

Devil’s on the Edges: Selective Quad Attention for Scene Graph Generation

Deunsol Jung Sanghyun Kim Won Hwa Kim Minsu Cho

Pohang University of Science and Technology (POSTECH), South Korea

<http://cvlab.postech.ac.kr/research/SQUAT>

Abstract

Scene graph generation aims to construct a semantic graph structure from an image such that its nodes and edges respectively represent objects and their relationships. One of the major challenges for the task lies in the presence of distracting objects and relationships in images; contextual reasoning is strongly distracted by irrelevant objects or backgrounds and, more importantly, a vast number of irrelevant candidate relations. To tackle the issue, we propose the Selective Quad Attention Network (SQUAT) that learns to select relevant object pairs and disambiguate them via diverse contextual interactions. SQUAT consists of two main components: edge selection and quad attention. The edge selection module selects relevant object pairs, i.e., edges in the scene graph, which helps contextual reasoning, and the quad attention module then updates the edge features using both edge-to-node and edge-to-edge cross-attentions to capture contextual information between objects and object pairs. Experiments demonstrate the strong performance and robustness of SQUAT, achieving the state of the art on the Visual Genome and Open Images v6 benchmarks.

1. Introduction

The task of scene graph generation (SGG) is to construct a visually-grounded graph from an image such that its nodes and edges respectively represent objects and their relationships in the image [30, 48, 51]. The scene graph provides a semantic structure of images beyond individual objects and thus is useful for a wide range of vision problems such as visual question answering [41, 42], image captioning [59], image retrieval [14], and conditional image generation [13], where a holistic understanding of the relationships among objects is required for high-level reasoning. With recent advances in deep neural networks for visual recognition, SGG has been actively investigated in the computer vision community. A vast majority of existing methods tackle SGG by first detecting candidate objects and then performing contextual reasoning between the objects via message pass-

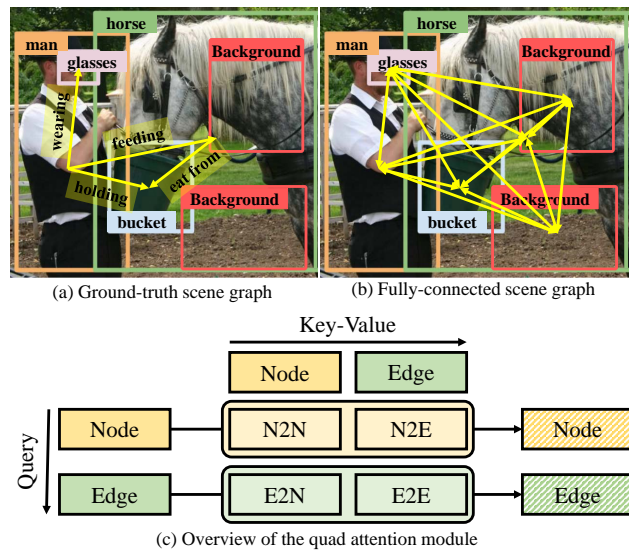


Figure 1. (a) The ground-truth scene graph contains only 4 ground-truth objects and 4 relations between the objects. (b) Only 13% of edges in a fully-connected graph have the actual relationships according to the ground-truths. (c) The overview of the quad attention. The node features are updated by node-to-node (N2N) and node-to-edge (N2E) attentions, and the edge features are updated by edge-to-node (E2N) and edge-to-edge (E2E) attentions.

ing [22, 24, 48] or sequential modeling [31, 41, 55]. Despite these efforts, the task of SGG remains extremely challenging, and even the state-of-the-art methods do not produce reliable results for practical usage.

While there exist a multitude of challenges for SGG, the intrinsic difficulty may lie in the presence of distracting objects and relationships in images; there is a vast number of potential but irrelevant relations, i.e., edges, which quadratically increase with the number of candidate objects, i.e., nodes, in the scene graph. The contextual reasoning for SGG in the wild is thus largely distracted by irrelevant objects and their relationship pairs. Let us take a simple example as in Fig. 1, where 4 objects and 4 relations in its ground-truth scene graph exist in the given image. If our

object detector obtains 6 candidate boxes, 2 of which are from the background (red), then the contextual reasoning, *e.g.*, message passing or self-attention, needs to consider 30 potential relations, 87% of which are not directly related according to the ground-truth and most of them may thus act as distracting outliers. In practice, the situation is far worse; in the Visual Genome dataset, the standard benchmark for SGG, an image contains 38 objects and 22 relationships on average [48], which means that only around 1% of object pairs have direct and meaningful relations even when object detection is perfect. As will be discussed in our experiments, we find that existing contextual reasoning schemes obtain only a marginal gain at best and often degrade the performance. The crux of the matter for SGG may lie in developing a robust model for contextual reasoning against irrelevant objects and relations.

To tackle the issue, we propose the *Selective Quad Attention Network (SQUAT)* that learns to select relevant object pairs and disambiguate them via diverse contextual interactions with objects and object pairs. The proposed method consists of two main components: edge selection and quad attention. The edge selection module removes irrelevant object pairs, which may distract contextual reasoning, by predicting the relevance score for each pair. The quad attention module then updates the edge features using edge-to-node and edge-to-edge cross-attentions as well as the node features using node-to-node and node-to-edge cross-attentions; it thus captures contextual information between all objects and object pairs, as shown in Figure 1 (c). Compared to previous methods [22, 24], which perform either node-to-node or node-to-edge interactions, our quad attention provides more effective contextual reasoning by capturing diverse interactions in the scene graph. For example, in the case of Fig. 1 (a), [‘man’, ‘feeding’, ‘horse’] relates to [‘man’, ‘holding’, ‘bracket’] and [‘horse’, ‘eat from’, ‘bracket’], and vice versa; node-to-node or node-to-edge interactions are limited in capturing such relations between the edges.

Our contributions can be summarized as follows:

- We introduce the edge selection module for SGG that learns to select relevant edges for contextual reasoning.
- We propose the quad attention module for SGG that performs effective contextual reasoning by updating node and edge features via diverse interactions.
- The proposed SGG model, SQUAT, outperforms the state-of-the-art methods on both Visual Genome and Open Images v6 benchmarks. In particular, SQUAT achieves remarkable improvement on the SGGDet settings, which is the most realistic and challenging.

2. Related work

Scene graph generation The vast majority of SGG methods [22, 24, 55] predict scene graphs in two stages: object detection and contextual reasoning. While the first stage is typically done by a pre-trained detection module [2, 36], contextual reasoning is performed by different types of message passing [3, 6, 15, 22–24, 26, 35, 46–48, 53, 54], which uses a graph neural network with node-to-edge and edge-to-node attentions, or sequential modeling [10, 31, 41, 55], which updates the node features with node-to-node attention and constructs edge features with edge-to-node attention. Unlike the previous methods, we propose quad attention, which comprises node-to-node, node-to-edge, edge-to-node, and edge-to-edge interactions, to capture all types of context exchange between candidate objects and their pairs for relational reasoning. In contextual reasoning, most of the methods consider all the candidate object pairs, *i.e.*, a fully-connected graph whose nodes are candidate objects. While Graph R-CNN [51] proposes a relation proposal network that prunes the edges from a fully-connected graph, it focuses on reducing the cost of message passing and does not analyze the effect of edge selection on the performance of scene graph generation. In contrast, we introduce an effective edge selection method and provide an in-depth analysis of it. On the other hand, since dataset imbalance/bias has recently emerged as a critical bottleneck for learning SGG¹, several methods [4, 19, 39, 45, 52] propose to adopt the techniques from long-tailed recognition, *e.g.*, data resampling [8, 22] and loss reweighting [11, 18, 50].

Transformers for vision tasks and graph structures

Transformers [43] have been adapted to the various computer vision tasks, *e.g.*, object classification [9, 29], object detection [2, 29, 37, 60] and segmentation [29, 58], and also extended for graph structures [16, 25, 32, 38]. Despite their success, vision transformer networks typically suffer from high complexity and memory consumption. Several variants of transformer networks [5, 17, 37, 44, 60] have been proposed to tackle the issue and showed that a proper sparsification technique, *e.g.*, Sparse DETR [37], can not only reduce the cost of computation and memory but also improve the task performance. Our transformer network is designed to perform contextual reasoning for scene graph generation by capturing the inherent relationships between objects and relevant object pairs, and unlike existing sparsification methods, which focus on token pruning [37] or local attention [17, 60], our edge selection module prunes not only the query edges to update but also the key-value edges used for updating.

¹For example, in the Visual Genome dataset, the most frequent entity class is 35 times larger than the least frequent one, and the most frequent predicate class is 8,000 times larger than the least frequent one.

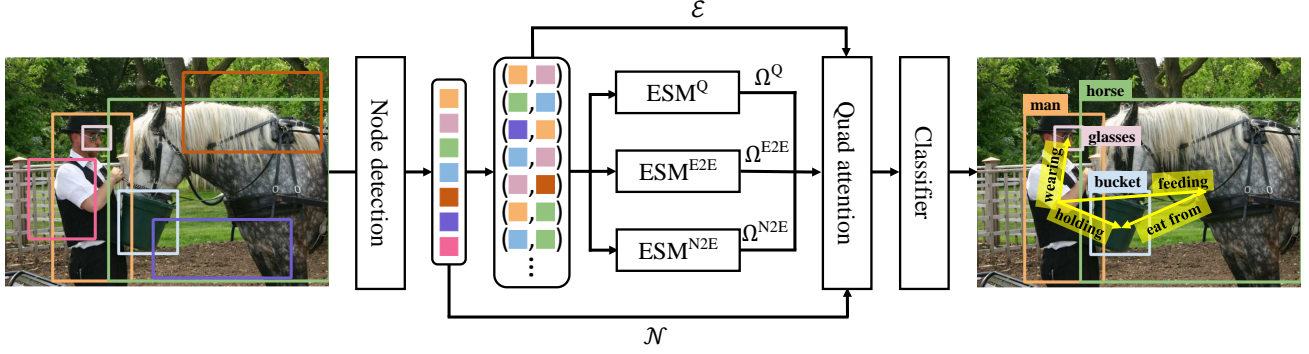


Figure 2. The overall architecture of Selective Quad Attention Networks (SQUAT). SQUAT consists of three components: the node detection module, the edge selection module, and the quad attention module. First, the node detection module extracts nodes \mathcal{N} by detecting object candidate boxes and extracting their features. Also, all possible pairs of the nodes are constructed as initial edges \mathcal{E} . Second, the edge selection module select valid edges ($\Omega^Q, \Omega^{E2E}, \Omega^{N2E}$) with high relatedness scores. Third, the quad attention module updates the node and edge features via four types of attention. Finally, the output features are passed into a classifier to predict the scene graph. See Sec. 4 for the details.

3. Problem Definition

Given an image I , the goal of SGG is to generate a visually grounded graph $G = (\mathcal{O}, \mathcal{R})$ that represents objects \mathcal{O} and their semantic relationships \mathcal{R} for object classes \mathcal{C} and predicate classes \mathcal{P} . An object $o_i \in \mathcal{O}$ is described by a pair of a bounding box $b_i \in [0, 1]^4$ and its class label $c_i \in \mathcal{C}$: $o_i = (b_i, c_i)$. A relationship $r_k \in \mathcal{R}$ is represented by a triplet of a subject $o_i \in \mathcal{O}$, an object $o_j \in \mathcal{O}$, and a predicate label $p_{ij} \in \mathcal{P}$: $r_k = (o_i, o_j, p_{ij})$, which represents relationship p_{ij} between subject o_i and object o_j .

4. Selective Quad Attention Networks

To generate semantically meaningful scene graphs as described in Section 3, we propose the Selective Quad Attention Network (SQUAT) that consists of three main components as shown in Fig. 2: the node detection module (Sec. 4.1), the edge selection module (Sec. 4.2), and the quad attention module (Sec. 4.3). First, the node detection module establishes nodes for a scene graph by detecting object candidate boxes and extracting their features. All possible pairs of the nodes are constructed as potential edges. Second, among all the potential edges, the edge selection module selects valid edges with high relatedness scores. Third, the quad attention module updates the features of nodes and valid edges via four types of attention: node-to-node (N2N), node-to-edge (N2E), edge-to-node (E2N), and edge-to-edge (E2E). For the quad attention module, we use three edge selection modules: query edge selection module for entire quad attention (ESM^Q) and key-value edge selection modules for N2E attention (ESM^{N2E}) and E2E attention (ESM^{E2E}). The nodes and edges may require different sets of edges to update their features, and some pruned edges may help to update nodes or selected edges.

For example, an edge between a person and a background, e.g., an ocean, is invalid but can help to predict the relationships between a person and other objects. Only the valid edges extracted by ESM^{N2E} and ESM^{E2E} are used to update the features of the nodes and valid edges from ESM^Q . Finally, the output features are then passed into a classifier to predict relationship classes. The remainder of this section presents the details of each component and training procedure (Sec. 4.4). In this section, the calligraphic font, i.e., \mathcal{N} and \mathcal{E} , denotes a set of features while the italic, i.e., N and E , denotes a matrix of stacked features of the set.

4.1. Node detection for object candidates

Given an image I , we use a pre-trained object detector, i.e., Faster R-CNN [36] in our experiments, to extract object bounding boxes and their class labels. Let $b_i \in [0, 1]^4$ be the i -th object box coordinate and $v_i \in \mathbb{R}^{d_v}$ its visual feature where d_v is the dimension of the visual feature. We construct a node feature f_i by transforming b_i and v_i via

$$f_i = W_o[W_v v_i; W_g b_i], \quad (1)$$

where W_o , W_v , and W_g are linear transformation matrices and $[\cdot; \cdot]$ is a concatenation operation. The edge feature f_{ij} is formed by concatenating two node features f_i and f_j and performing a linear transformation as

$$f_{ij} = W_p[f_i; f_j], \quad (2)$$

where W_p is the linear transformation matrix. As in Fig. 2, the set of entire node features $\mathcal{N} = \{f_i | 1 \leq i \leq n\}$ and the set of all possible edge features $\mathcal{E} = \{f_{ij} | 1 \leq i, j \leq n, i \neq j\}$ are passed into the edge selection and quad attention modules, whose details are described below. We denote the stacks of the features in \mathcal{N} and \mathcal{E} as N and E for the sake of simplicity.

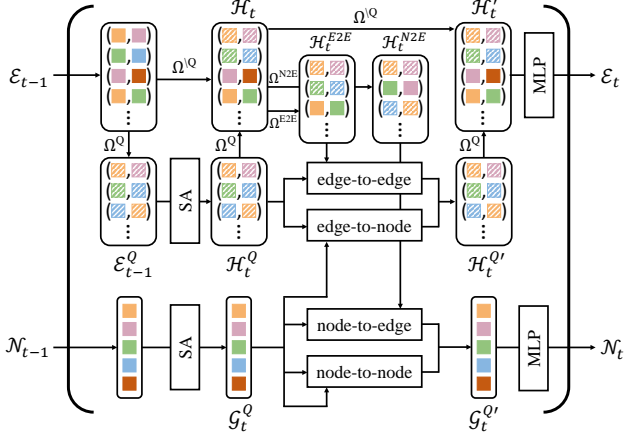


Figure 3. Detailed architecture of the quad attention. The node features are updated by node-to-node and node-to-edge attentions, and the valid edge features, selected by ESM^Q , are updated by edge-to-node and edge-to-edge attentions. The key-value of node-to-edge and edge-to-edge attentions are selected by ESM^{N2E} and ESM^{E2E} . See Sec. 4.3 for the details. Best viewed in color.

4.2. Edge selection for relevant object pairs

While node features N and edge features E can be updated via attentive message passing, a large number of irrelevant edges in E interferes with the attention process. We thus propose to prune invalid edges (*i.e.*, non-existing/false relationships) before proceeding to the quad attention module, which will be described in the next subsection.

In order to remove such distracting edges, we introduce an edge selection module (ESM) that takes an edge feature f_{ij} between nodes i and j and predicts its relatedness score s_{ij} using a simple multi-layer perceptron. We choose the pairs with top- $\rho\%$ highest relatedness scores as valid edges to use in the following quad attention module.

As mentioned earlier, we use three edge selection modules: ESM^Q , ESM^{N2E} , and ESM^{E2E} . Each edge selection module ESM^a takes the initial edge features \mathcal{E} as inputs and outputs the valid edge index set Ω for each module, resulting in Ω^Q , Ω^{E2E} , and Ω^{N2E} .

4.3. Quad attention for relationship prediction

To capture contextual interactions between the nodes and the edges, we propose a quad attention scheme inspired by the transformer decoder [43]. The main component of the quad attention is multi-head attention:

$$\text{MHA}(Q, K, V) = [\text{HA}_1; \dots; \text{HA}_h]W^O \quad (3)$$

$$\text{HA}_i = \text{softmax} \left(\frac{(QW_i^Q)(KW_i^K)^T}{\sqrt{d_k}} \right) VW_i^V, \quad (4)$$

where Q, K , and V are query, key, and value matrices. W_i^Q, W_i^K , and W_i^V are learnable transformation parameters for Q, K , and V , respectively, d_k is the dimension of the query vector, and W^O is a learnable transformation parameter for the output. Each attention head HA_i captures the information from different representation subspaces in parallel, and the multi-head attention aggregates them.

Fig. 3 shows the architecture of our quad attention layer. Following the transformer decoder layer, the t -th quad attention layer takes output edge features E_{t-1} and node features N_{t-1} from the $(t-1)$ -th layer as its input and update them with a self-attention first. Instead of updating all possible edge features E_{t-1} , we only update the valid edge features E_{t-1}^Q , whose indices are in Ω^Q extracted from ESM^Q :

$$G_t = \text{LN}(N_{t-1} + \text{MHA}(N_{t-1}, N_{t-1}, N_{t-1})), \quad (5)$$

$$H_t^Q = \text{LN}(E_{t-1}^Q + \text{MHA}(E_{t-1}^Q, E_{t-1}^Q, E_{t-1}^Q)), \quad (6)$$

where LN is layer normalization, G_t and H_t are the output of the self-attention layer for node features and valid edge features, respectively. For key-value edge features of N2E and E2E attentions, we extract the key-value set from the updated entire edge set $\mathcal{H}_t = \mathcal{H}_t^Q \cup \mathcal{E} \setminus \mathcal{E}^Q$, where \mathcal{H}_t^Q is the set of updated valid edges for query and $\mathcal{E} \setminus \mathcal{E}^Q$. Then, H_t^Q are refined by E2N and E2E attentions and G_t are refined by N2N and N2E attentions:

$$G'_t = \text{LN}(G_t + \underbrace{\text{MHA}(G_t, G_t, G_t)}_{\text{node-to-node attention}} + \underbrace{\text{MHA}(G_t, H_t^{\text{N2E}}, H_t^{\text{N2E}})}_{\text{node-to-edge attention}}), \quad (7)$$

$$H_t^{Q'} = \text{LN}(H_t^Q + \underbrace{\text{MHA}(H_t^Q, G_t, G_t)}_{\text{edge-to-node attention}} + \underbrace{\text{MHA}(H_t^Q, H_t^{\text{E2E}}, H_t^{\text{E2E}})}_{\text{edge-to-edge attention}}), \quad (8)$$

where H_t^{N2E} and H_t^{E2E} are selected by the indices Ω^{N2E} and Ω^{E2E} from the stack of \mathcal{H}_t , *i.e.*, H_t . Each attention explicitly represents a particular type of relationship between edges and nodes and helps to construct contextual information for the scene graph generation. Lastly, G'_t and H'_t are further updated by multi-layer perceptron (MLP) followed by the residual connection and a layer normalization:

$$N_t = \text{LN}(G'_t + \text{MLP}(G'_t)) \quad (9)$$

$$E_t = \text{LN}(H'_t + \text{MLP}(H'_t)), \quad (10)$$

where H'_t is the stack of $\mathcal{H}'_t = \mathcal{H}^{Q'} \cup \mathcal{E} \setminus \mathcal{E}^Q$, and the quad attention layer outputs N_t and E_t .

The inputs N_0 and E_0 of the first quad attention layer are the entire node features N and all possible edge features E , which are defined in Sec. 4.1. Every quad attention layer

uses the same valid edge sets to update the node features and valid edge features by four types of attention. Given the output edge features E_T of the last T -th quad attention layer, each edge feature $e_{ij} \in \mathcal{E}_T$ is passed into a feedforward MLP to produce a probability distribution y_{ij} over the predicate classes \mathcal{P} .

4.4. Training objective

To train SQUAT, we use a combination of two loss functions: a cross-entropy loss for the predicate classification and a binary cross-entropy loss for the edge selection module. The first predicate classification loss is defined as:

$$\mathcal{L}_{\text{PCE}} = \frac{1}{|\mathcal{E}|} \sum_{i,j=1,i \neq j}^{|\mathcal{N}|} \mathcal{L}_{\text{CE}}(y_{ij}, \hat{y}_{ij}), \quad (11)$$

where \mathcal{L}_{CE} is the cross-entropy loss and \hat{y}_{ij} is a one-hot vector of ground-truth relationship labels \hat{p}_{ij} between object i and object j . To train the edge selection module, we use auxiliary binary cross-entropy defined as:

$$\mathcal{L}_{\text{ESM}}^a = \frac{1}{|\mathcal{E}|} \sum_{i,j=1,i \neq j}^{|\mathcal{N}|} \mathcal{L}_{\text{BCE}}(s_{ij}^a, \hat{s}_{ij}), \quad (12)$$

where \mathcal{L}_{BCE} is the binary cross-entropy loss, \hat{s}_{ij} is the binary indicator of whether object i and object j have a relationship or not, and $a \in \mathcal{A} = \{\text{Q}, \text{E2E}, \text{N2E}\}$. The entire loss is defined as:

$$\mathcal{L} = \mathcal{L}_{\text{PCE}} + \lambda \frac{1}{|\mathcal{A}|} \sum_{a \in \mathcal{A}} \mathcal{L}_{\text{ESM}}^a, \quad (13)$$

where $\lambda > 0$ is a hyper-parameter. In training, \mathcal{L}_{CE} does not affect the parameters of ESM directly due to the hard selection of ESM, and the gradient passes on to train the edge feature extraction; ESM is mainly trained by \mathcal{L}_{ESM} .

5. Experiments

In this section, we perform a diverse set of experiments to evaluate the proposed model. We use two datasets: 1) Visual Genome (VG) [21] and 2) OpenImages v6 [20] datasets to train and evaluate model performances. We intend to show that our model can be generalized over heterogeneous cases by demonstrating competitive results on the two independent datasets.

5.1. Datasets and evaluation metrics

5.1.1 Visual Genome [21]

The Visual Genome dataset is composed of 108k images with an average of 38 objects and 22 relationships per image. However, most of the predicate classes have less than 10 samples. Therefore, we adopt the widely-used VG

split [24, 55] to select the most frequent 150 object classes and 50 predicate classes. Following the [55], we first split the dataset into a training set (70%) and a test set (30%). Then, we sample 5k validation images from the training set to tune the hyperparameters. We evaluate SQUAT on three subtasks: Predicate Classification (PredCls), Scene Graph Classification (SGCls), and Scene Graph Detection (SGDet). The PredCls predicts the relationships given the ground-truth bounding boxes and object labels, the SGCls aims to predict the object labels and the relationships given the ground-truth bounding boxes, and the SGDet targets predicting the object bounding boxes, object labels, and relationships without any ground-truth. As the evaluation metrics, we adopt the mean recall@K (mR@K), as previously used in scene graph generation literature [3, 40]. mR@K is the average of recall@K for each relation. Following [48], we apply the graph-constraint, in which each object pair can have only one relationship, for evaluation.

5.1.2 OpenImages v6 [20]

The OpenImages v6 dataset has 126,368 images for the training, 1,813 images for the validation, and 5,322 images for the test. Each image in the dataset has 4.1 objects and 2.8 relationships on average. The dataset has 301 object classes and 31 predicate classes. Compared with the Visual Genome dataset, the quality of annotation is far more robust and complete. For OpenImages v6, following [20, 57], we calculate Recall@50 (R@50), weighted mean AP of relationships (wmAP_{rel}), and weighted mean AP of phrases (wmAP_{phr}) as evaluation metrics. AP_{rel} evaluates the two object bounding boxes, the subject box and the object box, and three labels, the triplets of the subject, the object, and the predicate. AP_{phr} evaluates a union bounding box of subject and object and three labels as the same as AP_{rel}. To reduce the dataset bias in evaluation, we calculate wmAP_{rel} and wmAP_{phr} with weighted average of per-relationship AP_{phr} and AP_{phr}, respectively. The weight of each relationship is calculated by their relative ratios in the validation set. The final score score_{wtd} is obtained as $0.2 \times \text{R@50} + 0.4 \times \text{wmAP}_{\text{rel}} + 0.4 \times \text{wmAP}_{\text{phr}}$.

5.2. Implementation details

As in the previous work [22, 41], we adopt ResNeXt-101-FPN as a backbone network and Faster R-CNN as an object detector. The model parameters of the pre-trained object detector are frozen during the training time. We use a bi-level sampling [22] to handle the long-tailed distribution of the datasets. The hyperparameters of bi-level sampling are set the same as in [22]. We set the hyper-parameter $\lambda = 0.1$ for the loss function. The keeping ratio ρ is set to 70% for the SGDet setting on both the Visual Genome dataset and the OpenImages v6 dataset in the training. In the early stages of training, the edge selection model is not

Methods	PredCls			SGCls			SGDet		
	mR@20	mR@50	mR@100	mR@20	mR@50	mR@100	mR@20	mR@50	mR@100
IMP+ [‡] [48]	8.9	11.0	11.8	5.2	6.2	6.5	2.8	4.2	5.3
Motifs [‡] [55]	11.5	14.6	15.8	6.5	8.0	8.5	4.1	5.5	6.8
RelDN [57]	-	15.8	17.2	-	9.3	9.6	-	6.0	7.3
VCTree [‡] [41]	12.4	15.4	16.6	6.3	7.5	8.0	4.9	6.6	7.7
MSDN [24]	-	15.9	17.5	-	9.3	9.7	6.1	7.2	
GPS-Net [26]	-	15.2	16.6	-	8.5	9.1	-	6.7	8.6
RU-Net [28]	-	-	24.2	-	-	14.6	-	-	10.8
HL-Net [27]	-	-	22.8	-	-	13.5	-	-	9.2
VCTree-TDE [40]	18.4	25.4	28.7	8.9	12.2	14.0	6.9	9.3	11.1
Seq2Seq [31]	21.3	26.1	30.5	11.9	14.7	16.2	7.5	9.6	12.1
GPS-Net [†]	21.5	27.1	29.1	6.4	10.1	12.3	6.6	9.4	11.9
JMSGG [49]	-	24.9	28.0	-	13.1	14.7	-	9.8	11.8
BGNN [†] [22]	-	30.4	32.9	-	14.3	16.5	-	10.7	12.6
SQUAT [†] (Ours)	25.6	30.9	33.4	14.4	17.5	18.8	10.6	14.1	16.5

Table 1. The scene graph generation performance of three subtasks on Visual Genome (VG) dataset with graph constraints. † denotes that the bi-level sampling [22] is applied for the model. ‡ denotes that the results are reported from the [40].

Methods	R@50	wmAP		score _{wtd}
		rel	phr	
RelDN [57]	73.1	32.2	33.4	40.9
VCTree [41]	76.1	34.2	33.1	42.1
Motifs [55]	71.6	29.9	31.6	38.9
VCTree+TDE [40]	69.3	30.7	32.8	39.3
GPS-Net [26]	74.8	32.9	34.0	41.7
GPS-Net [†] [26]	74.7	32.8	33.9	41.6
BGNN [†] [22]	75.0	33.5	34.2	42.1
HL-Net [27]	76.5	35.1	34.7	43.2
RU-Net [28]	76.9	35.4	34.9	43.5
SQUAT [†]	75.8	34.9	35.9	43.5

Table 2. The scene graph generation performance on OpenImages v6 dataset with graph constraints. † denotes that the bi-level sampling [22] is applied for the model.

reliable, causing instability during training. To tackle the issue, we pre-trained the edge selection module for a few thousand iterations using \mathcal{L}_{ESM} while freezing all other parameters and then trained the entire SQUAT except for the object detector. Complete implementation details are specified in the supplementary material.

5.3. Comparison with state-of-the-art models

As shown in Table 1, on the Visual Genome dataset, SQUAT outperforms the state-of-the-art models on every setting, PredCls, SGCls and SGDet. Especially, SQUAT outperforms the state-of-the-art models by a large margin of 3.9 in mR@100 on the SGDet setting, which is the most realistic and important setting in practice as there is no perfect object detector. There are more invalid pairs in the SGDet setting than in other settings since the detected ob-

ject bounding boxes from the pre-trained object detector includes many background boxes. This means that previous work for contextual modeling was most likely distracted by the invalid pairs; thus, SQUAT shows significant performance improvement on the SGDet setting. BGNN [22] also leverage a scoring function to scale the messages of the invalid edges, however, SQUAT shows better results with our edge selection module to discard invalid object pairs. This becomes doable with the quad attention mechanism with edge selection which helps to reduce noise and outliers from invalid pairs more effectively. Also, SQUAT shows the performance improvements by 2.3 and 0.5 on the SGCls and the PredCls settings with mR@100, respectively; the more complex and realistic the task, the more noticeable the performance improvement of SQUAT becomes. It shows that SQUAT, composed of edge selection and quad attention, is appropriate for contextual reasoning to generate scene graphs even in a complex scene.

Also, as shown in Table 2, SQUAT achieve competitive results or even outperform compared with the state-of-the-art models on the OpenImages v6 dataset with score_{wtd}. Since there are fewer objects and relationships in the images of the OpenImages v6 dataset than of the Visual Genome, the edge selection module seems less effective for the OpenImages v6 dataset. As there is a trade-off in recall and mean recall when bi-level sampling is utilized [3, 31], the result of SQUAT in Table 2 is a compromised metric for R@50. But still, the R@50 of SQUAT is still competitive with that from RU-Net [28] and outperforms other recent baselines, and we achieve the best performance in wmAP_{phr} by a large margin. It shows SQUAT is effective in improving the scene graph generation performance and also in simple scenes.

Qualitative visualizations of SQUAT and more exper-

Q	Variants		SGDet		
	E2E	N2E	mR@20	mR@50	mR@100
	BGNN [22]		7.49	10.31	12.46
			9.12	12.45	15.00
✓			9.92	13.22	15.66
	✓	✓	9.84	13.04	15.60
✓	✓	✓	10.57	14.12	16.47

Table 3. The ablation study on model variants on edge selection. We remove the edge selection module for query selection and key-value selection.

Q	Variants		SGDet		
	E2E	N2E	mR@20	mR@50	mR@100
shared	shared	shared	9.61	12.70	14.85
distinct	shared	shared	9.63	12.54	14.64
distinct	distinct	distinct	10.57	14.12	16.47

Table 4. The ablation study on model variants on edge selection. We share the edge selection module for query selection and key-value selection. ‘shared’ denotes the edge selection modules share the parameters.

iments on SQUAT with 1) additional measurement, *e.g.*, recall and non-graph constraint measurement, on Visual Genome, 2) performance with plug-and-play long-tailed recognition techniques and 3) additional qualitative results are given in Supplementary.

5.4. Ablation study

5.4.1 Model variants on edge selection

To verify the effectiveness of the edge selection module, we evaluate the model from which each component of edge selection is removed on the Visual Genome dataset. As shown in Table 3, we observe that the quad attention module without the edge selection module shows much lower performance at mR@100 (-8.9%) than the full model which has the edge selection module; thus, to select the valid edges is important for the scene graph generation. On the other hand, the quad attention module without the edge selection module shows 20.4% higher performance than the BGNN and achieves 15.00 on mR@100. It shows the effectiveness of the quad attention module itself without the edge selection module. We also observe that the query selection is more critical than the key-value selection for the scene graph generation; it shows that selecting what to update is important for the scene graph generation.

To evaluate the effectiveness of three distinct edge selection modules, we evaluate the models, some of which edge selection modules are shared. In Table 4, ESM^a, of

Method				SGDet		
N2N	N2E	E2N	E2E	mR@20	mR@50	mR@100
✓	✓			7.02	9.74	11.57
✓	✓		✓	9.76	12.98	15.30
✓	✓	✓		9.70	12.27	15.03
		✓	✓	9.90	13.05	15.28
	✓	✓	✓	9.77	12.93	15.42
✓		✓	✓	9.99	13.02	15.54
✓	✓	✓	✓	10.57	14.12	16.47

Table 5. The ablation study on model variants on quad attention. N2N, N2E, E2N, E2E denote the node-to-node, node-to-edge, edge-to-node, and edge-to-edge attentions, respectively.

which $a \in \{Q, E2E, N2E\}$ are denoted ‘shared’ in the column, share the same parameters. We observe that the three fully-distinct edge selection modules boost the scene graph generation performances. It shows there exist differences between the edges needed to update both features and the edges that need to be updated.

5.4.2 Model variants on quad attention

To verify the effectiveness of the quad attention module, we evaluate the model from which each attention is removed on the Visual Genome dataset. As shown in Table 5, the full model with quad attention outperforms the other model variants. We also observe that the SQUAT without updating edges, *i.e.*, edge-to-node and edge-to-edge attentions, performs poorer than the SQUAT without updating nodes, *i.e.*, node-to-node and node-to-edge attentions. It shows that updating edge features with context information is important for context reasoning. Especially, SQUAT without edge-to-edge attention shows worse performance than without edge-to-node attention since edge-to-edge attention, which is neglected in the previous work, can capture high-level information.

5.5. The effect of the edge selection module

We applied the edge selection module to the scene graph generation model with message passing methods. Since it is not known which pairs of objects have a relationship, the message passing methods use the fully-connected graph in the inference time. However, we empirically observe that message passing on fully-connected graph is meaningless or even harmful for the scene graph generation. We use three baselines, IMP [48], BGNN [22], and SQUAT, for ablation study on the Visual Genome dataset. Table 6 shows that message passing through fully-connected graph is harmful for the scene graph generation. Even though BGNN uses a gating function to scale down the messages from invalid edges, it does not work well.

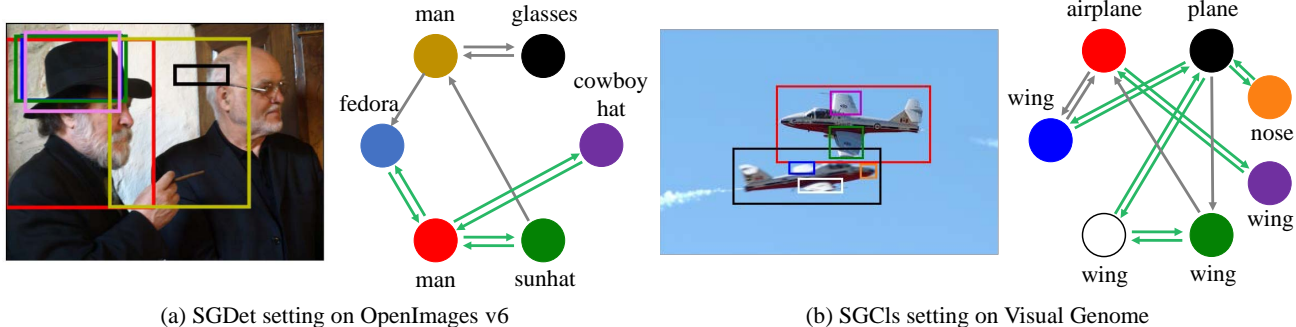


Figure 4. Qualitative results for edge selection module ESM^Q for query selection. The selected edges after edge selection are drawn in the right graph. The green arrows denote the valid pairs, and the gray arrows denote the invalid pairs. The keeping ratio for the two settings is the same $\rho = 35\%$. All of the valid edges remain, and most of the invalid edges are removed.

model	Graph	SGDet		
		mR@20	mR@50	mR@100
IMP [48]	No	4.09	5.56	6.53
	Full	2.87	4.24	5.42
BGNN [22]	No	8.99	11.84	13.56
	Full	7.49	10.31	12.46
	ES	9.00	11.86	14.20
	GT	14.15	16.41	17.09
SQUAT	No	8.68	11.52	13.99
	Full	9.12	12.45	15.00
	ES	10.57	14.12	16.47
	GT	17.95	19.21	19.51

Table 6. The ablation study on message passing for the scene graph generation. There are four settings depending on which graphs are used in the message passing: No, Full, ES, and GT. Every result is reproduced with the authors’ code.

To further investigate the effect of edge selection, we applied message passing through ground-truth scene graphs to each model. In table 6, ‘No’ represents that message passing is not used, ‘Full’ and ‘GT’ indicate that message passing is used with the complete graph and the ground-truth scene graph, respectively; ‘ES’ means that the proposed edge selection module is used with message passing. As shown in Table 6, every model with the message passing through ground truth outperforms state-of-the-art models by a substantial margin, showing that removing the invalid edges is crucial for scene graph generation. The edge selection module clearly improves not only the performance of SQUAT but also that of BGNN, the previous state-of-the-art model. It indicates that the edge selection module effectively removes the invalid edges and can be used as a plug-and-play module for message-passing-based scene graph methods.

5.6. Qualitative results

Qualitative results for the edge selection module are shown in Fig. 4. As shown in Fig. 4 (a), the object detection module extracts 6 bounding boxes, then the fully-connected graph has 30 edges in total, where only 6 valid edges are in the ground-truth. After edge selection with keeping ratio $\rho = 35\%$, only 10 edges remain where 6 valid edges all remain. It significantly reduces noises from invalid edges. The other example in Fig. 4 (b) shows the same tendency.

6. Conclusion

We presented a novel scene graph generation model that predicts a scene graph within an image. The method is designed to selectively utilize valid edges using our proposed quad attention module, and update the model from the valid edges only. The edge selection module effectively filters out invalid edges to sparsify a noisy scene graph, and thus it removes uncertainties brought by invalid edges. The quad attention module, which is composed of four components — node-to-node, node-to-edge, edge-to-node, and edge-to-edge attentions — captures the high-level information for accurately predicting relationships among different objects. We have shown the effectiveness of the SQUAT, and each component under various settings was properly validated in the experiments to demonstrate stability.

Acknowledgements. This work was supported by the IITP grants (2021-0-00537: Visual common sense through self-supervised learning for restoration of invisible parts in images (50%), 2022-0-00959: Few-shot learning of causal inference in vision and language (40%), and 2019-0-01906: AI graduate school program at POSTECH (10%)) funded by the Korea government (MSIT).

Supplementary Material

In this supplementary material, we provide additional results and details of our method, Selective Quad Attention Networks (SQUAT).

S1. Implementation details

S1.1 Code base and GPUs.

We implemented SQUAT using Pytorch [34] and some of the official code-base for BGNN [22]². SQUAT was trained for ~ 8 hours on 4 RTX 3090 GPUs with batch size 12.

S1.2 Edge selection module.

Following [37], we use simple MLP with 4 linear layers and Layer Normalization [1] with GeLU [12] activation. To capture the global statistics of the edge features $\mathcal{E} = \{f_{ij}\}_{i,j}$, we average half of the output dimensions of the first layer as a global feature g :

$$[h_{ij}^l; h_{ij}^g] = l^1(f_{ij}) \quad (14)$$

$$g = \frac{1}{|\mathcal{E}|} \sum_i \sum_j h_{ij}^g, \quad (15)$$

where l^1 is the first layer of the edge selection module and $[\cdot; \cdot]$ is the concatenation operation. The dimensions of the local part h_{ij}^l and the global part h_{ij}^g are the same. We concatenate the global feature g with each of the remaining local parts h_{ij}^l and pass into the remaining 3-layer MLP to calculate the relatedness scores s_{ij} :

$$s_{ij} = l^2([h_{ij}^l; g]), \quad (16)$$

where l^2 is the remaining 3-layer MLP. In order to remove the invalid edges, we choose top- $\rho\%$ highest relatedness score pairs \mathcal{E}^ρ as the valid edges.

S1.3 Training details.

To train SQUAT, we use Stochastic Gradient Descent (SGD) optimizer with a learning rate 10^{-3} . In the early stages of training, notice that the edge selection model is too naive to select the valid edges to construct feasible scene graphs and therefore causes instability during training. To make the training stable, we pre-trained the edge selection module for 2000 iterations with a learning rate of 10^{-4} freezing all other parameters, and then we trained the entire SQUAT without the node detection module.

We use the keeping ratio $\rho = 0.7$ and $\rho = 0.35$ in training time and inference time, respectively, for all the SGM settings on the Visual Genome and the Open Images v6 datasets. Also, we use the keeping ratio $\rho = 0.9$ for the

²<https://github.com/SHTUPLUS/PySGG>

SGCs and the PredCls settings on Visual Genome. Since the background proposals do not exist in the SGCs and the PredCls settings, there are fewer invalid edges than in the SGM setting; thus, we use a smaller keeping ratio. We use three quad attention layers for the SGM setting and two quad attention layers for the SGCs and the PredCls settings.

S2. Additional evaluations on Visual Genome

S2.1 Trade-off between recall and mean recall

Since the Visual Genome dataset³ has extremely long-tailed distribution, there is the trade-off between recall and mean recall [31, 40]. To evaluate various trade-offs of the scene graph generation methods, Zhang *et al.* [56] propose the F@K measure, the harmonic mean of recall and mean-recall, recently. Table S2 shows the R@50/100, mR@50/100, and F@50/100 on the Visual Genome dataset. SQUAT outperforms all of the state-of-the-art methods at F@50/100 measurements. It shows that although the recall of SQUAT degrades, the trade-off between the recall and the mean recall is the best in the state-of-the-art methods.

S2.2 Mean recall with no-graph constraints

Following [33, 55], we also evaluate SQUAT without the graph constraint, *i.e.*, each edge can have multiple relationships. For each edge, while mR@K evaluates only one predicate with the highest score, ng-mR@K evaluates all 50 predicates. As shown in Table. S3, on the Visual Genome dataset, SQUAT outperforms the state-of-the-art models. Especially, SQUAT outperforms the state-of-the-art models by a large margin of ng-mR@K on the SGM settings as it does in the evaluation of mR@K.

S2.3 Recall for head, body, and tail classes

Following [22], we split the relationship classes into three sets according to the number of relationship instances: head (more than 10k), body (0.5k~10k), and tail (less than 0.5k) classes. Table S4 shows the mR@100 for each group. SQUAT outperforms the state-of-the-art methods for body and tail classes by a large margin. Especially for the tail classes, SQUAT achieves twice mR@100 as that of BGNN. It shows that the scene graphs from SQUAT have more meaningful predicates, *i.e.*, tail classes such as ‘walking in’, instead of general predicates, *i.e.*, head classes such as ‘on’.

³The most frequent entity class is 35 times larger than the least frequent one and the most frequent predicate class is 8,000 times larger than the least frequent one.

Methods	PredCls			SGCls			SGDet		
	R@50 / 100	mR@50/100	F@50 / 100	R@50 / 100	mR@50/100	F@50 / 100	R@50 / 100	mR@50/100	F@50 / 100
IMP+ [‡] [24]	61.1 / 63.1	11.0 / 11.8	18.6 / 19.9	37.5 / 38.5	6.2 / 6.5	10.6 / 11.1	25.9 / 31.2	4.2 / 5.2	7.2 / 8.9
Motifs [‡] [55]	66.0 / 67.9	14.6 / 15.8	23.9 / 25.6	39.1 / 39.9	8.0 / 8.5	13.3 / 14.0	32.1 / 36.9	5.5 / 6.8	9.4 / 11.5
Motifs ^{†‡} [55]	64.6 / 66.7	18.5 / 20.0	28.8 / 30.8	37.9 / 38.8	11.1 / 11.8	17.2 / 18.1	30.5 / 35.4	8.2 / 9.7	12.9 / 15.2
RelDN [57]	64.8 / 66.7	15.8 / 17.2	25.4 / 27.3	38.1 / 39.3	9.3 / 9.6	15.0 / 15.4	31.4 / 35.9	6.0 / 7.3	7.2 / 8.9
VCTree [‡] [41]	65.5 / 67.4	15.4 / 16.6	24.9 / 26.6	38.9 / 39.8	7.4 / 7.9	12.4 / 13.2	31.8 / 36.1	6.6 / 7.7	10.9 / 12.7
MSDN [24]	64.6 / 66.6	15.9 / 17.5	25.5 / 27.7	38.4 / 39.8	9.3 / 9.7	15.0 / 15.6	31.9 / 36.6	6.1 / 7.2	10.2 / 12.0
GPS-Net [26]	65.2 / 67.1	15.2 / 16.6	24.7 / 26.6	39.2 / 37.8	8.5 / 9.1	14.0 / 14.7	31.1 / 35.9	6.7 / 8.6	11.0 / 13.9
RU-Net [28]	<u>67.7 / 69.6</u>	- / 24.2	- / 35.9	42.4 / 43.3	- / 14.6	- / 21.8	<u>32.9 / 37.5</u>	- / 10.8	- / 16.8
HL-Net [27]	60.7 / 67.0	- / 22.8	- / 34.0	<u>42.6 / 43.5</u>	- / 13.5	- / 20.6	33.7 / 38.1	- / 9.2	- / 14.8
VCTree-TDE [40]	47.2 / 51.6	25.4 / 28.7	33.0 / 36.9	25.4 / 27.9	12.2 / 14.0	16.5 / 18.6	19.4 / 23.2	9.3 / 11.1	12.6 / 15.0
Seq2Seq [31]	66.4 / 68.5	26.1 / 30.5	37.5 / 42.2	38.3 / 39.0	<u>14.7 / 16.2</u>	<u>21.2 / 22.9</u>	30.9 / 34.4	9.6 / 12.1	14.6 / 17.9
GPS-Net ^{†‡} [26]	64.4 / 66.7	19.2 / 21.4	29.6 / 32.4	37.5 / 38.6	11.7 / 12.5	17.8 / 18.9	27.8 / 32.1	7.4 / 9.5	11.7 / 14.7
JMSGG [49]	70.8 / 71.7	24.9 / 28.0	36.8 / 40.3	43.4 / 44.2	13.1 / 14.7	20.1 / 22.1	29.3 / 32.3	9.8 / 11.8	14.7 / 17.3
BGNN [22] [†]	59.2 / 61.3	<u>30.4 / 32.9</u>	40.2 / 42.8	37.4 / 38.5	14.3 / <u>16.5</u>	20.7 / <u>23.1</u>	31.0 / 35.8	<u>10.7 / 12.6</u>	<u>15.9 / 18.7</u>
SQUAT [†] (Ours)	55.7 / 57.9	30.9 / 33.4	<u>39.7 / 42.4</u>	33.1 / 34.4	17.5 / 18.8	22.9 / 24.3	24.5 / 28.9	14.1 / 16.5	17.9 / 21.0

Table S2. Recall, mean recall and F score of three subtasks on Visual Genome (VG) dataset with graph constraints. † denotes that the bi-level sampling [22] is applied for the model. ‡ denotes that the results are reported from the [3]. Bold numbers indicate the best performances and underlined numbers indicate the second best performances.

Methods	PredCls			SGCls			SGDet		
	ng-mR@20	ng-mR@50	ng-mR@100	ng-mR@20	ng-mR@50	ng-mR@100	ng-mR@20	ng-mR@50	ng-mR@100
IMP+ ^{†*} [24]	-	20.3	28.9	-	12.1	16.9	-	5.4	8.0
Frequency ^{†*} [55]	-	24.8	37.3	-	13.5	19.6	-	5.9	8.9
Motifs ^{†*} [55]	-	27.5	37.9	-	15.4	20.6	-	9.3	12.9
KERN [3]	-	36.3	49.0	-	19.8	26.2	-	11.7	16.0
GB-NET- β [54]	-	44.5	58.7	-	25.6	32.1	-	11.7	16.6
Motifs [55]	19.9	32.8	44.7	11.3	19.0	25.0	7.5	12.5	16.9
VCTree [41]	21.4	35.6	47.8	14.3	23.3	31.4	7.5	12.5	16.7
VCTree-TDE [40]	20.9	32.4	41.5	12.4	19.1	25.5	7.8	11.5	15.2
GPS-Net ^{†*} [26]	29.4	45.4	57.1	8.3	15.9	23.1	7.9	12.1	16.7
SQUAT [†]	31.8	46.0	57.8	18.7	27.1	32.6	12.1	17.9	22.5

Table S3. The scene graph generation performance of three subtasks on the Visual Genome (VG) dataset without graph constraints. † denotes that the bi-level sampling [22] is applied for the model. * denotes that the model is reproduced with the authors' code. ‡ denotes that the results are reported from the [3]. Models in the first group use pre-trained Faster R-CNN with VGG16 backbone. Bold numbers indicate the best performances.

model	head	body	tail	mR@100	R@100	F@100
VCTree-TDE [40]	24.5	13.9	0.1	11.1	23.2	15.0
GPSNet [†] [26]	30.4	8.5	3.8	9.5	32.1	14.7
BGNN [†] [22]	34.0	12.9	6.0	12.6	35.8	18.6
SQUAT [†] (Ours)	29.5	16.4	12.4	16.5	28.9	21.0

Table S4. mR@100 on the SGDet setting for head, body, and tail classes. † denotes that the bi-level sampling is applied on the model to achieve these results. Bold numbers indicate the best performances.

S2.4 Recall on simple, moderate, and complex scenes

As shown in Tables 1 and 2 in the main paper, the SQUAT shows exceptionally high performance on the most complicated task, *i.e.*, SGDet, and the most complex dataset, *i.e.*,

Visual Genome. Furthermore, to analyze the performance on the complexity of the scene, we divide the image sets in the Visual Genome into three disjoint sets according to the number of objects in the scene: simple (≤ 9), moderate (10 \sim 16), and complex (≥ 17). As shown in Table S5, the SQUAT shows a higher performance gain on the more complex images; the SQUAT is more effective for realistic and complex scenes.

S3. SQUAT with off-the-shelf method

To reduce the biases of the scene graph generation datasets, many off-the-shelf methods [4, 7, 8, 11, 18, 19, 39, 45, 50, 52, 56] are proposed. For a fair comparison, we do not compare the off-the-shelf methods with SQUAT in the main paper. We applied Internal and External Data

model	simple	moderate	complex	mR@100
BGNN [18]	15.52	12.71	9.87	12.46
SQUAT	19.54	16.80	13.28	16.47
Gain (%)	25.90	32.18	34.55	32.18

Table S5. mR@100 on the simple, moderate, and complex sets.

model	SGDet		
	mR@20	mR@50	mR@100
VCTree [41]	4.9	6.6	7.7
VCTree+TDE [40]	6.3	8.6	10.3
VCTree+PCPL ‡ [50]	8.1	10.8	12.6
VCTree+DLFE [4]	8.6	11.8	13.8
VCTree+TDE+EBM [39]	7.1	9.69	11.6
Transformer+BPL+SA [11]	10.7	13.5	15.6
Transformer+HML [7]	11.4	15.0	17.7
GPSNet+IETrans+Rwt [56]	-	16.2	18.8
SQUAT +IETrans+Rwt [56]	12.0	16.3	19.1

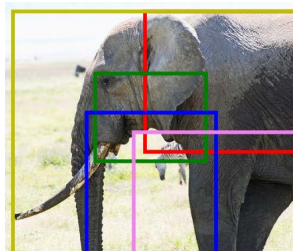
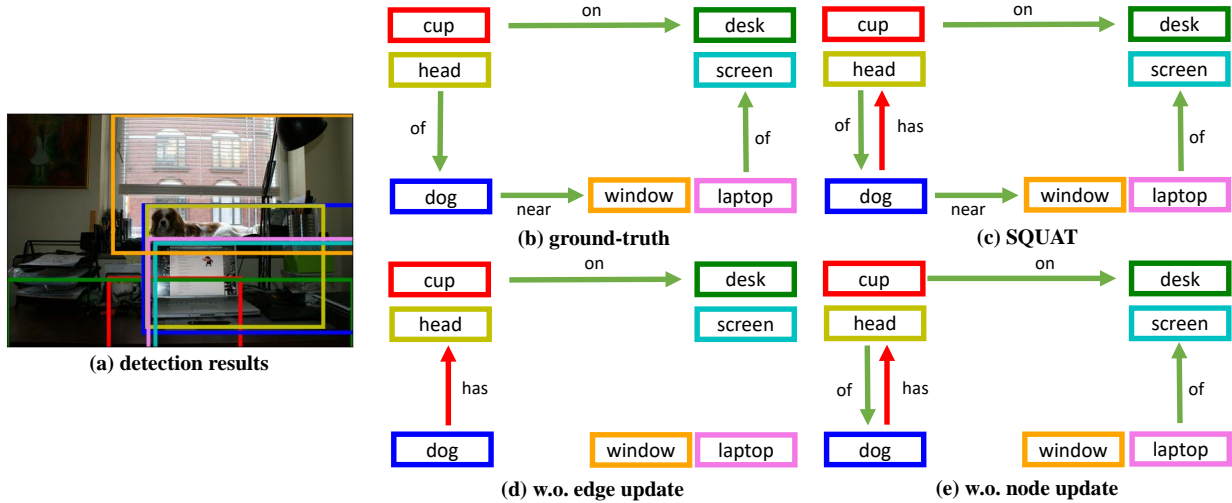
Table S6. The ablation study with the off-the-shelf learning methods on Visual Genome (VG) dataset with graph constraint. ‡ denotes that the results are reported from the [4]. The other results are from each of the original papers.

Transfer (IETrans) and reweighting (Rwt) [56], which are the state-of-the-art off-the-shelf learning methods for scene graph generation, to the SQUAT. For efficiency, we only report a model with the best performance for each off-the-shelf method. As shown in Table S6, without careful hyperparameter search, SQUAT+IETrans+Rwt model outperforms VCTree+IETrans+Rwt model and outperforms other off-the-shelf methods with Motifs [55], Transformer [40], and VCTree [41]. It shows that other off-the-shelf learning methods can be adopted for SQUAT to improve its performance.

S4. Additional qualitative results

In Fig. S2 and S3, we show the qualitative results for SQUAT model. We also compare the results of SQUAT with the results from ablated models: model without node updates and model without edge updates. The full SQUAT model shows the most informative scene graph compared to the other ablated models. There are some false positives, such as ('mouth of elephant', 'eye of elephant') in Fig. S2 bottom and ('glasses on man', 'man and woman') in Fig. S3 top, however, such errors are often caused by the incompleteness of the dataset, and hence it can be seen as a true positive. 'dog near window' in Fig. S2 top, 'zebra behind elephant' in Fig. S2 bottom, and 'man standing on sidewalk' in Fig. S3 bottom are predicted by only the full SQUAT model. It shows that quad attention modules can capture more informative contextual information.

In Fig. S4, we show the qualitative results for the edge selection module in SQUAT. The edge selection module successfully selects the valid edges. In particular, the edge selection module removes the edges between the background and the foreground, *e.g.*, most of the edges of 'sunhat' and 'scarf' are removed in Fig. S4 (a) and (b), respectively. Also, the edges between the boxes which denote the same objects are removed. For example, the edges of ('tea', 'coffee') and ('mug', 'coffee cup') are removed in Fig. S4 (d). It shows that the edge selection module successfully removed invalid edges and helps the informative message passing in the quad attention module.



(a) detection results

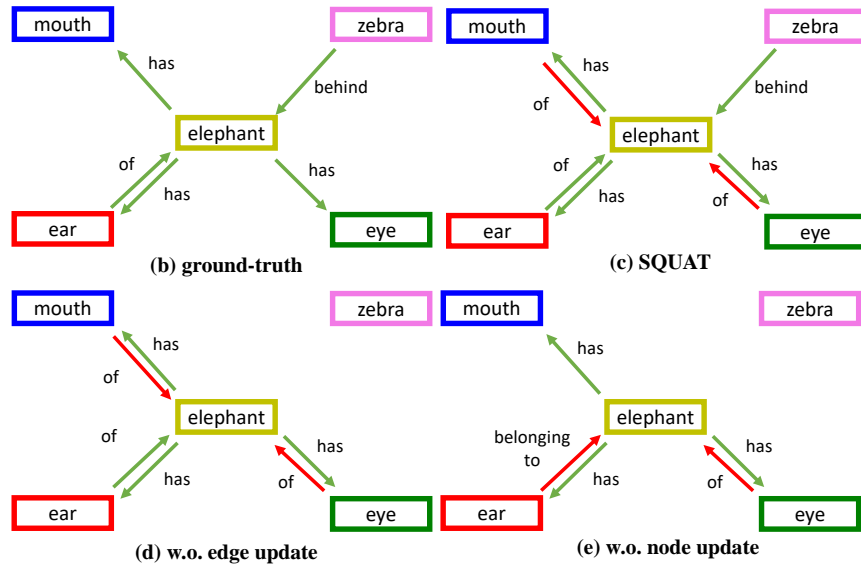


Figure S2. The qualitative results for SQUAT. (a) The detection results from pre-trained Faster R-CNN [36]. (b) The ground-truth scene graph. (c) The results from full SQUAT. (d) The results from SQUAT without edge update, *i.e.*, the edge-to-edge and the edge-to-node attentions. (e) The results from SQUAT without node update, *i.e.*, the node-to-edge and the node-to-node attentions. Full SQUAT shows more informative scene graphs than the other ablated models. The green arrows denote the true positives and the red arrows denote the false positives.

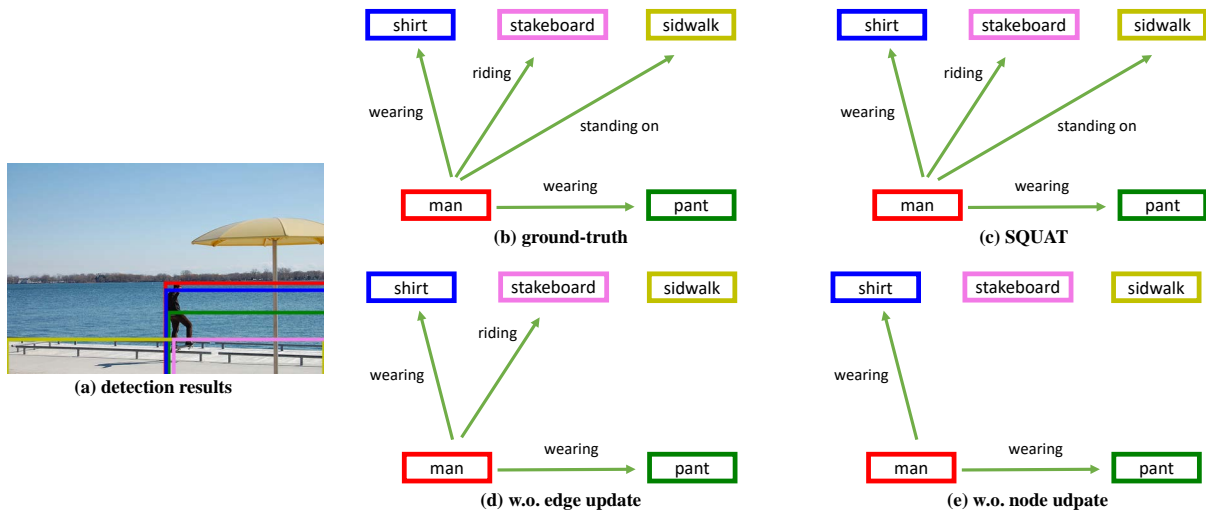
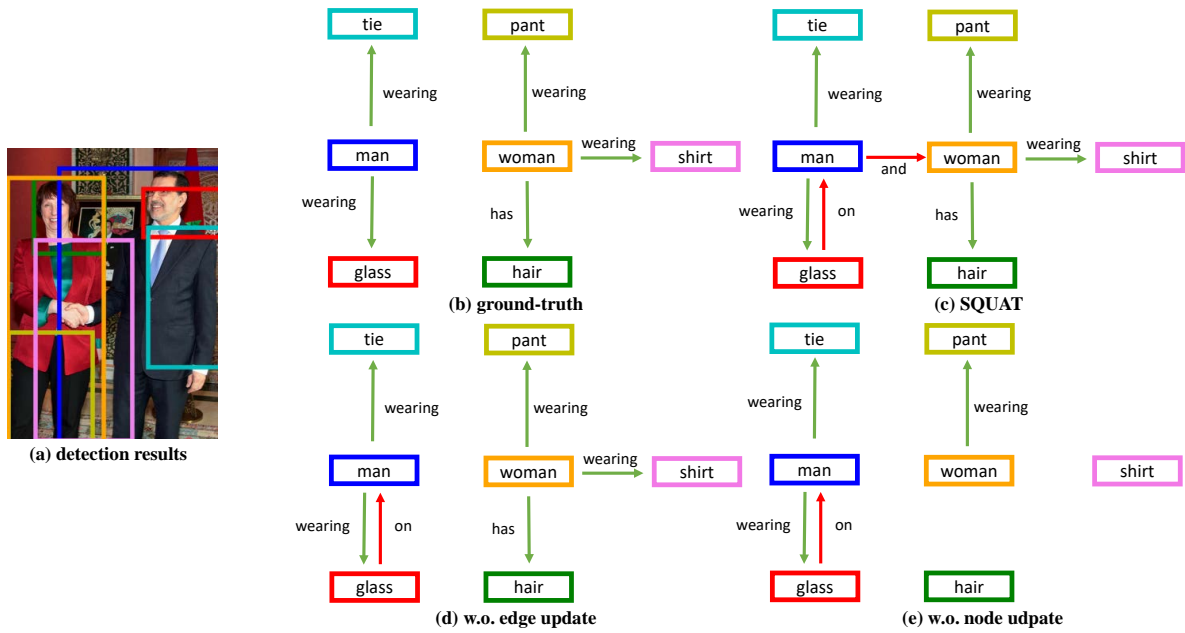


Figure S3. The qualitative results for SQUAT. (a) The detection results from pre-trained Faster R-CNN [36]. (b) The ground-truth scene graph. (c) The results from full SQUAT. (d) The results from SQUAT without edge update, *i.e.*, the edge-to-edge and the edge-to-node attentions. (e) The results from SQUAT without node update, *i.e.*, the node-to-edge and the node-to-node attentions. Full SQUAT shows more informative scene graphs than the other ablated models. The green arrows denote the true positives and the red arrows denote the false positives.

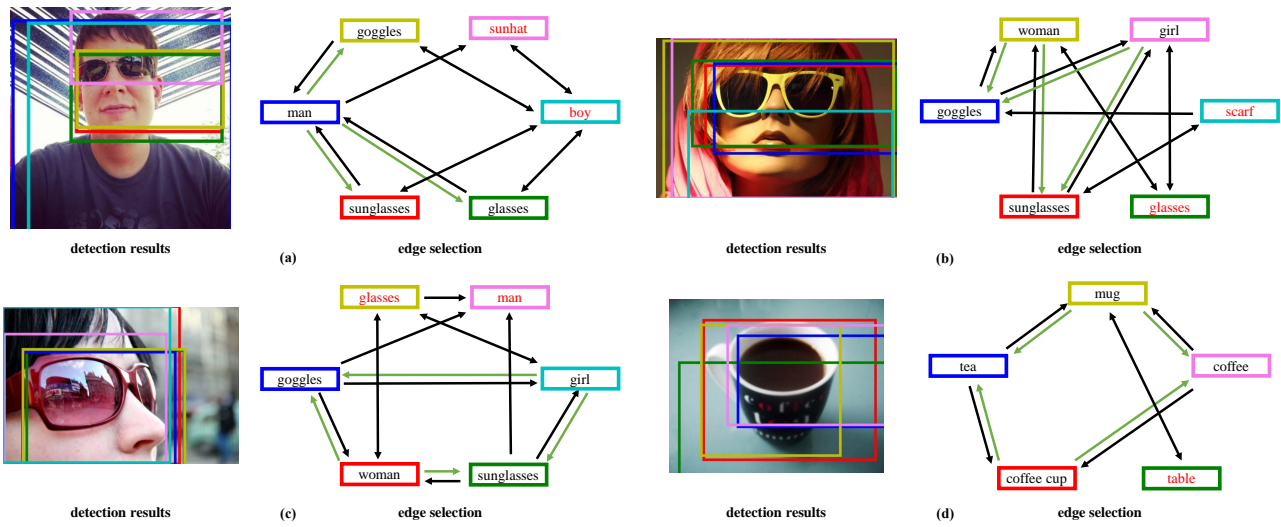


Figure S4. The qualitative results for the edge selection module on the Open Images v6 dataset. The graph denotes the results of the ESM^Q and the green arrows denote the valid edges. The boxes with the red class denote the incorrect prediction or the background.

References

- [1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. **9**
- [2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020. **2**
- [3] Tianshui Chen, Weihao Yu, Riquan Chen, and Liang Lin. Knowledge-embedded routing network for scene graph generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6163–6171, 2019. **2, 5, 6, 10**
- [4] 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, 10, 11**
- [5] Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking attention with performers. *International Conference on Learning Representations*, 2021. **2**
- [6] 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**
- [7] Youming Deng, Yansheng Li, Yongjun Zhang, Xiang Xiang, Jian Wang, Jingdong Chen, and Jiayi Ma. Hierarchical memory learning for fine-grained scene graph generation. *arXiv preprint arXiv:2203.06907*, 2022. **10, 11**
- [8] 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 (ICCV)*, pages 15404–15413, October 2021. **2, 10**
- [9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. **2**
- [10] 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/CVF conference on computer vision and pattern recognition*, pages 1969–1978, 2019. **2**
- [11] 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 (ICCV)*, pages 16383–16392, October 2021. **2, 10, 11**
- [12] Dan Hendrycks and Kevin Gimpel. Bridging nonlinearities and stochastic regularizers with gaussian error linear units. 2016. **9**
- [13] 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**
- [14] Justin Johnson, Ranjay Krishna, Michael Stark, Li-Jia Li, David Shamma, Michael Bernstein, and Li Fei-Fei. Image retrieval using scene graphs. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3668–3678, 2015. **1**
- [15] Siddhesh Khandelwal, Mohammed Suhail, and Leonid Sigal. Segmentation-grounded scene graph generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15879–15889, 2021. **2**
- [16] Jinwoo Kim, Saeyoon Oh, and Seunghoon Hong. Transformers generalize deepsets and can be extended to graphs & hypergraphs. *Advances in Neural Information Processing Systems*, 34:28016–28028, 2021. **2**
- [17] Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. *International Conference on Learning Representations*, 2020. **2**
- [18] Boris Knyazev, Harm de Vries, Cătălina Cangea, Graham W Taylor, Aaron Courville, and Eugene Belilovsky. Graph density-aware losses for novel compositions in scene graph generation. *arXiv preprint arXiv:2005.08230*, 2020. **2, 10**
- [19] Boris Knyazev, Harm de Vries, Cătălina Cangea, Graham W. Taylor, Aaron Courville, and Eugene Belilovsky. Generative compositional augmentations for scene graph prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 15827–15837, October 2021. **2, 10**
- [20] Ivan Krasin, Tom Duerig, Neil Alldrin, Vittorio Ferrari, Sami Abu-El-Haija, Alina Kuznetsova, Hassan Rom, Jasper Uijlings, Stefan Popov, Shahab Kamali, Matteo Mallocci, Jordi Pont-Tuset, Andreas Veit, Serge Belongie, Victor Gomes, Abhinav Gupta, Chen Sun, Gal Chechik, David Cai, Zheyun Feng, Dhyanesh Narayanan, and Kevin Murphy. Openimages: A public dataset for large-scale multi-label and multi-class image classification. *Dataset available from <https://storage.googleapis.com/openimages/web/index.html>*, 2017. **5**
- [21] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, Michael Bernstein, and Li Fei-Fei. Visual genome: Connecting language and vision using crowdsourced dense image annotations. 2016. **5**
- [22] 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. **1, 2, 5, 6, 7, 8, 9, 10**
- [23] Yikang Li, Wanli Ouyang, Bolei Zhou, Jianping Shi, Chao Zhang, and Xiaogang Wang. Factorizable net: an efficient subgraph-based framework for scene graph generation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 335–351, 2018. **2**
- [24] Yikang Li, Wanli Ouyang, Bolei Zhou, Kun Wang, and Xiaogang Wang. Scene graph generation from objects, phrases

- and region captions. In *Proceedings of the IEEE international conference on computer vision*, pages 1261–1270, 2017. [1](#), [2](#), [5](#), [6](#), [10](#)
- [25] Kevin Lin, Lijuan Wang, and Zicheng Liu. Mesh graphormer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 12939–12948, 2021. [2](#)
- [26] Xin Lin, Changxing Ding, Jinqian 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. [2](#), [6](#), [10](#)
- [27] Xin Lin, Changxing Ding, Yibing Zhan, Zijian Li, and Dacheng Tao. Hl-net: Heterophily learning network for scene graph generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19476–19485, 2022. [6](#), [10](#)
- [28] Xin Lin, Changxing Ding, Jing Zhang, Yibing Zhan, and Dacheng Tao. Ru-net: Regularized unrolling network for scene graph generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19457–19466, 2022. [6](#), [10](#)
- [29] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021. [2](#)
- [30] Cewu Lu, Ranjay Krishna, Michael Bernstein, and Li Fei-Fei. Visual relationship detection with language priors. In *European Conference on Computer Vision*, 2016. [1](#)
- [31] Yichao Lu, Himanshu Rai, Jason Chang, Boris Knyazev, Guangwei Yu, Shashank Shekhar, Graham W Taylor, and Maksims Volkovs. Context-aware scene graph generation with seq2seq transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15931–15941, 2021. [1](#), [2](#), [6](#), [9](#), [10](#)
- [32] Erxue Min, Runfa Chen, Yatao Bian, Tingyang Xu, Kangfei Zhao, Wenbing Huang, Peilin Zhao, Junzhou Huang, Sophia Ananiadou, and Yu Rong. Transformer for graphs: An overview from architecture perspective. *arXiv preprint arXiv:2202.08455*, 2022. [2](#)
- [33] Alejandro Newell and Jia Deng. Pixels to graphs by associative embedding. *Advances in neural information processing systems*, 30, 2017. [9](#)
- [34] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimeshain, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. [9](#)
- [35] Mengshi Qi, Weijian Li, Zhengyuan Yang, Yunhong Wang, and Jiebo Luo. Attentive relational networks for mapping images to scene graphs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3957–3966, 2019. [2](#)
- [36] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28, 2015. [2](#), [3](#), [12](#), [13](#)
- [37] Byungseok Roh, JaeWoong Shin, Wuhyun Shin, and Sae-hoon Kim. Sparse detr: Efficient end-to-end object detection with learnable sparsity. *arXiv preprint arXiv:2111.14330*, 2021. [2](#), [9](#)
- [38] Yu Rong, Yatao Bian, Tingyang Xu, Weiyang Xie, Ying Wei, Wenbing Huang, and Junzhou Huang. Self-supervised graph transformer on large-scale molecular data. *Advances in Neural Information Processing Systems*, 33:12559–12571, 2020. [2](#)
- [39] 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. [2](#), [10](#), [11](#)
- [40] 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. [5](#), [6](#), [9](#), [10](#), [11](#)
- [41] 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. [1](#), [2](#), [5](#), [6](#), [10](#), [11](#)
- [42] Damien Teney, Lingqiao Liu, and Anton van Den Hengel. Graph-structured representations for visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2017. [1](#)
- [43] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. [2](#), [4](#)
- [44] Sinong Wang, Belinda Z Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with linear complexity. *arXiv preprint arXiv:2006.04768*, 2020. [2](#)
- [45] Tzu-Jui Julius Wang, Selen Pehlivan, and Jorma Laaksonen. Tackling the unannotated: Scene graph generation with bias-reduced models. *arXiv preprint arXiv:2008.07832*, 2020. [2](#), [10](#)
- [46] Wenbin Wang, Ruiping Wang, Shiguang Shan, and Xilin Chen. Exploring context and visual pattern of relationship for scene graph generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8188–8197, 2019. [2](#)
- [47] Sanghyun Woo, Dahun Kim, Donghyeon Cho, and In So Kweon. Linknet: Relational embedding for scene graph. *Advances in Neural Information Processing Systems*, 31, 2018. [2](#)
- [48] Danfei Xu, Yuke Zhu, Christopher B Choy, and Li Fei-Fei. Scene graph generation by iterative message passing. In *Pro-*

- ceedings of the IEEE conference on computer vision and pattern recognition*, pages 5410–5419, 2017. [1](#), [2](#), [5](#), [6](#), [7](#), [8](#)
- [49] 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. [6](#), [10](#)
- [50] Shaotian Yan, Chen Shen, Zhongming Jin, Jianqiang Huang, Rongxin Jiang, Yaowu Chen, and Xian-Sheng Hua. Pcp1: Predicate-correlation perception learning for unbiased scene graph generation. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 265–273, 2020. [2](#), [10](#), [11](#)
- [51] 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. [1](#), [2](#)
- [52] Yuan Yao, Ao Zhang, Xu Han, Mengdi Li, Cornelius Weber, Zhiyuan Liu, Stefan Wermter, and Maosong Sun. Visual distant supervision for scene graph generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 15816–15826, October 2021. [2](#), [10](#)
- [53] Guojun Yin, Lu Sheng, Bin Liu, Nenghai Yu, Xiaogang Wang, Jing Shao, and Chen Change Loy. Zoom-net: Mining deep feature interactions for visual relationship recognition. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 322–338, 2018. [2](#)
- [54] 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. [2](#), [10](#)
- [55] Rowan Zellers, Mark Yatskar, Sam Thomson, and Yejin Choi. Neural motifs: Scene graph parsing with global context. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5831–5840, 2018. [1](#), [2](#), [5](#), [6](#), [9](#), [10](#), [11](#)
- [56] Ao Zhang, Yuan Yao, Qianyu Chen, Wei Ji, Zhiyuan Liu, Maosong Sun, and Tat-Seng Chua. Fine-grained scene graph generation with data transfer. In *ECCV*, 2022. [9](#), [10](#), [11](#)
- [57] Ji Zhang, Kevin J Shih, Ahmed Elgammal, Andrew Tao, and Bryan Catanzaro. Graphical contrastive losses for scene graph generation. 2019. [5](#), [6](#), [10](#)
- [58] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng Feng, Tao Xiang, Philip HS Torr, et al. Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6881–6890, 2021. [2](#)
- [59] Yiwu Zhong, Liwei Wang, Jianshu Chen, Dong Yu, and Yin Li. Comprehensive image captioning via scene graph decomposition. In *European Conference on Computer Vision*, pages 211–229. Springer, 2020. [1](#)
- [60] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *International Conference on Learning Representations*, 2021. [2](#)