} R_j^* &= \arg \max_R P(R_j | Z_r = 1, S = s_j, O = o_j) \\ &= \arg \max_R \frac{P_{I,\theta_m}(R_j)}{P(R_j)} P(R_j | S = s_j, O = o_j) \end{aligned} \quad (9)$$

Detailed derivation and pseudo-code are provided in Appendix C of the supplementary material. We denote this as *constraint optimization* in our overview in Figure 2.

5. Learning BN with embedding similarity

The within-triplet priors $P(S) \in \mathbb{R}^{N_s}$, $P(O) \in \mathbb{R}^{N_o}$, and conditional distribution $P(R|S, O) \in \mathbb{R}^{N_s \times N_o \times N_r}$ are learned from annotations of training data where N_s , N_o , and N_r denote the number of categories of subject, object, and relationship. A training dataset of $\mathbb{Z}^{3 \times N}$ is created by collecting total N ground truth (GT) triplets from the training images. Afterward, we apply MLE to estimate $P(R|S, O)$ and $P(R)$ as follows

$$\begin{aligned} P(R = r | S = s, O = o) &= \frac{N_{s,r,o}^c}{\sum_{r'} N_{s,r',o}^c} \\ P(R) &= \sum_{S,O} P(R|S, O) P(S) P(O) \end{aligned} \quad (10)$$

where $N_{s,r,o}^c$ is the count of triplet with $S = s, O = o$ and $R = r$. However, because of selecting only top N_r relationships from the training data, many semantically similar triplets whose relationships lie outside these top N_r , are ignored in this count (e.g. ‘man-consuming-pizza’ is ignored whereas ‘man-eating-pizza’ is considered as a valid triplet). Hence, we propose a novel sample augmentation method using off-the-shelf sentence embedding models [31] where all the ignored invalid triplets, lying within a ϵ -neighbourhood of a valid triplet in the embedding space, are counted as augmented samples of that valid triplet. For any subject and object pair with $S = s, O = o$, we denote any valid triplet as $T = \{s, r, o\}$ and invalid triplet as $T_i = \{s, r_i, o\}$ where $r \in N_r$ and $r_i \notin N_r$. Now, if we denote the original count of T as $N_{s,r,o}^c$ and that of T_i as $N_{s,r_i,o}$, we can augment the original count as following

$$N_{s,r,o}^a = \begin{cases} N_{s,r,o}^c + \sum_{T_i \in \mathcal{N}_\epsilon(T)} N_{s,r_i,o} & \text{if } \mathcal{N}_\epsilon(T) \neq \emptyset \\ N_{s,r,o}^c & \text{if } \mathcal{N}_\epsilon(T) = \emptyset \end{cases} \quad (11)$$

Here, ϵ is a hyper-parameter and $\mathcal{N}_\epsilon(T)$ is defined as the ϵ -neighbourhood of a valid triplet T in the embedding space using the following criteria

$$\mathcal{N}_\epsilon(T) = \{T_i : \phi(f(T), f(T_i)) < \epsilon\} \quad (12)$$

where $\phi(f(T), f(T_i))$ represents the distance between two embedding vectors of two triplet $f(T)$ and $f(T_i)$. Since the

embeddings lie on a unit-sphere, we employ cosine similarity to measure the angular distance between two embedding vectors of two triplets. We visualize the ϵ -neighbourhood of a valid triplet T in Figure 5.

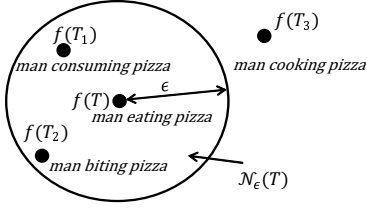


Figure 5. $\mathcal{N}_\epsilon(T)$ in the embedding space where $f(T_1), f(T_2) \in \mathcal{N}_\epsilon(T)$ and $f(T_3) \notin \mathcal{N}_\epsilon(T)$.

6. Experimental settings

6.1. Dataset

We evaluate our proposed method on two datasets: (1) Visual Genome (VG), and (2) GQA.

(1) Visual Genome: For SGG, the most commonly used database is the Visual Genome (VG) [11] and from the original database, the most frequent 150 object and 50 predicate categories are retained [33,43]. We adopt the standard train-test split ratio of 70 : 30. The number of prior triplets from the original training dataset is around 323K and after augmenting with $\epsilon = 0.05$, the number rises to around 391K.

(2) GQA: GQA [8] is a refined dataset derived from the VG images. We retain images only with the most frequent 150 object and 50 relationship categories. We train on around 50k valid images and perform the evaluation on the validation dataset of valid 7k images. We collect over 190K prior triplets and after embedding-based augmentation with $\epsilon = 0.05$, the number rises to around 200K.

6.2. Task description

A triplet is considered correct if the subject, relationship, and object label matches with the ground truth labels and the boundary boxes of subject and object have an Intersection over Union (IoU) of at least 50% with the ground truth annotations. We consider three test-time tasks, defined by [33], - (1) **PredCls**: known object labels and locations, (2) **SGCls**: known object locations, and (3) **SGDet**: where no information about objects are known. We apply graph constraints for all the tasks where for each pair of objects, only one relationship is allowed.

6.3. Evaluation metrics

The performance of our debiased SGGs is evaluated through recall (**R@K**) and mean recall (**mR@K**). **R@K** of an image is computed as the fraction of ground truth triplets in top@K predicted triplets [19] whereas **mR@K** computes

recall for each relationship separately and then the average over all relationships are computed [2,28].

6.4. Implementation details

We perform training and testing of the baseline models released by [27], [3], and [14]. We collect measurement results and ground truth annotations of the testing database using Python. For sample augmentation, we employ the sentence transformer model ‘all-mpnet-base-v2’ released by HuggingFace [31]. We learn the within-triplet prior from the original and augmented annotations, and perform posterior inference in MATLAB on a computer with core i5 7th generation Intel processor running at 2.5 MHz with 8.00 GB RAM. The total training time for the prior probabilities of VG dataset is 968s and that of GQA is 430s. The inference task per image requires 0.13s.

7. Experimental results

7.1. Quantitative results

We generate triplet measurements from four classical SGG models - (1) IMP [33], (2) MOTIF [43], (3) VC-Tree [28], and (4) Unb-MOTIF [27] from codebase [27] (\diamond), and two recent-most SOTA SGG models released by (1) DLFE-MOTIF [3] and (2) BGNN [14] (\blacklozenge). Considering these measurements as uncertain evidence, we perform within-triplet inference for all three settings in Sect. 6.2 and conduct two-step updating only for **SGCls** and **SGDet** settings. We denote the final results with the prefix ‘Inf-’. The Bayesian network is learned from the augmented counts of triplet derived by Eqn. (11). We report **R@K** and **mR@K** for all three tasks with VG and GQA in Table 1. We also include a separate comparison with other bias removal techniques (1) **Unb-** [27], (2) **DLFE-** [3], and (3) **NICE-** [12] in Table 2. Our method performs better in balancing the *head* and *tail* classes without any retraining of the biased model.

7.2. Analysis

7.2.1 Ablation study on prior and uncertain evidence

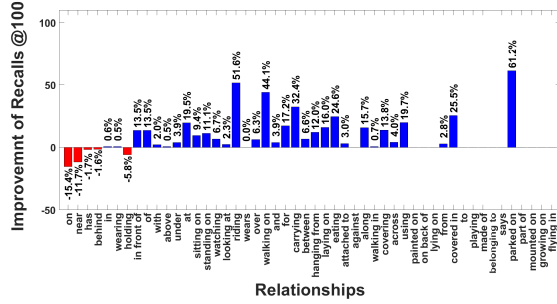
We perform an ablation study on the measurement results of the PredCls task by VCTree [28] in Table 3. We visualize the improvement of tail classes with BN learned from original and augmented samples in Figure 6a and 6b. In the former case, *mid* relationships are improving after inference and in the latter case both *mid* and *tail* are improving.

7.2.2 Ablation study on conflict resolution

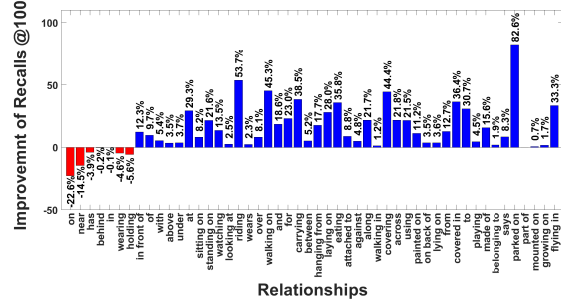
As discussed in Sect. 4.4, the potential conflicts of object labels after within-triplet inference can be resolved with naive mode-selection or by our proposed constraint optimization. We observe the effectiveness of the proposed optimization method over mode selection in Table 4.

DS	Method	Recall and Mean Recall @K					
		PredCls		SGCls		SGDet	
		R@50/100	mR@50/100	R@50/100	mR@50/100	R@50/100	mR@50/100
VG	IMP [◊] [33]	61.63/ 63.63	11.53/ 12.38	36.14/ 37.09	5.69/ 5.98	28.04/ 31.30	4.89/ 5.84
	Inf-IMP	59.93/ 62.02 ↓	25.14/ 28.34 ↑	36.02/ 37.09 ↓	12.57/ 14.06 ↑	26.50/ 29.51 ↓	8.58/ 10.67 ↑
	MOTIF [◊] [43]	59.57/ 63.95	12.88/ 15.47	36.45/ 38.47	7.66/ 8.83	26.85/ 30.50	5.61/ 6.73
	Inf-MOTIF	51.49/ 55.07 ↓	24.67/ 30.71 ↑	32.15/ 33.79 ↓	14.50/ 17.40 ↑	23.94/ 27.09 ↓	9.36/ 11.71 ↑
	VCTree [◊] [28]	65.46/ 67.18	15.36/ 16.61	44.15/ 45.11	9.17/ 9.83	29.94/ 32.57	6.21/ 6.96
	Inf-VCTree	59.50/ 60.97 ↓	28.14/ 30.72 ↑	40.69/ 41.55 ↓	17.31/ 19.40 ↑	27.74/ 30.10 ↓	10.40/ 11.86 ↑
	Unb-MOTIF [◊] [27]	45.87/ 51.24	24.75/ 28.69	26.30/ 28.78	13.21/ 15.06	16.25/ 19.53	8.65/ 10.47
	Inf-Unb-MOTIF	42.40/ 46.82 ↓	28.64/ 35.65 ↑	24.13/ 26.28 ↓	15.85/ 18.88 ↑	15.06/ 18.03 ↓	9.60/ 11.94 ↑
	DLFE-MOTIF [◆] [3]	51.63/ 53.28	26.87/ 28.75	28.79/ 29.66	15.61/ 16.38	24.22/ 27.95	10.62/ 12.61
	Inf-DLFE-MOTIF	43.27/ 44.82 ↓	35.25/ 38.20 ↑	24.34/ 25.14 ↓	19.74/ 20.66 ↑	20.61/ 23.80 ↓	14.07/ 16.76 ↑
GQA	BGNN [◆] [14]	58.15/ 60.41	29.46/ 31.83	-/-	-/-	30.26/ 34.98	10.37/ 12.31
	Inf-BGNN	55.42/ 57.47 ↓	32.18/ 34.27 ↑	-/-	-/-	26.16/ 30.11 ↓	13.24/ 16.10 ↑
	IMP [◊] [33]	61.94/ 63.68	13.04/ 13.74	34.25/ 34.83	7.46/ 7.80	25.39/ 27.42	5.77/ 6.58
	Inf-IMP	61.87/ 63.98 ↑	35.14/ 37.54 ↑	33.19/ 34.01 ↓	19.06/ 20.33 ↑	23.46/ 25.56 ↓	12.17/ 14.13 ↑
	MOTIF [◊] [43]	68.29/ 69.65	20.67/ 21.56	34.89/ 35.43	10.90/ 11.31	27.83/ 29.38	7.38/ 8.32
	Inf-MOTIF	62.95/ 64.23 ↓	37.93/ 40.07 ↑	31.82/ 32.35 ↓	19.09/ 20.00 ↑	25.51/ 26.86 ↓	14.34/ 15.84 ↑
	VCTree [◊] [28]	68.83/ 70.14	22.07/ 23.01	35.04/ 35.58	10.59/ 10.97	27.21/ 28.79	7.03/ 7.75
	Inf-VCTree	62.80/ 64.05 ↓	39.44/ 41.63 ↑	32.23/ 32.80 ↓	19.18/ 20.03 ↑	25.04/ 26.44 ↓	13.58/ 15.11 ↑
	Unb-MOTIF [◊] [27]	51.87/ 55.87	27.81/ 32.30	26.10/ 28.06	14.09/ 16.33	18.22/ 21.63	10.78/ 12.89
	Inf-Unb-MOTIF	49.86/ 53.60 ↓	34.45/ 40.80 ↑	24.76/ 26.74 ↓	17.18/ 20.59 ↑	16.84/ 19.94 ↓	12.35/ 14.76 ↑

Table 1. **R@K** and **mR@K** results of inference with prefix ‘Inf-’. We observe a significant increase in **mR@K** with a slight decrease in **R@K** for all baseline models in both datasets. Graph constraint is applied in all settings. Baseline results are generated by codebase released by [27] (◊) and by respective authors (◆)



(a) Recall improvement of each relationship with original samples.



(b) Recall improvement of each relationship with augmented samples.

Figure 6. Improvement of mean recalls with VCTree [28] evidence for PredCls task in VG. Relationships are ordered in descending order of their frequencies. In (a), BN learned with original samples improves the *mid* region whereas in (b), the embedding-based augmentation improves both the *mid* and the *tail* ones. In both cases, the *head* relationships are worsened after debiasing.

7.2.3 Effect of ϵ on training data augmentation

We observe the effect of neighborhood radius ϵ in the embedding space on **R@K** and **mR@K** in Table 5. Larger ϵ tends to hurt the majority classes more. We choose $\epsilon = 0.05$ for our final experiments.

7.2.4 Effect on zero-shot recall

The zero-shot prediction of the measurement model is compromised after posterior inference due to the MLE-based learning of BN. However, since unseen triplets will have higher entropy than the seen ones in the prediction phase,

Method	Re-train	R@K	mR@K
		@50/100	@50/100
DT2-ACBS [5]	Yes	23.3/ 25.6	35.9/ 39.7
BGNN [14]	Yes	59.2/ 61.3	30.4/ 32.9
IMP [33]	-	61.6/ 63.6	11.5/ 12.4
Inf-IMP (Ours)	No	59.9/ 62.0	25.1/ 28.3
MOTIF [33]	-	59.6/ 64.0	12.9/ 15.5
Unb-MOTIF [27]	No	45.9/ 51.2	24.8/ 28.7
DLFE-MOTIF [3]	Yes	51.6/ 53.2	26.9/ 28.8
NICE-MOTIF [12]	Yes	55.1/ 57.2	29.9/ 32.3
Inf-MOTIF (Ours)	No	51.5/ 55.1	24.7/ 30.7
VCTree [28]	-	65.5/ 67.2	15.4/ 16.6
Unb-VCTree [27]	No	47.2/ 51.6	25.4/ 28.7
DLFE-VCTree [3]	Yes	51.8/ 53.5	25.3/ 27.1
NICE-VCTree [12]	Yes	55.0/ 56.9	30.7/ 33.0
Inf-VCTree (Ours)	No	59.5/ 61.0	28.1/ 30.7

Table 2. Comparison with other de-biasing methods in PredCls. Without re-training, our loss in **R@K** is significantly lower, and gain in **mR@K** is higher or competitive than other SOTA debiasing methods.

Method	R@K	Mean R@K
	@50/100	@50/100
Uncertain evidence (VCTree)	65.5/ 67.2	15.4/ 16.6
WT BN only (org) (FREQ)	64.3/ 65.8 ↓	16.1/ 17.5 ↑
WT BN only (aug) (Ours)	62.7/ 64.2 ↓	16.3/ 17.8 ↑
WT BN (org) + Unc. Evi (Ours)	62.5 /64.1 ↓	22.7/ 24.8 ↑
WT BN (aug) + Unc. Evi (Ours)	59.5 /61.0 ↓	28.1/ 30.7 ↑

Table 3. Ablation study on PredCls performance for VG. We observe the consistent improvement of **mR@K** starting with uncertain evidence and ending in posterior inference with uncertain evidence with BN learned from augmented samples.

Method	R@K	mR@K
	@50/100	@50/100
VCTree [28]	44.2/ 45.1	9.2/ 9.8
Inf-VCTree (Conflict res. by mode)	40.3/ 41.2	16.9/ 18.7
Inf-VCTree (Conflict res. by opt.)	40.7/ 41.6 ↑	17.3/ 19.4 ↑

Table 4. Ablation study on conflict resolution for VCTree SGCLs performance for Visual Genome. The optimization algorithm performs better than the naive mode-selection version.

we can filter out the high-entropy triplets so that they do not get refined by the BN. An ablation study with respect to entropy threshold is shown in Table 6.

Method	ϵ	Recall@K	Mean Recal@K
		@50/@100	@50/@100
Inf- VCTree (org)	-	62.48/ 64.06	22.74/ 24.78
Inf- VCTree (aug)	0.03	59.44/ 60.88	28.01/ 30.55
	0.05	59.50/ 60.97	28.14/ 30.72
	0.07	59.27/ 60.73	28.28/ 30.92

Table 5. Effect of ϵ on PredCls performance for VG with VCTree baseline [28]. We find that using larger ϵ tends to drop **R@K** more. Based on this study, we choose $\epsilon = 0.05$ for our final experiments.

Method	Entropy Th.	ZS R@100	ZS mR@100
		org/aug	org/aug
VCTree [28]	-	6.02	16.61
Inf-VCTree	0	6.02/ 6.02	16.61/ 16.61
	1.5	5.81/ 6.13	20.20/ 22.71
	2.5	5.47/ 6.03	24.42/ 29.28
	3.912	1.06/ 2.30	24.78/ 30.70

Table 6. Zero-shot recall on PredCls for VG for VCTree [28]. Lower thresholds restore the zero-shot capability of measurement models by hurting the mean recall. We choose not to use any threshold (last row) to maximize performance in the mean recall.

8. Limitations

While improving **mR@K** by our proposed method, we lose performance in **R@K**. This phenomenon is prevailing in SGG debiasing works since we are perturbing the ‘head’ classes to gain improvement in the ‘tail’. Moreover, the augmentation hyper-parameter ϵ may vary from dataset to dataset and need to be chosen carefully. Another weakness is the compromise in the zero-shot prediction capability of a measurement model due to MLE-based learning of BN.

9. Conclusion

We proposed a debiasing strategy of scene graphs by combining prior and uncertain evidence of triplets in the Bayesian framework. We performed MAP inference and optimally solved the conflict between object entities to predict the debiased triplet from biased evidence. We augmented the count of valid triplets with semantically similar invalid triplets to alleviate sample insufficiency. Our method showcased significant improvement in mean recall with baseline measurements. We also attained a better balance between majority and minority performances of the relationships. In the future, we will extend the MAP inference for multiply connected triplets and we will explore well-defined criteria for zero-shot refinement to restore the zero-shot recall of the measurement models.

References

- [1] Hei Chan and Adnan Darwiche. On the revision of probabilistic beliefs using uncertain evidence. *Artificial Intelligence*, 163(1):67–90, 2005. [3](#), [4](#)
- [2] Tianshui Chen, Weihao Yu, Riquan Chen, and Liang Lin. Knowledge-embedded routing network for scene graph generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6163–6171, 2019. [3](#), [6](#)
- [3] Meng-Jiun Chiou, Henghui Ding, Hanshu Yan, Changhu Wang, Roger Zimmermann, and Jiashi Feng. Recovering the unbiased scene graphs from the biased ones. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 1581–1590, 2021. [2](#), [3](#), [6](#), [7](#), [8](#)
- [4] Bo Dai, Yuqi Zhang, and Dahua Lin. Detecting visual relationships with deep relational networks. In *Proceedings of the IEEE conference on computer vision and Pattern recognition*, pages 3076–3086, 2017. [2](#), [3](#)
- [5] Alakh Desai, Tz-Ying Wu, Subarna Tripathi, and Nuno Vasconcelos. Learning of visual relations: The devil is in the tails. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15404–15413, 2021. [1](#), [2](#), [3](#), [8](#)
- [6] Jiuxiang Gu, Handong Zhao, Zhe Lin, Sheng Li, Jianfei Cai, and Mingyang Ling. Scene graph generation with external knowledge and image reconstruction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1969–1978, 2019. [3](#)
- [7] Yuyu Guo, Lianli Gao, Xuanhan Wang, Yuxuan Hu, Xing Xu, Xu Lu, Heng Tao Shen, and Jingkuan Song. From general to specific: Informative scene graph generation via balance adjustment. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16383–16392, 2021. [2](#), [3](#)
- [8] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6700–6709, 2019. [6](#)
- [9] Zih-Siou Hung, Arun Mallya, and Svetlana Lazebnik. Contextual translation embedding for visual relationship detection and scene graph generation. *IEEE transactions on pattern analysis and machine intelligence*, 43(11):3820–3832, 2020. [3](#)
- [10] Justin Johnson, Agrim Gupta, and Li Fei-Fei. Image generation from scene graphs. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1219–1228, 2018. [1](#)
- [11] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International Journal of Computer Vision*, 123(1):32–73, 2017. [1](#), [6](#)
- [12] Lin Li, Long Chen, Yifeng Huang, Zhimeng Zhang, Songyang Zhang, and Jun Xiao. The devil is in the labels: Noisy label correction for robust scene graph generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18869–18878, June 2022. [3](#), [6](#), [8](#)
- [13] Rongjie Li, Songyang Zhang, and Xuming He. Sgtr: End-to-end scene graph generation with transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 19486–19496, June 2022. [3](#)
- [14] Rongjie Li, Songyang Zhang, Bo Wan, and Xuming He. Bipartite graph network with adaptive message passing for unbiased scene graph generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11109–11119, 2021. [3](#), [6](#), [7](#), [8](#)
- [15] Wei Li, Haiwei Zhang, Qijie Bai, Guoqing Zhao, Ning Jiang, and Xiaojie Yuan. Ppdl: Predicate probability distribution based loss for unbiased scene graph generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 19447–19456, June 2022. [3](#)
- [16] Xiaodan Liang, Lisa Lee, and Eric P Xing. Deep variation-structured reinforcement learning for visual relationship and attribute detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 848–857, 2017. [3](#)
- [17] Xin Lin, Changxing Ding, Jinquan Zeng, and Dacheng Tao. Gps-net: Graph property sensing network for scene graph generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3746–3753, 2020. [3](#)
- [18] Hugo Liu and Push Singh. Conceptnet—a practical commonsense reasoning tool-kit. *BT technology journal*, 22(4):211–226, 2004. [3](#)
- [19] Cewu Lu, Ranjay Krishna, Michael Bernstein, and Li Fei-Fei. Visual relationship detection with language priors. In *European conference on computer vision*, pages 852–869. Springer, 2016. [1](#), [3](#), [6](#)
- [20] Judea Pearl. *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Elsevier, 2014. [3](#), [4](#)
- [21] Julia Peyre, Ivan Laptev, Cordelia Schmid, and Josef Sivic. Detecting unseen visual relations using analogies. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1981–1990, 2019. [3](#)
- [22] Brigit Schroeder, Subarna Tripathi, and Hanlin Tang. Triplet-aware scene graph embeddings. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 0–0, 2019. [3](#)
- [23] Sahand Sharifzadeh, Sina Moayed Baharlou, and Volker Tresp. Classification by attention: Scene graph classification with prior knowledge. *arXiv preprint arXiv:2011.10084*, 2020. [3](#)
- [24] Jiaxin Shi, Hanwang Zhang, and Juanzi Li. Explainable and explicit visual reasoning over scene graphs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8376–8384, 2019. [1](#)
- [25] Krishna Kumar Singh, Dhruv Mahajan, Kristen Grauman, Yong Jae Lee, Matt Feiszli, and Deepti Ghadiyaram. Don’t judge an object by its context: learning to overcome contextual bias. In *Proceedings of the IEEE/CVF Conference*

- on *Computer Vision and Pattern Recognition*, pages 11070–11078, 2020. [1](#)
- [26] Mohammed Suhail, Abhay Mittal, Behjat Siddiquie, Chris Broaddus, Jayan Eledath, Gerard Medioni, and Leonid Sigal. Energy-based learning for scene graph generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13936–13945, 2021. [3](#)
- [27] Kaihua Tang, Yulei Niu, Jianqiang Huang, Jiabin Shi, and Hanwang Zhang. Unbiased scene graph generation from biased training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3716–3725, 2020. [1](#), [2](#), [3](#), [6](#), [7](#), [8](#)
- [28] Kaihua Tang, Hanwang Zhang, Baoyuan Wu, Wenhan Luo, and Wei Liu. Learning to compose dynamic tree structures for visual contexts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6619–6628, 2019. [3](#), [6](#), [7](#), [8](#)
- [29] Antonio Torralba and Alexei A Efros. Unbiased look at dataset bias. In *CVPR 2011*, pages 1521–1528. IEEE, 2011. [1](#), [2](#)
- [30] Hai Wan, Yonghao Luo, Bo Peng, and Wei-Shi Zheng. Representation learning for scene graph completion via jointly structural and visual embedding. In *IJCAI*, pages 949–956. Stockholm, Sweden, 2018. [3](#)
- [31] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online, Oct. 2020. Association for Computational Linguistics. [5](#), [6](#)
- [32] Sanghyun Woo, Dahun Kim, Donghyeon Cho, and In So Kweon. Linknet: Relational embedding for scene graph. *arXiv preprint arXiv:1811.06410*, 2018. [3](#)
- [33] Danfei Xu, Yuke Zhu, Christopher B Choy, and Li Fei-Fei. Scene graph generation by iterative message passing. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5410–5419, 2017. [3](#), [6](#), [7](#), [8](#)
- [34] Minghao Xu, Meng Qu, Bingbing Ni, and Jian Tang. Joint modeling of visual objects and relations for scene graph generation. *Advances in Neural Information Processing Systems*, 34, 2021. [3](#)
- [35] Gengcong Yang, Jingyi Zhang, Yong Zhang, Baoyuan Wu, and Yujiu Yang. Probabilistic modeling of semantic ambiguity for scene graph generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12527–12536, 2021. [3](#)
- [36] Jianwei Yang, Jiasen Lu, Stefan Lee, Dhruv Batra, and Devi Parikh. Graph r-cnn for scene graph generation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 670–685, 2018. [3](#)
- [37] Xu Yang, Kaihua Tang, Hanwang Zhang, and Jianfei Cai. Auto-encoding scene graphs for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10685–10694, 2019. [1](#)
- [38] Yuan Yao, Ao Zhang, Xu Han, Mengdi Li, Cornelius Weber, Zhiyuan Liu, Stefan Wermter, and Maosong Sun. Visual distant supervision for scene graph generation. *arXiv preprint arXiv:2103.15365*, 2021. [3](#)
- [39] Keren Ye and Adriana Kovashka. Linguistic structures as weak supervision for visual scene graph generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8289–8299, 2021. [3](#)
- [40] Ruichi Yu, Ang Li, Vlad I Morariu, and Larry S Davis. Visual relationship detection with internal and external linguistic knowledge distillation. In *Proceedings of the IEEE international conference on computer vision*, pages 1974–1982, 2017. [3](#)
- [41] Alireza Zareian, Svebor Karaman, and Shih-Fu Chang. Bridging knowledge graphs to generate scene graphs. In *European Conference on Computer Vision*, pages 606–623. Springer, 2020. [3](#)
- [42] Alireza Zareian, Zhecan Wang, Haoxuan You, and Shih-Fu Chang. Learning visual commonsense for robust scene graph generation. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIII 16*, pages 642–657. Springer, 2020. [3](#)
- [43] Rowan Zellers, Mark Yatskar, Sam Thomson, and Yejin Choi. Neural motifs: Scene graph parsing with global context. *arXiv:1711.06640*, 2017. [1](#), [3](#), [6](#), [7](#)
- [44] Hanwang Zhang, Zawlin Kyaw, Shih-Fu Chang, and Tat-Seng Chua. Visual translation embedding network for visual relation detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5532–5540, 2017. [3](#)
- [45] Ji Zhang, Kevin J Shih, Ahmed Elgammal, Andrew Tao, and Bryan Catanzaro. Graphical contrastive losses for scene graph parsing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11535–11543, 2019. [2](#)
- [46] Yi Zhou, Shuyang Sun, Chao Zhang, Yikang Li, and Wanli Ouyang. Exploring the hierarchy in relation labels for scene graph generation. *arXiv preprint arXiv:2009.05834*, 2020. [3](#)