# Scene Graph Generation with External Knowledge and Image Reconstruction

Jiuxiang Gu[1]*, Handong Zhao[2], Zhe Lin[2], Sheng Li[3], Jianfei Cai[1], Mingyang Ling[4]

[1] ROSE Lab, Interdisciplinary Graduate School, Nanyang Technological University, Singapore
[2] Adobe Research, USA [3] University of Georgia, USA [4] Google Cloud AI, USA

{jgu004, asjfcai}@ntu.edu.sg, {hazhao, zlin}@adobe.com
sheng.li@uga.edu, mingyangling@google.com

## Abstract

*Scene graph generation has received growing attention with the advancements in image understanding tasks such as object detection, attributes and relationship prediction, etc. However, existing datasets are biased in terms of object and relationship labels, or often come with noisy and missing annotations, which makes the development of a reliable scene graph prediction model very challenging. In this paper, we propose a novel scene graph generation algorithm with external knowledge and image reconstruction loss to overcome these dataset issues. In particular, we extract commonsense knowledge from the external knowledge base to refine object and phrase features for improving generalizability in scene graph generation. To address the bias of noisy object annotations, we introduce an auxiliary image reconstruction path to regularize the scene graph generation network. Extensive experiments show that our framework can generate better scene graphs, achieving the state-of-the-art performance on two benchmark datasets: Visual Relationship Detection and Visual Genome datasets.*

## 1. Introduction

With recent breakthroughs in deep learning and image recognition, higher-level visual understanding tasks, such as visual relationship detection, has been a popular research topic [9, 19, 15, 40, 44]. Scene graph, as an abstraction of objects and their complex relationships, provides rich semantic information of an image. It involves the detection of all ⟨*subject-predicate-object*⟩ triplets in an image and the localization of all objects. Scene graph provides a structured representation of an image that can support a wide range of high-level visual tasks, including image captioning [12, 14, 13, 43], visual question answering [36, 38, 47], image retrieval [11, 21], and image generation [20]. How-



Figure 1: Conceptual illustration of our scene graph learning model. The left (*green*) part illustrates the image to scene graph generation, the right (*blue*) part illustrates the image-level regularizer that reconstructs the image based on object labels and bounding boxes. The commonsense knowledge reasoning (*top*) is introduced to the scene graph generation process.

ever, it is not easy to extract scene graphs from images, since it involves not only detecting and localizing pairs of interacting objects but also recognizing their pairwise relationships. Currently, there are two categories of approaches for scene graph generation. Both categories group object proposals into pairs and use the phrase features (features of their union area) for predicate inference. The difference of the two categories lies in the different procedures. The first category detects the objects first and then recognizes the relationships between those objects [5, 28, 29]. The second category jointly identifies the objects and their relationships based on the object and relationship proposals [27, 25, 37].

Despite the promising progress introduced by these approaches, most of them suffer from the limitations of existing scene graph datasets. First, to comprehensively depict an image using the scene graph, it requires a wide variety of relation triplets ⟨*subject-predicate-object*⟩. Unfortunately, current datasets only capture a small portion of the knowledge [29], *e.g.*, Visual Relationship Detection (VRD) dataset. Training on such a dataset with long-tail

---

*This work was done during the author's internship at Adobe Research.

distributions will cause the prediction model bias towards those most-frequent relationships. Second, predicate labels are highly determined by the identification of object pairs [46]. However, due to the difficulty of exhaustively labeling bounding boxes of all instances of each object, the current large-scale crowd-sourced datasets like Visual Genome (VG) [22] are contaminated by noises (*e.g.*, missing annotations and meaningless proposals). Such a noisy dataset will inevitably result in a poor performance of the trained object detector [3], which further hinders the performance of predicate detection.

For human beings, we are capable of reasoning over visual elements of an image based on our commonsense knowledge. For example, in Figure 1, humans have the background knowledge: the subject (*woman*) appears / stands on something; the object (*snow*) enhances the evidence of the predicate (*skiing*). Commonsense knowledge can also help correct object detection. For example, the specific external knowledge for *skiing* benefits inference of the object (*snow*) as well. This motivates us to leverage commonsense knowledge to help scene graph generation.

Meanwhile, despite the crucial role of object labels for relationship prediction, existing datasets are very noisy due to the significant amount of missing object annotations. However, our goal is to obtain scene graphs with more complete scene representation. Motivated by this goal, we regularize our scene graph generation network by reconstructing the image from detected objects. Considering the case in Figure 1, a method might recognize *snow* as *grass* by mistake. If we generate an image based on the falsely predicted scene graph, this minor error would be heavily penalized, even though most of the *snow*'s relationships might be correctly identified.

The contributions of this paper are threefold. 1) We propose a knowledge-based feature refinement module to incorporate commonsense knowledge from an external knowledge base. Specifically, the module extracts useful information from ConceptNet [35] to refine object and phrase features before scene graph generation. We exploit Dynamic Memory Network (DMN) [23] to implement multi-hop reasoning over the retrieved facts and infer the most probable relations accordingly. 2) We introduce image-level supervision module by reconstructing the image to regularize our scene graph generation model. We view this auxiliary branch as a regularizer, which is only present during training. 3) We conduct extensive experiments on two benchmark datasets: VRD and VG datasets. Our empirical results demonstrate that our approach can significantly improve the state-of-the-art on scene graph generation.

## 2. Related Works

**Incorporating Knowledge in Neural Networks.** There has been growing interest in improving data-driven mod-els with external Knowledge Bases (KBs) in natural language processing [17, 4] and computer vision communities [24, 1, 6]. Large-scale structured KBs are constructed either by manual effort (*e.g.*, Wikipedia, DBpedia [2]), or by automatic extraction from unstructured or semi-structured data (*e.g.*, ConceptNet). One direction to improve the data-driven model is to distill external knowledge into Deep Neural Networks [39, 45, 18]. Wu *et al.* [38] encode the mined knowledge from DBpedia [2] into a vector and combine it with visual features to predict answers. Instead of aggregating the textual vectors with average-pooling operation [38], Li *et al.* [24] distill the retrieved context-relevant external knowledge triplet through a DMN for open-domain visual question answering. Unlike [38, 24], Yu *et al.* [45] extract linguistic knowledge from training annotations and Wikipedia, and distill knowledge to regularize training and provide extra cues for inference. A teacher-student framework is adopted to minimize the KL-divergence of the prediction distributions of teacher and student.

**Visual Relationship Detection.** Visual relationship detection has been investigated by many works in the last decade [21, 8, 7, 31]. Lu *et al.* [29] introduce generic visual relationship detection as a visual task, where they detect objects first, and then recognize predicates between object pairs. Recently, some works have explored the message passing for context propagation and feature refinement [41, 27]. Xu *et al.* [41] construct the scene graph by refining the object and relationship features jointly with message passing. Dai *et al.* [5] exploit the statistical dependencies between objects and their relationships and refine the posterior probabilities iteratively with a Conditional Random Field (CRF) network. More recently, Zeller *et al.* [46] achieve a strong baseline by predicting relationships with frequency priors. To deal with the large number of potential relations between objects, Yang *et al.* [42] propose a relation proposal network that prunes out uncorrelated object pairs, and captures the contextual information with an attentional graph convolutional network. In [25], they propose a clustering method which factorizes the full graph into subgraphs, where each subgraph is composed of several objects and a subset of their relationships.

Most related to our work are the approaches proposed by Li *et al.* [25] and Yu *et al.* [45]. Unlike [25], which focuses on the efficient scene graph generation, our approach addresses the long tail distribution of relationships by commonsense cues along with visual cues. Unlike [45], which leverages linguistic knowledge to regularize the network, our knowledge-based module improves the feature refining procedure by reasoning over a basket of commonsense knowledge retrieved from ConceptNet.
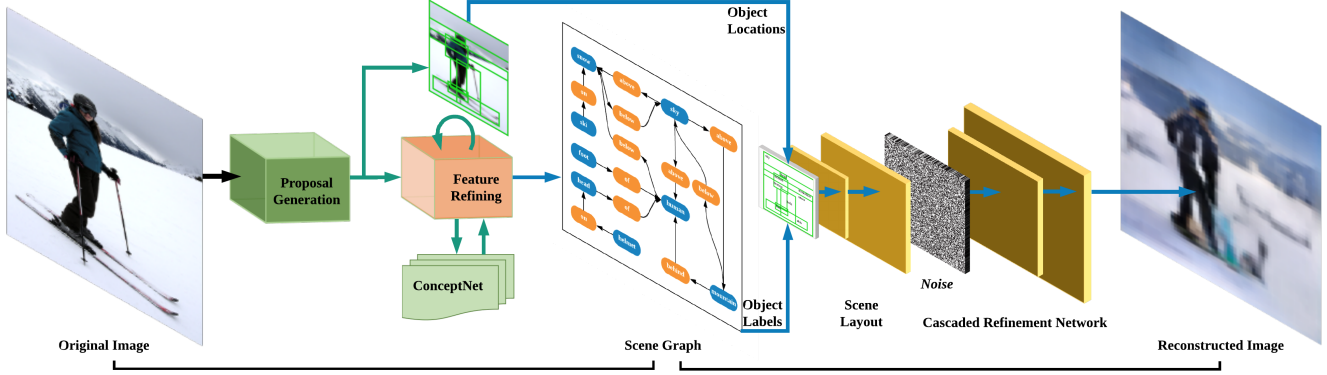
Figure 2: Overview of the proposed scene graph generation framework. The left part generates a scene graph from the input image. The right part is an auxiliary image-level regularizer which reconstructs the image based on the detected object labels and bounding boxes. After training, we discard the image reconstruction branch.

## 3. Methodology

Figure 2 gives an overview of our proposed scene graph generation framework. The entire framework can be divided into the following steps: (1) generate object and subgraph proposals for a given image; (2) refine object and subgraph features with external knowledge; (3) generate the scene graph by recognizing object categories with object features and recognizing object relations by fusing subgraph features and object feature pairs; (4) reconstruct the input image via an additional generative path. During training, we use two types of supervisions: scene graph level supervision and image-level supervision. For scene graph level supervision, we optimize our model by guiding the generated scene graph with the ground truth object and predicate categories. The image-level supervision is introduced to overcome the aforementioned missing annotations by reconstructing the image from objects and enforcing the reconstructed image close to the original image.

### 3.1. Proposal Generation

**Object Proposal Generation.** Given an image $\mathbf{I}$, we first use the Region Proposal Network (RPN) [33] to extract a set of object proposals:

$$[\mathbf{o}_0, \cdots, \mathbf{o}_{N-1}] = f_{\text{RPN}}(\mathbf{I}) \tag{1}$$

where $f_{\text{RPN}}(\cdot)$ stands for the RPN module, and $o_i$ is the $i$-th object proposal represented by a bounding box $r_i = [x_i, y_i, w_i, h_i]$ with $(x_i, y_i)$ being the coordinates of the top left corner and $w_i$ and $h_i$ being the width and the height of the bounding box, respectively. For any two different objects $\langle o_i, o_j \rangle$, there are two possible relationships in opposite directions. Thus, for $N$ object proposals, there are totally $N(N-1)$ potential relations. Although more object proposals lead to a bigger scene graph, the number of potential relations will increase dramatically, which significantly increases the computational cost and deteriorates the

inference speed. To address this issue, subgraph is introduced in [25] to reduce the number of potential relations by clustering.

**Subgraph Proposal Construction.** We adopt the clustering approach proposed in [25]. In particular, for a pair of object proposals, a subgraph proposal is constructed as the union box with the confidence score being the product of the scores of the two object proposals. Then, subgraph proposals are suppressed by non-maximum-suppression (NMS). In this way, a candidate relation can be represented by two objects and one subgraph: $\langle o_i, o_j, s_k^i \rangle$, where $i \neq j$ and $s_k^i$ is the $k$-th subgraph of all the subgraphs associated with $o_i$, which contains $o_j$ as well as some other object proposals. Following [25], we represent a subgraph and an object as a feature map, $\mathbf{s}_k^i \in \mathbb{R}^{D \times K_s \times K_s}$, and a feature vector, $\mathbf{o}_i \in \mathbb{R}^D$, respectively, where $D$ and $K_s$ are the dimensions.

### 3.2. Feature Refinement with External Knowledge

**Object and Subgraph Inter-refinement.** Considering that each object $\mathbf{o}_i$ is connected to a set of subgraphs $\mathbf{S}^i$ and each subgraph $\mathbf{s}_k$ is associated with a set of objects $\mathbf{O}^k$, we refine the object vector (resp. the subgraph) by attending the associated subgraph feature maps (resp. the associated object vectors):

$$\bar{\mathbf{o}}_i = \mathbf{o}_i + f_{s \to o}\left( \sum_{\mathbf{s}_k^i \in \mathbf{S}^i} \alpha_k^{s \to o} \cdot \mathbf{s}_k^i \right) \tag{2}$$

$$\bar{\mathbf{s}}_k = \mathbf{s}_k + f_{o \to s}\left( \sum_{\mathbf{o}_i^k \in \mathbf{O}^k} \alpha_i^{o \to s} \cdot \mathbf{o}_i^k \right) \tag{3}$$

where $\alpha_k^{s \to o}$ (resp. $\alpha_i^{o \to s}$) is the output of a softmax layer indicating the weight for passing $\mathbf{s}_k^i$ (resp. $\mathbf{o}_i^k$) to $\mathbf{o}_i$ (resp. $\mathbf{s}_k$), and $f_{s \to o}$ and $f_{o \to s}$ are non-linear mapping functions. This part is similar to [25]. Note that due to different dimensions of $\mathbf{o}_i$ and $\mathbf{s}_k$, pooling or spatial location based attention needs to be respectively applied for $s \to o$ or $o \to s$
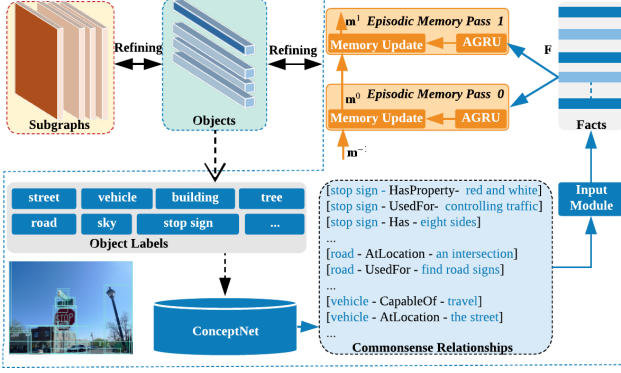
Figure 3: Illustration of our proposed knowledge-based feature refinement module. Given the object labels, we retrieve the facts (or symbolic triplets) from the ConceptNet (*bottom*), and then reason those facts with dynamic memory network using two passes (*top right*).

refinement. Interested readers are referred to [25] for details.

**Knowledge Retrieval and Embedding.** To address the relationship distribution bias of the current visual relationship datasets, we propose a novel feature refinement network to further improve the feature representation by taking advantage of the commonsense relationships in external knowledge base (KB). In particular, we predict the object label $a_i$ from the refined object vector $\bar{\mathbf{o}}_i$, and match $a_i$ with the corresponding semantic entities in KB. Afterwards, we retrieve the corresponding commonsense relationships from KB using the object label $a_i$:

$$a_i \xrightarrow{\text{retrieve}} \langle a_i, a_{i,j}^r, a_j^o, w_{i,j} \rangle, j \in [0, K-1] \quad (4)$$

where $a_{i,j}^r$, $a_j^o$ and $w_{i,j}$ are the top-$K$ corresponding relationships, the object entity and the weight, respectively. Note that the weight $w_{i,j}$ is provided by KB (*i.e.*, ConceptNet [35]), indicating how common a triplet $\langle a_i, a_{i,j}^r, a_j^o \rangle$ is. Based on the weight $w_{i,j}$, we can identify the top-$K$ most common relationships for $a_i$. Figure 3 illustrates the process of our proposed knowledge-based feature refinement module.

To encode the retrieved commonsense relationships, we first transform each symbolic triplet $\langle a_i, a_{i,j}^r, a_j^o \rangle$ into a sequence of words: $[X^0, \cdots, X^{T_a-1}]$, and then map each word in the sentence into a continuous vector space with word embedding $\mathbf{x}^t = \mathbf{W}_e X^t$. The embedded vectors are then fed into an RNN-based encoder [39] as

$$\mathbf{h}_k^t = \text{RNN}_{\text{fact}}(\mathbf{x}_k^t, \mathbf{h}_k^{t-1}), \ t \in [0, T_a - 1] \quad (5)$$

where $\mathbf{x}_k^t$ is the $t$-th word embedding of the $k$-th sentence, and $\mathbf{h}_k^t$ is the hidden state of the encoder. We use a bi-directional Gated Recurrent Unit (GRU) for $\text{RNN}_{\text{fact}}$ and the final hidden state $\mathbf{h}_k^{T_a-1}$ is treated as the vector repre-

sentation for the $k$-th retrieved sentence or fact, denoted as $\mathbf{f}_k^i$ for object $\mathbf{o}_i$.

**Attention-based Knowledge Fusion.** The knowledge units are stored in memory slots for reasoning and updating. Our target is to incorporate the external knowledge into the procedure of feature refining. However, for $N$ objects, we have $N \times K$ relevant fact vectors in memory slots. This makes it difficult to distill the useful information from the candidate knowledge when $N \times K$ is large. DMN [23] provides a mechanism to pick out the most relevant facts by using an episodic memory module. Inspired by this, we adopt the improved DMN [39] to reason over the retrieved facts $\mathbf{F}$, where $\mathbf{F}$ denotes the set of fact embedding $\{\mathbf{f}_k\}$. It consists of an attention component which generates a contextual vector using the episode memory $\mathbf{m}^{t-1}$. Specifically, we feed the object vector $\bar{\mathbf{o}}$ to a non-linear fully-connected layer and attend the facts as follows:

$$\mathbf{q} = \tanh(\mathbf{W}_q \bar{\mathbf{o}} + \mathbf{b}_q) \quad (6)$$

$$\mathbf{z}^t = [\mathbf{F} \circ \mathbf{q}; \mathbf{F} \circ \mathbf{m}^{t-1}; |\mathbf{F} - \mathbf{q}|; |\mathbf{F} - \mathbf{m}^{t-1}|] \quad (7)$$

$$\mathbf{g}^t = \text{softmax}(\mathbf{W}_1 \tanh(\mathbf{W}_2 \mathbf{z}^t + \mathbf{b}_2) + \mathbf{b}_1) \quad (8)$$

$$\mathbf{e}^t = \text{AGRU}(\mathbf{F}, \mathbf{g}^t) \quad (9)$$

where $\mathbf{z}^t$ is the interactions between the facts $\mathbf{F}$, the episode memory $\mathbf{m}^{t-1}$ and the mapped object vector $\mathbf{q}$, $\mathbf{g}^t$ is the output of a softmax layer, $\circ$ is the element-wise product, $|\cdot|$ is the element-wise absolute value, and $[\ ;\ ]$ is the concatenation operation. Note that $\mathbf{q}$ and $\mathbf{m}$ need to be expanded via duplication in order to have the same dimension as $\mathbf{F}$ for the interactions. In (9), AGRU($\cdot$) refers to the Attention based GRU [39] which replaces the update gate in GRU with the output attention weight $\mathbf{g}_k^t$ for fact $k$:

$$\mathbf{e}_k^t = g_k^t \text{GRU}(\mathbf{f}_k, \mathbf{e}_{k-1}^t) + (1 - g_k^t)\mathbf{e}_{k-1}^t \quad (10)$$

where $\mathbf{e}_K^t$ is the final state of the episode which is the state of the GRU after all the $K$ sentences have been seen.

After one pass of the attention mechanism, the memory is updated using the current episode state and the previous memory state:

$$\mathbf{m}^t = \text{ReLU}(\mathbf{W}_m[\mathbf{m}^{t-1}; \mathbf{e}_K^t; \mathbf{q}] + \mathbf{b}_m). \quad (11)$$

where $\mathbf{m}^t$ is the new episode memory state. By the final pass $T_m$, the episodic memory $\mathbf{m}^{T_m-1}$ can memorizes useful knowledge information for relationship prediction.

The final episodic memory $\mathbf{m}^{T_m-1}$ is passed to refine the object feature $\bar{\mathbf{o}}$ as

$$\tilde{\mathbf{o}} = \text{ReLU}(\mathbf{W}_c[\bar{\mathbf{o}}; \mathbf{m}^{T_m-1}] + \mathbf{b}_c) \quad (12)$$

where $\mathbf{W}_c$ and $\mathbf{b}_c$ are parameters to be learned. In particular, we refine objects with KB via (12) as well as jointly refining objects and subgraphs by replacing $\{\mathbf{o}_i, \mathbf{s}_i\}$ with $\{\tilde{\mathbf{o}}_i, \bar{\mathbf{s}}_i\}$ in (2) and (3), in an iterative fashion (see Alg. 1).
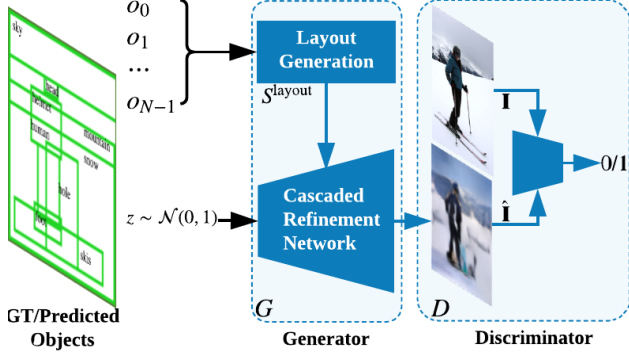
Figure 4: Illustration of our proposed object-to-image generation module Gen$_{\text{o2i}}$.

## 3.3. Scene Graph Generation

**Relation Prediction.** After the feature refinement, we can predict object labels as well as predicate labels with the refined object and subgraph features. For object label, we can predict it directly with the object features. For relationship label, as the subgraph feature is related to several object pairs, we predict the label based on subject and object feature vectors along with their corresponding subgraph feature map. We formulate the inference process as

$$\text{P}_{i,j} \sim \text{softmax}(f_{\text{rel}}([\tilde{\mathbf{o}}_i \otimes \bar{\mathbf{s}}_k; \tilde{\mathbf{o}}_j \otimes \bar{\mathbf{s}}_k; \bar{\mathbf{s}}_k])) \quad (13)$$

$$\text{V}_i \sim \text{softmax}(f_{\text{node}}(\tilde{\mathbf{o}}_i)) \quad (14)$$

where $f_{\text{rel}}(\cdot)$ and $f_{\text{node}}(\cdot)$ denote the mapping layers for predicate and object recognition, respectively, and $\otimes$ denotes the convolution operation [25]. Then, we can construct the scene graph as: $\mathcal{G} = \langle V_i, P_{i,j}, V_j \rangle, i \neq j$.

**Scene Graph Level Supervision.** Like other approaches [26, 25, 37], during training we want the generated scene graph close to the ground-truth scene graph by optimizing the scene graph generation process with object detection loss and relationship classification loss

$$\mathcal{L}_{\text{im2sg}} = \lambda_{\text{pred}}\mathcal{L}_{\text{pred}} + \lambda_{\text{obj}}\mathcal{L}_{\text{obj}} + \lambda_{\text{reg}}\mathbf{1}_{u \geq 1}\mathcal{L}_{\text{reg}} \quad (15)$$

where $\mathcal{L}_{\text{pred}}$, $\mathcal{L}_{\text{obj}}$ and $\mathcal{L}_{\text{reg}}$ are the predicate classification loss, the object classification loss and the bounding box regression loss, respectively, $\lambda_{\text{obj}}$, $\lambda_{\text{pred}}$ and $\lambda_{\text{reg}}$ are hyperparameters, and $\mathbf{1}$ is the indicator function with $u$ being the object label, $u \geq 1$ for object categories and $u = 0$ for background.

For the predicate detection, the output is the probability over all the candidate predicates. $\mathcal{L}_{\text{pred}}$ is defined as the softmax loss. Like the predicate classification, the output of the object detection is the probability over all the object categories. $\mathcal{L}_{\text{cls}}$ is also defined as the softmax loss. For the bounding box regression loss $\mathcal{L}_{\text{reg}}$, we use smooth $L_1$ loss [33].

## 3.4. Image Generation

To better regularize the networks, an object-to-image generative path is added. Figure 4 depicts our proposed object-to-image generation module Gen$_{\text{o2i}}$. In particular, we first compute a scene layout based on the object labels and their corresponding locations. For each object $i$, we expand the object embedding vectors $\mathbf{o}_i \in \mathbb{R}^D$ to shape $D \times 8 \times 8$, and then wrap it to the position of the bounding box $r_i$ using bilinear interpolation to give an object layout $o_i^{\text{layout}} \in \mathbb{R}^{D \times H \times W}$, where $D$ is the dimension of the embedding vectors for objects and $H \times W = 64 \times 64$ is the output image resolution. We sum all object layouts to obtain the scene layout $S^{\text{layout}} = \sum_i o_i^{\text{layout}}$.

Given the scene layout, we synthesize an image that respects the object positions with an image generator $G$. Here, we adopt a cascaded refinement network [20] which consists of a series of convolutional refinement modules to generate the image. The spatial resolution doubles between the convolutional refinement modules. This allows the generation to proceed in a coarse-to-fine manner. For each module, it takes two inputs. One is the output from the previous module (the first module takes Gaussian noise), and the other one is the scene layout $S^{\text{layout}}$, which is downsampled to the input resolution of the module. These inputs are concatenated channel-wisely and passed to a pair of $3 \times 3$ convolution layers. The outputs are then upsampled using nearest-neighbor interpolation before being passed to the next module. The output from the last module is passed to two final convolution layers to produce the output image.

**Image-level Supervision.** In addition to the common pixel reconstruction loss $\mathcal{L}_{\text{pixel}}$, we also adopt a conditional GAN loss [32], considering the image is generated based on the objects. In particular, we train the discriminator $D_i$ and the generator $G_i$ by alternatively maximizing $\mathcal{L}_{D_i}$ in Eq. (16) and $\mathcal{L}_{G_i}$ in Eq. (17):

$$\mathcal{L}_{D_i} = \mathbb{E}_{I \sim p_{\text{real}}}[\log D_i(\mathbf{I})] \quad (16)$$

$$\mathcal{L}_{G_i} = \mathbb{E}_{\hat{I} \sim p_{\text{G}}}[\log(1 - D_i(\hat{\mathbf{I}})] + \lambda_p\mathcal{L}_{\text{pixel}} \quad (17)$$

where $\lambda_p$ is the tuning parameter. For the generator loss, we maximize $\log D_i(G_i(z|S^{\text{layout}}))$ rather than minimizing the original $\log(1 - D_i(G_i(z|S^{\text{layout}})))$ for better gradient behavior. For the pixel reconstruction loss, we calculate the $\ell_1$ distance between the real image $\mathbf{I}$ and a corresponding synthetic image $\hat{\mathbf{I}}$ as $||\mathbf{I} - \hat{\mathbf{I}}||_1$.

As shown in Figure 2, we view the object-to-image generation branch as a regularizer. It can be seen as a corrective model for scene graph generation by improving the performance of object detection. During training, backpropagation from losses (15), (16), and (17) influences the model parameter updates. This image-level supervision can be seen as a corrective model for scene graph generation by improving the performance of object detection. The gradi-

ents back-propagated from the object-to-image branch update the parameters of our object detector and the feature refinement module which is followed by the relation prediction.

Alg. 1 summarizes the entire training procedure.

---

**Algorithm 1** Training procedure.

---

**Input:** Image $\mathbf{I}$, number of training steps $T_s$.
 1: Pretrain image generation module $\text{Gen}_{\text{o2i}}$ (GT objects)
 2: **for** $t = 0 : T_m - 1$ **do**
 3:     Get objects and relationship triples.
 4:     Proposal Generation: $(\mathbf{O}, \mathbf{S}) \leftarrow \mathbf{I}$ {RPN}
 5:     /*Knowledge-based Feature Refining*/
 6:     **for** $r = 0 : T_r - 1$ **do**
 7:         $\bar{\mathbf{o}}_i \leftarrow \{\mathbf{o}_i, \mathbf{S}^i\}$ /*Refining using (2)*/
 8:         $\bar{\mathbf{s}}_k \leftarrow \{\mathbf{s}_k, \mathbf{O}^k\}$ /*Refining using (3)*/
 9:         $\tilde{\mathbf{o}}_i \leftarrow \{\mathbf{F}, \bar{\mathbf{o}}_\mathbf{i}\}$ /*Refining using (12)*/
10:         $\mathbf{o}_i \leftarrow \tilde{\mathbf{o}}_i, \mathbf{s}_i \leftarrow \bar{\mathbf{s}}_i$
11:     **end for**
12:     Update parameters with $\text{Gen}_{\text{o2i}}$ (predicted objects)
13:     Update parameters with (15)
14: **end for**
**Function:** $\text{Gen}_{\text{o2i}}$
**Input:** Real image $\mathbf{I}$, objects (GT / predicted).
 1: Object Layout Generation: $o_i^{\text{layout}} \leftarrow \{\mathbf{o}_i, r_i\}$
 2: Scene Layout Generation: $S^{\text{layout}} = \sum_i o_i^{\text{layout}}$
 3: Image Reconstruction: $\hat{\mathbf{I}} = G_i(z, S^{\text{layout}})$
 4: Update image generator $G_i$ parameters using (17).
 5: Update image discriminator $D_i$ parameters using (16).

---

# 4. Experiments

## 4.1. Datasets

We evaluate our approach on two datasets: VRD [29] and VG [26]. VRD is the most widely used benchmark dataset for visual relationship detection. Compared with VRD, the raw VG [22] contains a large number of noisy labels. In our experiment, we use a cleansed-version VG-MSDN in [26]. Detailed statistics of both datasets are shown in Table 1.

For the external KB, we employ the English subgraph of ConceptNet [35] as our knowledge graph. ConceptNet is a large-scale graph of general knowledge which aims to align its knowledge resources on its core set of 40 relations. A large portion of these relation types can be considered as visual relations, such as, spatial co-occurrence (*e.g.*, *AtLocation*, *LocatedNear*), visual properties of objects (*e.g.*, *HasProperty*, *PartOf* ), and actions (*e.g.*, *CapableOf*, *UsedFor*).

## 4.2. Implementation Details

As shown in Alg. 1, we train our model in two phrases. The initial phase looks only at the object annotations of

Table 1: Dataset statistics. #Img and #Rel denote the number of images and relation pairs respectively, #Obj denotes the number of object categories, and #Pred denotes the number of predicate categories.

| Dataset | Training Set | | Testing Set | | #Obj | #Pred |
|---|---|---|---|---|---|---|
| | #Img | #Rel | #Img | #Rel | | |
| VRD [29] | 4,000 | 30,355 | 1,000 | 7,638 | 100 | 70 |
| VG-MSDN [26] | 46,164 | 507,296 | 10,000 | 111,396 | 150 | 50 |

the training set, ignoring the relationship triplets. For each dataset, we filter the objects according to the category and relation vocabularies in Table 1. We then learn an image-level regularizer that reconstructs the image based on the object labels and bounding boxes. The output size of the image generator is $64 \times 64 \times 3$, and the real image is resized before inputting to the discriminator. We train the regularizer with learning rate $10^{-4}$ and batch size 32. For each mini-batch we first update $G_i$, and then update $D_i$.

The second phase jointly trains the scene graph generation model and the auxiliary reconstruction branch. We adopt the Faster R-CNN [33] associated with VGG-16 [34] as the backbone. During training, the number of object proposals is 256. For each proposal, we use ROI align [16] pooling to generate object and subgraph features. The subgraph regions are pooled to $5 \times 5$ feature maps. The dimension $D$ of the pooled object vector and the subgraph feature map is set to 512. For the knowledge-based refinement module, we set the dimension of word embedding to 300 and initialize it with the GloVe 6B pre-trained word vectors [30]. We keep the top-8 commonsense relationships. The number of hidden units of the fact encoder is set to 300, and the dimension of episodic memory is set to 512. The iteration number $T_m$ of DMN update is set to 2. For the relation inference module, we adopt the same bottleneck layer as [25]. All the newly introduced layers are randomly initialized except the auxiliary regularizer. We set $\lambda_{\text{pred}} = 2.0$, $\lambda_{\text{cls}} = 1.0$, and $\lambda_{\text{reg}} = 0.5$ in Eq (15). The hyperparameter $\lambda_p$ in Eq (17) is set to 1.0. The iteration number $T_r$ of the feature refinement is set to 2. We first train RPNs and then jointly train the entire network. The initial learning rate is 0.01, decay rate is 0.1, and stochastic gradient descent (SGD) is used as the optimizer. We deploy weight decay and dropout to prevent over-fitting.

During testing, the image reconstruction branch will be discarded. We respectively set the RPN non-maximum suppression (NMS) [33] threshold to 0.6 and subgraph clustering [25] threshold to 0.5. We output all the predicates and use the top-1 category as the prediction for objects and relations. Models are evaluated on two tasks: Visual Phrase Detection (**PhrDet**) and Scene Graph Generation (**SGGen**). **PhrDet** is to detect the $\langle$subject-predicate-object$\rangle$ phrases. **SGGen** is to detect the objects within the image and recognize their pairwise relationships. Following [29, 25],

the Top-$K$ Recall (denoted as Rec@$K$) is used as the performance metric; it calculates how many labeled relationships are hit in the top K predictions. In our experiments, Rec@50 and Rec@100 are reported. Note that, Li *et al*. [26] and Yang *et al*. [42] reported the results on two more metrics: *Predicate Recognition* and *Phrase Recognition*. These two evaluation metrics are based on ground-truth object locations, which is not the case we consider. In our setting, we use detected objects for image reconstruction and scene graph generation. To be consistent with the training, we choose *PhrDet* and *SGGen* as the evaluation metrics, which is also more practical.

### 4.3. Baseline Approaches for Comparisons

**Baseline.** This baseline model is the re-implementation of Factorizable Net [25]. We re-train it based on our backbone. Specifically, we use the same RPN model, and jointly train the scene graph generator until convergence.
**KB.** This model is a KB-enhanced version of the baseline model. External knowledge triples are incorporated in DMN. The explicit knowledge-based reasoning is incorporated in the feature refining procedure.
**GAN.** This model improves the baseline model by attaching an auxiliary branch that generates the image from objects with GAN. We train this model in two phases. The first phase trains the image reconstruction branch only with the object annotations. Then we refine the model jointly with the scene graph generation model.
**KB-GAN.** This is our full model containing both KB and GAN. It is initialized with the trained parameters from KB and GAN, and fine-tuned with Alg. 1.

### 4.4. Quantitative Results

In this section, we present our quantitative results and analysis. To verify the effectiveness of our approach and analyze the contribution of each component, we first compare different baselines in Table 2, and investigate the improvement in recognizing objects in Table 3. Then, we conduct a simulation experiment on VRD to investigate the effectiveness of our auxiliary regularizer in Table 4. The comparison of our approach with the state-of-the-art methods is reported in Table 5.
**Component Analysis.** In our framework, we proposed two novel modules – KB-based feature refinement (KB) and auxiliary image generation (GAN). To get a clear sense of how these components affect the final performance, we perform ablation studies in Table 2. The left-most columns in Table 2 indicate whether or not we use KB and GAN in our approach. To further investigate the improvement of our approach on recognizing objects, we also report object detection performance mAP [10] in Table 3.

In Table 2, we observe that KB boosts **PhrDet** and **SGGen** significantly. This indicates our knowledge-based

Table 2: Ablation studies of individual components of our method on VRD.

| KB | GAN | PhrDet | | SGGen | |
|---|---|---|---|---|---|
| | | Rec@50 | Rec@100 | Rec@50 | Rec@100 |
| - | - | 25.57 | 31.09 | 18.16 | 22.30 |
| ✓ | - | 27.02 | 34.04 | 19.85 | 24.58 |
| - | ✓ | 26.65 | 34.06 | 19.56 | 24.64 |
| ✓ | ✓ | **27.39** | **34.38** | **20.31** | **25.01** |

Table 3: Ablation study of the object detection on VRD.

| Model | Faster R-CNN [33] | ViP-CNN [27] | Baseline | KB | GAN | KB-GAN |
|---|---|---|---|---|---|---|
| mAP | 14.35 | 20.56 | 20.70 | 22.26 | 22.10 | **22.49** |

feature refinement can effectively learn the commonsense knowledge of objects to achieve high recall for the correct relationships. By adding the image-level supervision to the baseline model, the performance is further improved. This improvement demonstrates that the proposed image-level supervision is capable of capturing meaningful context across the objects. These results align with our intuitions discussed in the introduction. With KB and GAN, our model can generate scene graphs with high recall.

Table 3 demonstrates the improvement in recognizing objects. We can see that our full model (KB-GAN) outperforms Faster R-CNN [33], ViP-CNN [27] measured by mAP. It is worth noticing that the huge gain of KB illustrates that the introduction of commonsense knowledge substantially contributes to the object detection task.

Table 4: Ablation study of image-level supervision on sub-sampled VRD.

| KB | GAN | PhrDet | | SGGen | |
|---|---|---|---|---|---|
| | | Rec@50 | Rec@100 | Rec@50 | Rec@100 |
| - | - | 15.44 | 20.96 | 10.94 | 14.53 |
| - | ✓ | 24.07 | 30.89 | 17.50 | 22.31 |
| ✓ | ✓ | **26.62** | **31.13** | **19.78** | **24.17** |

**Investigation on Image-level Supervision.** As aforementioned, our image-level supervision can exploit the instances of rare categories. To demonstrate that our introduced image-level supervision can help on this issue, we exaggerate the problem by randomly removing 20% object instances as well as their corresponding relationships from the dataset. In Table 4, we can see that training on such a sub-sampled dataset (with only 80% object instances), Rec@50 of the baseline model drops from 25.57 (resp. 18.16) to 15.44 (resp. 10.94) for PhrDet and SGGen. However, with the help of GAN, Rec@50 of our final model decreases only slightly from 27.39 (resp. 20.31) to 26.62 (resp. 19.78) for PhrDet and SGGen, respectively.

We give our explanation on this significant performance improvement as below. Too many low-frequency categories deteriorate the training gain when only utilizing the class la-

Table 5: Comparison with existing methods on **PhrDet** and **SGGen**.

| Dataset | Model | PhrDet | | SGGen | |
|---|---|---|---|---|---|
| | | Rec@50 | Rec@100 | Rec@50 | Rec@100 |
| VRD [29] | ViP-CNN [27] | 22.78 | 27.91 | 17.32 | 20.01 |
| | DR-Net [5] | 19.93 | 23.45 | 17.73 | 20.88 |
| | U+W+SF+LK: T+S [45] | 26.32 | 29.43 | 19.17 | 21.34 |
| | Factorizable Net [25] | 26.03 | 30.77 | 18.32 | 21.20 |
| | **KB-GAN** | **27.39** | **34.38** | **20.31** | **25.01** |
| VG-MSDN [26] | ISGG [41] | 15.87 | 19.45 | 8.23 | 10.88 |
| | MSDN [26] | 19.95 | 24.93 | 10.72 | 14.22 |
| | Graph R-CNN [42] | – | – | 11.40 | 13.70 |
| | Factorizable Net [25] | 22.84 | 28.57 | 13.06 | 16.47 |
| | **KB-GAN** | **23.51** | **30.04** | **13.65** | **17.57** |



Figure 5: Qualitative results from KB-GAN. In each example, the left image is the original input image; the scene graph is generated by KB-GAN; and the right image is reconstructed from the detected objects.

bel as training targets. With the explicit image-level supervision, the proposed image reconstruction path can utilize the large quantities of instances of rare classes. This image-level supervision idea is generic, which can apply to many potential applications such as object detection.

**Comparison with Existing Methods.** Table 5 shows the comparison of our approach with the existing methods. We can see that our proposed method outperforms all the existing methods in the recall on both datasets. Compared with these methods, our model recognizes the objects and their relationships not only in the graph domain but also in the image domain.

### 4.5. Qualitative Results

Figure 5 visualizes some examples of our full-model. We show the generated scene graph as well as the reconstructed image for each sample. It is clear that our method can generate high-quality relationship predictions in the generated scene graph. Also notable is that our auxiliary output images are reasonable. This demonstrates our model's capa-

bility to generate rich scene graph by learning with both external KB and auxiliary image-level regularizer.

## 5. Conclusion

In this work, we have introduced a new model for scene graph generation which includes a novel knowledge-base feature refinement network that effectively propagates contextual information across the graph, and an image-level supervision that regularizes the scene graph generation from image domain. Our framework outperforms state-of-the-art methods for scene graph generation on VRD and VG datasets. Our experiments show that it is fruitful to incorporate the commonsense knowledge as well as the image-level supervision into the scene graph generation. Our work shows a promising way to improve high-level image understanding via scene graph.

## Acknowledgments

# References

[1] Somak Aditya, Yezhou Yang, and Chitta Baral. Explicit reasoning over end-to-end neural architectures for visual question answering. In *AAAI*, 2018. 2

[2] Sören Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary Ives. Dbpedia: A nucleus for a web of open data. In *The semantic web*, pages 722–735. 2007. 2

[3] Ankan Bansal, Karan Sikka, Gaurav Sharma, Rama Chellappa, and Ajay Divakaran. Zero-shot object detection. In *ECCV*, 2018. 2

[4] Junwei Bao, Nan Duan, Ming Zhou, and Tiejun Zhao. Knowledge-based question answering as machine translation. In *ACL*, 2014. 2

[5] Bo Dai, Yuqi Zhang, and Dahua Lin. Detecting visual relationships with deep relational networks. In *CVPR*, 2017. 1, 2, 8

[6] Jia Deng, Nan Ding, Yangqing Jia, Andrea Frome, Kevin Murphy, Samy Bengio, Yuan Li, Hartmut Neven, and Hartwig Adam. Large-scale object classification using label relation graphs. In *ECCV*, 2014. 2

[7] Henghui Ding, Xudong Jiang, Bing Shuai, Ai Qun Liu, and Gang Wang. Context contrasted feature and gated multi-scale aggregation for scene segmentation. In *CVPR*, 2018. 2

[8] Henghui Ding, Xudong Jiang, Bing Shuai, Ai Qun Liu, and Gang Wang. Semantic correlation promoted shape-variant context for segmentation. In *CVPR*, 2019. 2

[9] Desmond Elliott and Frank Keller. Image description using visual dependency representations. In *EMNLP*, 2013. 1

[10] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *IJCV*, 2010. 7

[11] Jiuxiang Gu, Jianfei Cai, Shafiq Joty, Li Niu, and Gang Wang. Look, imagine and match: Improving textual-visual cross-modal retrieval with generative models. In *CVPR*, 2018. 1

[12] Jiuxiang Gu, Jianfei Cai, Gang Wang, and Tsuhan Chen. Stack-captioning: Coarse-to-fine learning for image captioning. In *AAAI*, 2018. 1

[13] Jiuxiang Gu, Shafiq Joty, Jianfei Cai, and Gang Wang. Unpaired image captioning by language pivoting. In *ECCV*, 2018. 1

[14] Jiuxiang Gu, Gang Wang, Jianfei Cai, and Tsuhan Chen. An empirical study of language cnn for image captioning. In *ICCV*, 2017. 1

[15] Jiuxiang Gu, Zhenhua Wang, Jason Kuen, Lianyang Ma, Amir Shahroudy, Bing Shuai, Ting Liu, Xingxing Wang, Gang Wang, Jianfei Cai, et al. Recent advances in convolutional neural networks. *Pattern Recognition*, 2017. 1

[16] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *ICCV*, 2017. 6

[17] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. In *NIPS Workshop*, 2015. 2

[18] Zhiting Hu, Zichao Yang, Ruslan Salakhutdinov, and Eric Xing. Deep neural networks with massive learned knowledge. In *EMNLP*, 2016. 2

[19] Hamid Izadinia, Fereshteh Sadeghi, and Ali Farhadi. Incorporating scene context and object layout into appearance modeling. In *CVPR*, 2014. 1

[20] Justin Johnson, Agrim Gupta, and Li Fei-Fei. Image generation from scene graphs. In *CVPR*, 2018. 1, 5

[21] Justin Johnson, Ranjay Krishna, Michael Stark, Li-Jia Li, David Shamma, Michael Bernstein, and Li Fei-Fei. Image retrieval using scene graphs. In *CVPR*, 2015. 1, 2

[22] 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. In *ICCV*, 2017. 2, 6

[23] Ankit Kumar, Ozan Irsoy, Peter Ondruska, Mohit Iyyer, James Bradbury, Ishaan Gulrajani, Victor Zhong, Romain Paulus, and Richard Socher. Ask me anything: Dynamic memory networks for natural language processing. In *ICML*, 2016. 2, 4

[24] Guohao Li, Hang Su, and Wenwu Zhu. Incorporating external knowledge to answer open-domain visual questions with dynamic memory networks. In *CVPR*, 2018. 2

[25] Yikang Li, Wanli Ouyang, Bolei Zhou, Yawen Cui, Jianping Shi, and Xiaogang Wang. Factorizable net: An efficient subgraph-based framework for scene graph generation. In *ECCV*, 2018. 1, 2, 3, 4, 5, 6, 7, 8

[26] Yikang Li, Wanli Ouyang, Bolei Zhou, Kun Wang, and Xiaogang Wang. Scene graph generation from objects, phrases and region captions. In *ICCV*, 2017. 5, 6, 7, 8

[27] Yikang Li, Wanli Ouyang, Bolei Zhou, Kun Wang, and Xiaogang Wang. Vip-cnn: Visual phrase guided convolutional neural network. In *CVPR*, 2017. 1, 2, 7, 8

[28] Wentong Liao, Lin Shuai, Bodo Rosenhahn, and Michael Ying Yang. Natural language guided visual relationship detection. *arXiv preprint arXiv:1711.06032*, 2017. 1

[29] Cewu Lu, Ranjay Krishna, Michael Bernstein, and Li Fei-Fei. Visual relationship detection with language priors. In *ECCV*, 2016. 1, 2, 6, 8

[30] Jeffrey Pennington, Richard Socher, and Christopher Manning. Glove: Global vectors for word representation. In *EMNLP*, 2014. 6

[31] Bryan A Plummer, Arun Mallya, Christopher M Cervantes, Julia Hockenmaier, and Svetlana Lazebnik. Phrase localization and visual relationship detection with comprehensive image-language cues. In *ICCV*, 2017. 2

[32] Scott Reed, Zeynep Akata, Xinchen Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. Generative adversarial text to image synthesis. In *ICML*, 2016. 5

[33] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NIPS*, 2015. 3, 5, 6, 7

[34] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 6

[35] Robert Speer and Catherine Havasi. Conceptnet 5: A large semantic network for relational knowledge. In *The Peoples Web Meets NLP*, pages 161–176. 2013. 2, 4, 6

[36] Peng Wang, Qi Wu, Chunhua Shen, Anthony Dick, and Anton van den Hengel. Fvqa: Fact-based visual question answering. *PAMI*, 2018. 1

[37] Yu-Siang Wang, Chenxi Liu, Xiaohui Zeng, and Alan Yuille. Scene graph parsing as dependency parsing. In *ACL*, 2018. 1, 5

[38] Qi Wu, Chunhua Shen, Peng Wang, Anthony Dick, and Anton van den Hengel. Image captioning and visual question answering based on attributes and external knowledge. *PAMI*, 2018. 1, 2

[39] Caiming Xiong, Stephen Merity, and Richard Socher. Dynamic memory networks for visual and textual question answering. In *ICML*, 2016. 2, 4

[40] Yuanjun Xiong, Kai Zhu, Dahua Lin, and Xiaoou Tang. Recognize complex events from static images by fusing deep channels. In *CVPR*, 2015. 1

[41] Danfei Xu, Yuke Zhu, Christopher B Choy, and Li Fei-Fei. Scene graph generation by iterative message passing. In *CVPR*, 2017. 2, 8

[42] Jianwei Yang, Jiasen Lu, Stefan Lee, Dhruv Batra, and Devi Parikh. Graph r-cnn for scene graph generation. *ECCV*, 2018. 2, 7, 8

[43] Xu Yang, Kaihua Tang, Hanwang Zhang Zhang, and Jianfei Cai. Auto-encoding scene graphs for image captioning. In *CVPR*, 2019. 1

[44] Xu Yang, Hanwang Zhang, and Jianfei Cai. Shuffle-then-assemble: Learning object-agnostic visual relationship features. In *ECCV*, 2018. 1

[45] Ruichi Yu, Ang Li, Vlad I Morariu, and Larry S Davis. Visual relationship detection with internal and external linguistic knowledge distillation. In *ICCV*, 2017. 2, 8

[46] Rowan Zellers, Mark Yatskar, Sam Thomson, and Yejin Choi. Neural motifs: Scene graph parsing with global context. In *CVPR*, 2018. 2

[47] Handong Zhao, Quanfu Fan, Dan Gutfreund, and Yun Fu. Semantically guided visual question answering. In *WACV*, 2018. 1