

Learning to Generate Scene Graph from Natural Language Supervision

Yiwu Zhong¹, Jing Shi², Jianwei Yang³, Chenliang Xu², Yin Li¹

¹University of Wisconsin-Madison ²University of Rochester ³Microsoft Research

{yzhong52, yin.li}@wisc.edu {j.shi, chenliang.xu}@rochester.edu jianwei.yang@microsoft.com

Abstract

Learning from image-text data has demonstrated recent success for many recognition tasks, yet is currently limited to visual features or individual visual concepts such as objects. In this paper, we propose one of the first methods that learn from image-sentence pairs to extract a graphical representation of localized objects and their relationships within an image, known as scene graph. To bridge the gap between images and texts, we leverage an off-the-shelf object detector to identify and localize object instances, match labels of detected regions to concepts parsed from captions, and thus create “pseudo” labels for learning scene graph. Further, we design a Transformer-based model to predict these “pseudo” labels via a masked token prediction task. Learning from only image-sentence pairs, our model achieves 30% relative gain over a latest method trained with human-annotated unlocalized scene graphs. Our model also shows strong results for weakly and fully supervised scene graph generation. In addition, we explore an open-vocabulary setting for detecting scene graphs, and present the first result for open-set scene graph generation. Our code is available at https://github.com/YiwuZhong/SGG_from_NLS.

1. Introduction

An image might have millions of pixels, yet its visual content can be often summarized using dozens of words. Images and their text descriptions (*i.e.* captions) are available in great abundance from the Internet [41], and offer a unique opportunity of image understanding aided by natural language. Learning visual knowledge from image-text pairs has been a long-standing problem [51, 8, 14, 7, 23, 55, 59, 16, 62, 38, 11, 36], with recent success on learning deep models for visual representation [7, 23, 38, 11, 36], and for recognizing and detecting individual visual concepts (*e.g.* objects) [55, 59, 62, 36, 16]. In this paper, we ask the question: *can we learn to detect visual relationships beyond individual concepts from image-text pairs?* Fig. 1 (a) illustrates an example of such relationships (“man drive boat”).

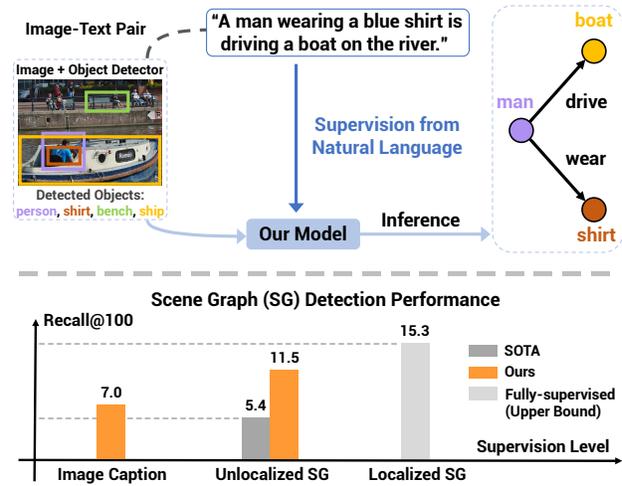


Figure 1. **Top (our setting)**: Our goal is learning to generate localized scene graphs from image-text pairs. Once trained, our model takes an image and its detected objects as inputs and outputs the image scene graph. **Bottom (our results)**: A comparison of results from our method and [61] with varying levels of supervision.

As a first step, we focus on learning scene graph generation (SGG) from image-sentence pairs. A scene graph is a symbolic and graphical representation of an image, with each graph node as a localized object and each edge as a relationship (*e.g.* a predicate) between a pair of objects. Scene graph has emerged as a structured representation for many vision tasks, including action recognition [20], 3D scene understanding [1, 50], image generation and editing [22, 13], and vision-language tasks (*e.g.* image captioning [56, 57, 67] and visual question answering [42, 48, 18]). Most previous scene graph methods [53, 29, 63, 28, 54, 4, 47, 46, 66] follow a fully supervised approach, relying on human annotations of object bounding boxes, object categories and their relationships. These annotations are very costly and difficult to scale. Recently, Zareian *et al.* [61] considered weakly supervised learning of scene graphs from image-level labels of unlocalized scene graphs. Nonetheless, learning scene graphs from images and their text descriptions remains unexplored.

A major challenge of learning scene graphs from image-

sentence pairs is the missing link between many candidate image regions and a few concepts (*e.g.* nouns and predicates) parsed from an image caption. To this end, we propose to leverage off-the-shelf object detectors, capable of identifying and localizing object instances from hundreds of common categories. Our key idea is that object labels of detected image regions can be further matched to sentence concepts, and thus provide “pseudo” labels for learning scene graphs, thereby bridging the gap between region-concept pairs. Our hypothesis is that these “pseudo” labels, coupled with a large-scale dataset, can be used for training a deep model to detect scene graph of an input image. Our language supervised setting is shown in Fig. 1 (a).

Inspired by the recent success of vision-language pre-training [9, 27, 68, 32, 44, 45, 31], we develop a Transformer-based model for learning to generate scene graphs supervised by image-sentence pairs. Specifically, our model takes inputs of visual features from a pair of detected object regions, text embeddings of their predicted categorical labels, and contextual features from other object regions, all provided by an off-the-shelf detector [37]. Our model then learns to recognize the visual relationship between the input object pair, represented as a localized subject-predicate-object (SPO) triplet. A scene graph can thus be generated by enumerating all pairs from a small set of detected objects. During training, our model learns from only image-sentence pairs using “pseudo” labels produced by matching the detected object labels to the parsed sentence concepts. During inference, our model generates a scene graph given an input image with its detection results.

Our model is trained on captioning datasets including COCO Caption [6] and Conceptual Caption [41], and evaluated on Visual Genome [24] — a widely used scene graph benchmark. Our results, summarized in Fig. 1 (b), significantly outperform the state of the art [61] on weakly supervised SGG by a relative margin of 30%, despite that our model only requires image-sentence pairs for training while [61] is trained using human-annotated unlocalized scene graphs. With the same supervision as [61], our model achieves a relative gain of 112% in recall. Further, our model also demonstrates strong results on fully supervised SGG. While these results are reported on closed-set setting with known target concepts during training, we also present promising results on open-set SGG where the concept vocabulary is crafted from image captions. Our work is among the first methods for learning to detect scene graphs from only image-sentence pairs, and presents the first result for open-set SGG. We believe our work provides a solid step towards structured image understanding.

This paper was accepted to ICCV 2021. In this arXiv version, we add additional comparison to a concurrent work from [58], provide more details during the discussion of our results, and include additional results in the appendix.

2. Related Work

We briefly review recent works on learning visual knowledge from natural language and scene graph generation, with a focus on the development of deep models.

Learning Visual Knowledge from Language. The availability of images and their text descriptions on the Internet has spurred a surge of interest in learning from image-text pairs. Early works focused on learning from image-hashtag pairs for visual representation learning [7, 23] and for recognizing objects, scenes, and actions [8, 14, 15, 26]. More recent works have shifted attention to learning from images and their sentence descriptions. For example, image-sentence pairs were used for visual representation learning via image captioning [11], image-text matching [36], or image-conditioned language modeling [38], and for visual and textual representation learning using context prediction tasks [9, 27, 68, 32, 44, 45, 31]. Image captions were exploited for object recognition [51] and object detection [59, 19, 55, 62]. Unlike these previous works, our work learns to detect localized scene graphs that encode objects and their relationships in an input image. Inspired by previous works on visual and textual representation learning, we propose a Transformer-based model for scene graph generation and formulate the problem as predicting masked tokens of subject, predicate and object.

Fully Supervised Scene Graph Generation. An image scene graph represents localized object instances as nodes and their relationships as edges on the graph. Scene graph generation (SGG) aims to extract this graphical representation from an input image. A related problem is visual relationship detection (VRD) [60, 30, 64, 10] that also localizes objects and recognizes their relationships yet without the notation of a graph. Thanks to the development of large-scale densely annotated image scene graph datasets, such as Visual Genome (VG) dataset [24], a large array of methods have been proposed for scene graph generation. Several different models have been explored, including iterative message passing [53, 29], recurrent network [63], tree structure encoding [47, 52], graph convolution and pruning [28, 54], casual reasoning [4, 46] and contrastive learning [66]. A major drawback of these approaches is the requirement of human-annotated, localized scene graphs with categorical labels and locations of all nodes and edges. Our work seeks to address this drawback by learning to detect scene graphs from only image-sentence pairs.

Weakly-supervised Scene Graph Generation. Several recent works have explored weakly supervised settings for VRD [35, 3, 65] and SGG [65, 61, 43]. Most of them addressed the task of VRD and seeks to learn from unlocalized SPO triplets. For example, Peyre *et al.* [35] proposed to assign image-level labels to pairs of detected objects via discriminative clustering. Baldassarre *et al.* [3] first predicted visual predicates given the detected objects, and then

retrieved the subjects and objects using backward explanation techniques. Zhang *et al.* [65] designed a fully convolutional network to jointly learn object detection and predicate prediction from image-level labels, using object proposals as model inputs. They reported results on both VRD and SGG. The most relevant work is from Zareian *et al.* [61]. They proposed learning from unlocalized scene graphs for SGG, and developed a message passing mechanism to update features of detected objects and to gradually refine labels of objects and predicates. Our recent work [43] presented a simple baseline for weakly supervised SGG using first-order graph matching. Similar to these approaches, our method explores learning using less labels for SGG. Unlike previous approaches, our method leverages image captions — a different type of labels that are easier to obtain than unlocalized SPO triplets or scene graphs. A concurrent work from Ye *et al.* [58] also explored learning scene graph from image-sentence pairs. They proposed to use visual grounding to iteratively match the detected image regions and the text entities parsed from captions. Unlike their method, we leverage an object detector to create the pseudo labels for SPO triplets, leading to significantly better empirical results. Our work is thus among the first methods to detect scene graphs by learning from only image-sentence pairs.

3. Scene Graph from Language Supervision

With a large collection of paired images $\{I\}$ and captions $\{S\}$, our goal is learning to detect an image scene graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ from an input image I . \mathcal{G} is a directed graph with nodes \mathcal{V} and edges \mathcal{E} . Each node $v_i \in \mathcal{V}$ denotes a localized object in I , represented by its bounding box b_i and object label o_i within a vocabulary \mathcal{C}_o^g . Each edge $e_{ij} \in \mathcal{E}$ denotes a predicate (e.g. “drive”) from a vocabulary \mathcal{C}_p^g pointing from node v_i to node v_j , where $T_{ij} = (v_i, e_{ij}, v_j)$ defines a triplet of subject-predicate-object (SPO). Scene graph generation is thus a challenging structured output prediction problem.

Similar to previous SGG methods [53, 64, 47, 63, 61], we assume a set of object regions $R = \{r_n\}$ provided by an detector. Each region $r_n = (\bar{b}_n, \bar{o}_n)$ consists of a bounding box \bar{b}_n and a predicted object category \bar{o}_n from a vocabulary \mathcal{C}_o^d given by the detector. r_n thus defines a candidate node of the target scene graph \mathcal{G} . It is worth noting that the vocabulary of the detector \mathcal{C}_o^d is different from the vocabulary of the scene graph \mathcal{C}_o^g (i.e. $\mathcal{C}_o^d \neq \mathcal{C}_o^g$). With object regions $R = \{r_n\}$, SGG is reduced to classify r_n into object categories ($\mathcal{C}_o^g \cup \{\text{background}\}$), and infer the predicate label ($\mathcal{C}_p^g \cup \{\text{background}\}$) between each subject-object region pair (r_k, r_l) . A main innovation of our model is to learn from only image-sentence pairs for SGG, without the need of ground-truth object labels nor their relationships.

Learning from Language Supervision. Our key idea is to extract SPO triplets from an image caption, and match these

Scene Graph Generation Settings	Required Annotation during Training		
	Image Description	Object&Predicate Category Labels	Object Boxes
Fully Supervised [53]		✓	✓
Weakly Supervised [61]		✓	
Language Supervised (ours)	✓		

Table 1. Our language supervised setting vs. fully supervised and weakly supervised settings. Our method learns from only image-sentence pairs to generate localized image scene graphs.

triplets to object categories of image regions given by the detector, thereby creating “pseudo” labels for these regions and their relationships. Specifically, we adopt a language parser to extract a set of triplets $\{T'\}$ from the caption S . We further link object region pairs $\{r_k, r_l\}$ in the image I provided by the detector to the parsed sentence triplets T' . This is done by matching detected object categories \bar{o}_k and \bar{o}_l from every region pair to the subject and object in each T' , respectively. If matched, the sentence triplet T' will define a “pseudo” label for the region pair (r_k, r_l) (subject, object) and their relationship e_{kl} (predicate). These “pseudo” labels can then be used to train our model.

Comparison to Fully and Weakly Supervised Settings.

Our setting of learning to generate scene graphs from image-sentence pairs is different from previous fully and weakly supervised settings, as shown in Table 1. Our setting provides a new opportunity of learning structured visual knowledge from natural language supervision.

Overview of Our Model. Our model, inspired by recent works in vision-language pretraining [9, 27, 68, 32, 44, 45, 31], seeks to label the SPO triplet given a pair of regions. Specifically, we design a Transformer-based model with its inputs as a region pair (r_k, r_l) and the contextual features from other regions $\{r_n\} - \{r_k, r_l\}$. Our model then predicts the category labels (o_k, e_{kl}, o_l) of a SPO triplet T_{kl} for the input region pair (r_k, r_l) . During training, our model is supervised by the “pseudo” labels T' parsed from caption S . During inference, our model takes inputs of the image I and its detection results $R = \{r_n\}$, labels every region pair (r_k, r_l) , and aggregates the SPO triplets into a full scene graph. Fig. 2 illustrates our model.

3.1. Triplet Transformer

Our proposed Triplet Transformer is a triplet labeling model based on an input region pair and its contextual features. Specifically, for each region $r_n = (\bar{b}_n, \bar{o}_n)$, we denote its visual, positional, and textual features as \mathbf{x}_n^r , \mathbf{x}_n^p , and \mathbf{x}_n^o , respectively. \mathbf{x}_n^r is the visual feature (ROI) pooled from the region \bar{b}_n . \mathbf{x}_n^p is a feature encoding the position of the bounding box, i.e., a 7-D vector with the normalized top/left/bottom/right coordinates, width, height and area for the region box \bar{b}_n . \mathbf{x}_n^o is the word embedding of the region object label \bar{o}_n . Given an input region pair (r_k, r_l) and all other detected regions, our model builds a composition

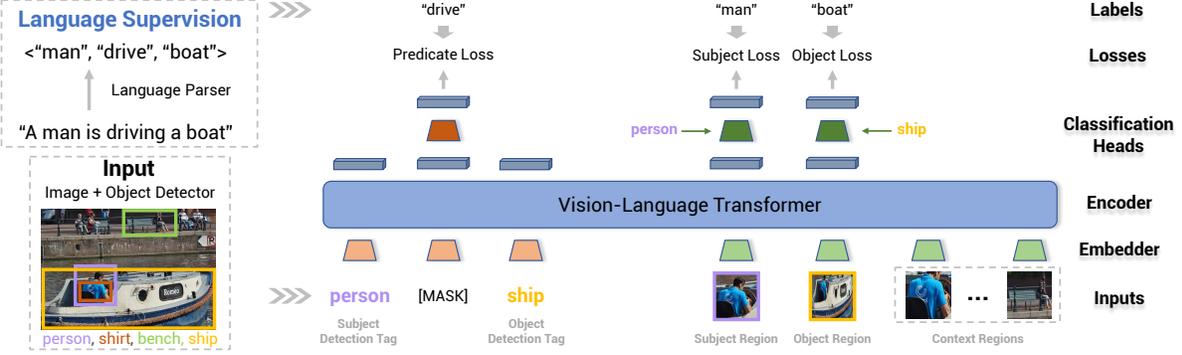


Figure 2. Overview of our proposed model for language supervised scene graph generation. Given an image, an object detector is first applied with the detected objects as the inputs to our model. Our model further embeds the detected region features and textual object categories (e.g., the tags of a pair of subject-object, the MASK representing the predicate) into token embeddings, followed by a multi-layer Transformer encoder. Finally, our model predicts the labels of the subject region, the object region and the predicate.

function $f = g \circ h$ to predict the labels (o_k, e_{kl}, o_l) of a SPO triplet, given by

$$o_k, e_{kl}, o_l = g \circ h \left(\underbrace{\mathbf{x}_k^o, \mathbf{x}_l^o}_{\text{Textual Embedder}}; \underbrace{\mathbf{x}_k^r, \mathbf{x}_l^r, \mathbf{x}_k^p, \mathbf{x}_l^p}_{\text{Visual Embedder}}; \underbrace{\{\mathbf{x}_u^r, \mathbf{x}_u^p\}_{u \neq k, l}}_{\text{Contextual Features}} \right)$$

where u indexes all regions except r_k and r_l .

Our model thus consists of: (1) a visual embedder that encodes visual and positional region features; (2) a textual embedder that embeds textual region features (from object labels); (3) a multi-layer Transformer h that conducts message passing among the input visual and textual embeddings; and (4) classification heads g that predict the labels of a triplet. We now present details for each component.

Visual Embedder. Our visual embedder transforms visual and positional features $(\mathbf{x}_n^r$ and $\mathbf{x}_n^p)$ of region r_n into an embedding \mathbf{v}_n , where n indexes all region features including k (subject), l (object) and u (context). This is given by

$$\mathbf{v}_n = \text{LN}(\text{LN}(\mathbf{W}_r \mathbf{x}_n^r) + \text{LN}(\mathbf{W}_p \mathbf{x}_n^p) + \mathbf{e}_n^t), \quad (1)$$

where \mathbf{W}_r and \mathbf{W}_p are trainable weights that project the features into the same dimension d . $\mathbf{e}_n^t \in \mathbb{R}^d$ is the type embedding of a region (subject vs. object vs. context). LN denotes Layer Normalization [2].

Textual Embedder. Our textual embedder accepts two inputs: (1) the word embeddings \mathbf{x}_k^o and \mathbf{x}_l^o of region labels for subject and object region, respectively; and (2) the word embedding of a special word “MASK”, denoted as \mathbf{x}_p^o , representing the missing predicate. The embedder encodes the input word embedding and the positional embedding into a textual embedding \mathbf{t}_m , given by

$$\mathbf{t}_m = \text{LN}(\mathbf{W}_e \mathbf{x}_m^o + \mathbf{e}_m^p), \quad (2)$$

where m indexes k (subject), p (predicate) or l (object). $\mathbf{e}_m^p \in \mathbb{R}^d$ is the positional embedding [12] of the current

token. \mathbf{W}_e represents the trainable weights projecting the word embedding into the dimension of d .

Transformer Encoder. The visual and textual embeddings $(\mathbf{v}_n$ and $\mathbf{t}_m)$ are further fed into a multi-layer Transformer encoder [49]. This encoder uses multi-head self-attention, coupled with multilayer perceptron (MLP) and layer normalization, to output a contextualized embedding $(\hat{\mathbf{v}}_n \in \mathbb{R}^d$ or $\hat{\mathbf{t}}_m \in \mathbb{R}^d)$ for each input \mathbf{v}_n or \mathbf{t}_m . This Transformer encoder can be considered as conducting message passing across all input tokens. Among all the outputs, the embeddings corresponding to the subject, predicate, object tokens will be further used for triplet label prediction, as shown in Fig. 2. For a region pair (r_k, r_l) , the embeddings $\hat{\mathbf{v}}_k / \hat{\mathbf{t}}_k$ correspond to the visual / textual feature of the subject region (*i.e.* the first input region), the predicate embedding $\hat{\mathbf{t}}_p$ is from the special word “MASK”, and the embeddings $\hat{\mathbf{v}}_l / \hat{\mathbf{t}}_l$ represent the visual / textual feature of the object region (*i.e.* the second input region).

Classification Heads. Our model further fuses the encoder outputs, and predicts labels of a SPO triplet (subject-predicate-object) for the input region pair (r_k, r_l) . The feature fusion is given by

$$\begin{aligned} \mathbf{s} &= \hat{\mathbf{v}}_k + \mathbf{W}_v \mathbf{x}_k^o, & \mathbf{o} &= \hat{\mathbf{v}}_l + \mathbf{W}_v \mathbf{x}_l^o, \\ \mathbf{p} &= \hat{\mathbf{t}}_p + \mathbf{W}_{ts} \hat{\mathbf{t}}_k + \mathbf{W}_{to} \hat{\mathbf{t}}_l + \mathbf{W}_{vs} \hat{\mathbf{v}}_k + \mathbf{W}_{vo} \hat{\mathbf{v}}_l, \end{aligned} \quad (3)$$

where \mathbf{W}_v , \mathbf{W}_{ts} , \mathbf{W}_{to} , \mathbf{W}_{vs} , \mathbf{W}_{vo} are learnable weights. The outputs $\mathbf{s} \in \mathbb{R}^d$, $\mathbf{o} \in \mathbb{R}^d$, $\mathbf{p} \in \mathbb{R}^d$ are further used to classify subject, predicate, and object labels, respectively. This is done using a two-layer MLP followed by softmax.

3.2. Learning from Language Supervision

Our key innovation is the use of image captions as the only supervisory signal for training our model. This is done by constructing “pseudo” labels of triplets from image captions. Concretely, we first parse a caption into a set of SPO triplets. Each triplet is further matched to every pair of re-

gions, by comparing subject and object tokens in the sentence triplet to the predicted categories of a region pair. Our model is then trained on the matched pairs of regions to predict their corresponding sentence triplets. We point out that our approach of learning from image-sentence pairs can be easily adapted by different SGG models.

Closed-Set vs. Open-Set. In this paper, we primarily consider a closed-set setting — the vocabulary of the subject, predicate, and object during evaluation is known in prior. In this setting, our learning is focused on the concepts of interest and our model only considers sentence triplets within the vocabulary. Nonetheless, our method does support the open-set setting, where there are no limits on the vocabulary. In this case, our model learns from all frequently appearing subject, predicate, and object tokens in the captions. Additional matching step is needed at inference time to identify concepts in the target vocabulary. We will explore this setting in our experiment.

Triplet Parsing and Filtering. We use an off-the-shelf rule-based language parser [21] based on Schuster *et al.* [39] to parse the triplets in the image captions. After parsing, the triplets with the lemmatized words for subject, predicate and object are obtained. We further perform an optional filtering step on the initial collection of triplets. For the closed-set setting, we only keep concepts that can be matched to the categories in the target vocabulary. Two concepts are matched if (1) there is overlapping between their synsets, lemmas or hypernyms in WordNet [33] (*e.g.* “tortoise” → “animal”), or (2) if their root forms can be matched (*e.g.* “baseball player” → “player”).

Pseudo Label Assignment. With the filtered triplets, our next step is to match sentence triplets to pairs of regions provided by the object detector. This is done by a greedy matching between every triplet from the caption and each region pair from the image. Specifically, we match the corresponding subject and object tokens between a triplet and a region pair, again using a token’s synsets, the synsets’ lemmas and hypernyms in WordNet [33] and its root form. If multiple triplets are matched to the same region pair, we randomly select one of them. We also filter out region pairs that does not overlap and far away from each other, following [63], as these pairs are less likely to contain a relationship. Once matched, the triplet is considered as the pseudo label of the region pair for training our model.

Model Training. Our model is trained by predicting the pseudo labels of the region pairs. We apply a multi-class cross-entropy loss for the subject, predicate, and object, respectively. Our final loss function is given by

$$\mathcal{L} = \lambda_s \mathcal{L}_s + \lambda_p \mathcal{L}_p + \lambda_o \mathcal{L}_o \quad (4)$$

where \mathcal{L}_s , \mathcal{L}_p , and \mathcal{L}_o is the loss for subject, predicate, and object respectively. And λ_s , λ_p , and λ_o are their corresponding loss weights. We set $\lambda_s = \lambda_o = 0.5$ and $\lambda_p = 1$

following previous work [63, 61].

Weighted Loss. One challenge for learning is the domain gap between (a) image-sentence pairs used for training and (b) images and their target scene graphs during inference. For example, the distributions of concepts might be quite different in image-sentence pairs vs. image scene graphs. In the closed-set setting, we might have an estimated frequency of the concepts on scene graphs. In this case, we apply a weighted loss during training, where the weight for each category is set to the ratio between the frequency of the token in image-sentence pairs and the estimated frequency of the matched tokens in scene graphs. If a category is not matched to any target category, no loss weight will be applied. This weighted loss function only requires an estimated frequency of concepts on the target dataset, and can be considered as a simple approach for domain adaption.

Model Inference. Once trained, our model takes a region pair and its contextual features, and predicts a SPO triplet. To obtain a scene graph, we enumerate all possible region pairs and feed them into our model. The predicted probabilities are further averaged for each region and thus each region is predicted to single category. In the open-set setting, an additional matching step is needed to infer the probability of target categories based on the predicted categories from image-sentence pairs. In this case, we apply the same matching step in our label assignment step.

Extension to Weakly and Fully Supervised Settings. Our model can be easily extended to weakly and fully supervised settings. In weakly supervised setting, we replace triplets parsed from captions with those from unlocalized scene graphs [61], and follow the same label assignment of our setting. For fully supervised setting, we simply replace our pseudo labels with ground-truth scene graph labels.

4. Experiments and Results

We now present our experiments and results. We start with our main results on learning SGG from image-sentence pairs, followed by our ablation studies. Further, we present results on fully supervised SGG, and explore open-set SGG.

Datasets. To evaluate our model, we used the standard split [53] of Visual Genome (VG) [24] (150 objects, 50 predicates, 75K/32K images for train/test). VG comes with human-annotated image captions and localized scene graphs, and is a widely used benchmark for SGG. We also considered image captions on VG for our ablation study, and localized scene graphs on VG for fully supervised SGG.

For training, we considered image captions from VG, COCO Caption (COCO) [6], and Conceptual Caption (CC) [41]. COCO contains 123K images with each labeled by 5 human-annotated captions. We selected 106k images in COCO for training by filtering out images that exist in the test set of VG. CC contains 3.3M image-caption pairs automatically collected from alt-text enabled images on the web.

Method	Training Setting					SGDet	
	Supervision	Level	Source	#Triplets	#Images	R@50	R@100
Ours+Full	Localized Scene Graph	Full	Visual Genome	406K	58K	13.8	15.3
VSPNet [61]	Unlocalized Scene Graph	Weak	Visual Genome	406K	58K	4.7	5.4
VSPNet†						6.7	7.4
Ours+Weak						10.0	11.5
Ours+MotifNet	Image Description	Weaker	CC + COCO	313K	210K	5.6	6.7
Ours						5.9	7.0

Table 2. Results of language supervised SGG. Different from all previous approaches, our model can learn from image-sentence pairs for SGG. With only image-sentence pairs as the supervisory signal, our model outperforms VSPNet — a latest method of weakly supervised SGG trained using human-annotated and unlocalized scene graphs.

Method	Training Setting		SGDet	
	Supervision	Source	R@50	R@100
LSWS[58]	Unlocalized Scene Graph	Visual Genome	7.3	8.7
Ours			10.0	11.5
LSWS[58]	Image Description	Visual Genome	3.9	4.0
Ours		Visual Genome	9.2	10.3
LSWS[58]		COCO	3.3	3.7
Ours		COCO	5.8	6.7

Table 3. Comparison to the concurrent work of LSWS [58].

For the closed-set setting where the target categories are known, we matched the parsed tokens from each dataset to target categories, and kept 148-52, 143-56, 148-64 object-predicate categories for VG, COCO and CC, respectively, leading to 673K/75K (triplets/images) on VG, 154K/64K on COCO, and 159K/145K on CC.

Evaluation Protocol and Metrics. For the majority of our experiments, we evaluate Scene Graph Detection (SGDet) following the protocol from [53]. SGDet captures both the localization and classification performance using metrics of Recall@K (R@K) [30, 53] and mean Recall@K (mR@K) [5, 47]. R@K computes the recall between the top K predicted triplets and ground-truth ones. A predicted triplet is considered as correct only when all requirements are met: (1) the predicted triplet labels match one of the ground-truth triplet, (2) the detected subject-object regions match the ground-truth subject-object regions with an IoU ≥ 0.5 , respectively. mR@K averages R@K across all predicate categories. We also included Scene Graph Classification (SGCls) and Predicate Classification (PredCls) in our experiment on fully SGG. Importantly, all experiments were conducted with graph constraint that limits each subject-object pair to have only one predicate prediction.

Implementation Details. We used a Faster R-CNN [37] detector pre-trained on OpenImages [25], capable of detecting 601 object categories. We kept the top 36 objects per image and extracted the 1536-D region features from the detector. The object tags were represented by the 300-D GloVe embeddings [34]. We adopted the Transformer implementation from UNITER [9] with 2 self-attention layers, 12 attention heads in each layer and hidden size $d = 768$. SGD optimizer was used in training with the image batch of 32, 16 sampled triplets per image, and the initial learn-

ing rate of 0.0032. We used the benchmark implementation provided by Tang *et al.* [46] for evaluation.

4.1. Language Supervised Scene Graph Generation

We now present our main results on learning to generate scene graphs from image-sentence pairs.

Setup and Baselines. We consider several baselines and variants of our model. A key feature of our model is the ability to learn from only image-sentence pairs.

- **VSPNet** [61] is designed for weakly supervised SGG and learns from unlocalized scene graph. As our close competitor, VSPNet takes the inputs of object proposals from the same OpenImage detector used by our model.
- **VSPNet†** further augments VSPNet with object box predictions from the detector. VSPNet† thus has the same input image regions as our model.
- **Ours+Weak** is our model trained using unlocalized scene graphs, same as the setting of VSPNet.
- **Ours+MotifNet** combines our pseudo label assignment with a supervised SGG model (MotifNet [63]). This model is thus trained using only image-sentence pairs.
- **Ours+Full** is our model trained with full supervision and using ground-truth scene graph labels. This should be considered as an upper bound of our model.

Results. Table 2 presents our main results. With image description (CC + COCO) as only supervision, our models (Ours/Ours+MotifNet) significantly outperform VSPNet trained using unlocalized scene graphs (7.0/6.7 vs. 5.4 R@100), despite that image-sentence pairs are much weaker supervisory signals. Our Transformer-based model also beats Ours+MotifNet, and performs on par with the improved version of VSPNet (VSPNet†) (7.0 vs. 7.4 R@100). When trained using unlocalized scene graphs, our model (Ours+Weak) again outperforms VSPNet variants by a large margin (11.5 vs. 5.4/7.4 R@100). These results provide convincing evidence that our model can learn from only image-sentence pairs to detect scene graph in an image with high quality. Finally, there is a noticeable gap between Ours and Ours+Weak (7.0 vs. 11.5 R@100), and between Ours+Weak and Ours+Full (11.5 vs. 15.3 R@100), suggesting ample room for future work.

Fig. 3 further visualizes the output scene graphs from

Source of Image Description			Weighted Loss	#Triplets	#Images	SGDet	
CC	COCO	VG				R@50	R@100
✓				159K	145K	3.4	4.1
	✓			154K	64K	3.8	4.5
✓			✓	159K	145K	5.3	6.4
	✓		✓	154K	64K	5.8	6.7
✓	✓		✓	313K	210K	5.9	7.0
		✓	-	673K	75K	9.2	10.3

Table 4. Ablation study on different sources of image descriptions and weight loss for training our model.

our models, including Ours+Full (left), Ours+Weak (middle) and Ours (right) in Table 2. Our model trained by image-sentence pairs produces scene graphs with a comparable quality as those trained using strong supervision (*e.g.* “man-on-motorcycle” and “man-wearing-helmet” in the 1st row). Further, our models trained using scene graphs tend to predict a different set of predicate when compared to our model trained using image-sentence pairs. This is best illustrated in the 3rd row of Fig. 3 (“on” vs. “has”). We conjecture that this is caused by different distributions of predicates in scene graph and in image captions. Finally, it is worth noting that similar to many previous approaches, our models fall short when common sense reasoning is needed. This is shown in the 4th row of Fig. 3, where our models predict two man wearing the same jacket or shirt.

Comparison to LSWS [58]. In addition, we compare our results to a concurrent work of LSWS [58]. LSWS also learns to generate scene graph from image-sentence pairs using iterative visual grounding. Table 3 summarizes the comparison. When trained with the same level of supervision and the same dataset, our models constantly outperform LSWS by a large margin. For example, when trained using image-sentence pairs on COCO, our method achieves 5.8 R@50 and 6.7 R@100 vs. 3.3 R@50, and 3.7 R@100 from LSWS — a relative gain of at least 75%. When trained with unlocalized scene graph as the setting in VSPNet [61], our model also outperforms LSWS by a noticeable margin (+2.7 R@50 and +2.8 R@100).

4.2. Ablation Studies

We now present ablation studies of our method.

Sources of Image Descriptions. Table 4 presents results of our model trained using different sources of captions. Not surprisingly, the model trained on VG performs better (10.3 R@100) than the model trained on CC (4.1 R@100) or COCO (4.5 R@100), since the evaluation dataset is also VG. Interestingly, the model trained on CC performs on par with the model trained on COCO with similar number of triplets, despite that captions on COCO are manually annotated and has higher quality than those harvested from the Internet on CC. We thus conjecture that the performance of our model is minorly influenced by the caption quality.

Effects of Weighted Loss. Table 4 also compares the use of different loss functions. Adding weighted loss improves

Model	Object Detector	Label Assignment	SGDet	
			R@50	R@100
VSPNet [61]	OpenImages	Iterative Alignment	4.7	5.4
VSPNet†	OpenImages	Iterative Alignment	6.7	7.4
MotifNet	OpenImages	Detection Tags (Ours)	9.3	10.7
Ours	OpenImages	Detection Tags (Ours)	10.0	11.5
Ours	Objects365	Detection Tags (Ours)	6.1	6.4

Table 5. Ablation study on object detectors and label assignment schemes. Results are reported using unlocalized scene graphs as supervision. Our proposed label assignment scheme provides consistent performance boost across methods, while the choice of object detectors has a major impact of the performance.

Model Inputs		Object Detection mAP	SGDet	
Text Input	Visual Input		R@50	R@100
✓	✓	10.7	10.0	11.5
	✓	10.6	3.9	4.7
✓		6.9	6.2	7.7

Table 6. Ablation study on different model inputs. Results are reported using unlocalized scene graphs as supervision. Visual and textual features complement to each other for SGG.

the recall from 4.1 R@100 to 6.4 R@100 when using CC as the training data. This result indicates that using weighted loss can effectively close the domain gap between datasets. For example, the predicate “wearing” appears frequently in VG yet occurs rarely in CC and COCO. With weighted loss, the recall of “wearing” is improved by 22.1 R@100.

Effects of Label Assignment. We evaluate our label assignment scheme in Table 5. Specifically, we consider the weakly supervised setting of learning from unlocalized scene graphs, apply our method to MotifNet [63], and present the results in 3rd row. With our scheme, MotifNet beats the latest VSPNet (10.7 vs. 7.4 R@100), suggesting the effectiveness of our label assignment scheme.

Effects of Object Detector. We additionally consider another object detector trained on Object365 dataset [40]. In Table 5, our model with the Objects365 detector has lower recall (6.4 R@100) than our model with the OpenImages detector (11.5 R@100). Upon a close inspection, we conclude that the recall drop is mainly caused by the mismatch between the object categories of the detector and those in VG. Particularly, we find only 94 (out of 150) VG objects can be matched to Objects365 categories while 123 VG objects can be matched to OpenImages categories. For example, “shirt” and “building” are most frequent objects in VG. Objects365 detector cannot detect them while OpenImages detector can. As a result, the triplets involving these objects will not be used to train the model, and the trained model fails to detect these concepts.

Textual vs. Visual Inputs. Finally, We study the contribution of model inputs. This is done by probing our trained model during inference and masking out one input at a time. For textual input, we substitute subject and object embeddings for a random vector. For visual input, we replace the original region features with the averaged region features in

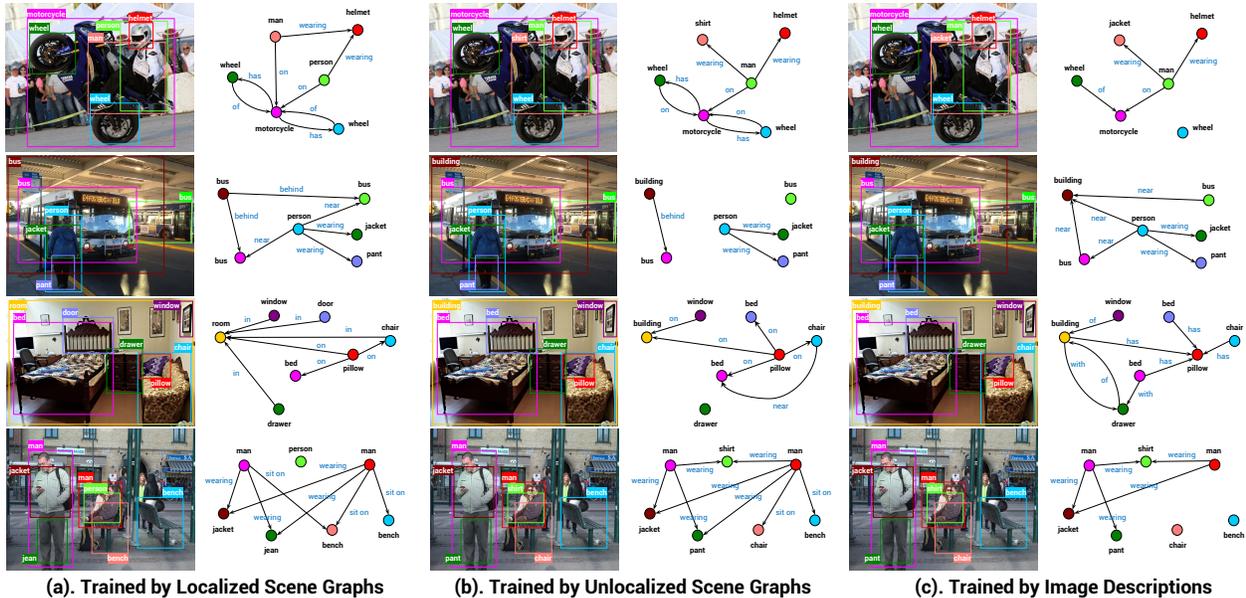


Figure 3. Qualitative results of our models on VG test set for SGG. All models take the same detected regions and predict the scene graph labels. In each row, we show 3 identical images and the corresponding scene graphs generated from the models trained by different levels of supervision. The visualized relationships are picked from the top 30 predicted triplets.

Model	Recall									Mean Recall								
	SGDet			SGCls			PredCls			SGDet			SGCls			PredCls		
	@20	@50	@100	@20	@50	@100	@20	@50	@100	@20	@50	@100	@20	@50	@100	@20	@50	@100
IMP [53]	18.1	25.9	31.2	34.0	37.5	38.5	54.3	61.1	63.1	2.8	4.2	5.3	5.2	6.2	6.5	8.9	11.0	11.8
VTransE [64]	23.0	29.7	34.3	35.4	38.6	39.4	59.0	65.7	67.6	3.7	5.0	6.0	6.7	8.2	8.7	11.6	14.7	15.8
VCtree [47]	24.7	31.5	36.2	37.0	40.5	41.4	59.8	66.2	68.1	4.2	5.7	6.9	6.2	7.5	7.9	11.7	14.9	16.1
MotifNet [63]	25.1	32.1	36.9	35.8	39.1	39.9	59.5	66.0	67.9	4.1	5.5	6.8	6.5	8.0	8.5	11.5	14.6	15.8
Ours	24.6	31.8	36.3	36.5	40.0	40.8	58.7	65.6	67.4	5.3	7.3	8.7	8.3	10.4	11.1	13.3	17.7	19.5

Table 7. Results of fully supervised SGG. All models use the same object detector pre-trained on the VG dataset, and the same codebase provided by Tang et al. [46] for evaluation. Results of previous models come from Tang et al. [46].

the current image. The results are presented in Table 6. Using only visual input leads to a minor drop in detection mAP (10.6 vs. 10.7) and a large drop in scene graph recall (4.7 vs. 11.5 R@100), indicating a major performance drop in predicate prediction. In contrast, using only text input has a large drop mAP (6.9 vs. 10.7) and a moderate drop in scene graph recall (7.7 vs. 11.5 R@100). These results suggest that visual and textual inputs compliment to each other — predicate prediction primarily relies on textual inputs while object prediction mainly depends on visual inputs.

4.3. Fully-supervised Scene Graph Generation

We further evaluate our model for fully supervised SGG. **Setup and Baselines.** To demonstrate the strength of our Transformer-based model, we present results on fully supervised SGG and compare to several latest methods [53, 64, 47, 63], following the standard protocol of training and testing on VG. Note that all models used the same object detector trained on VG and the same benchmark implementation provided by Tang *et al.* [46].

Results. We report recall and mean recall for SGDet, SG-

Cls, and PredCls in Table 7. The recalls of our model compare favorably to previous best results (SGDet: 36.3 vs. 36.9 R@100, SGCls: 40.8 vs. 41.4 R@100, PredCls: 67.4 vs. 68.1). More importantly, the mean recalls of our model are significantly higher than previous models across all evaluation protocols (SGDet: 8.7 vs. 6.9 R@100, SGCls: 11.1 vs. 8.7 R@100, PredCls: 19.5 vs. 16.1). Compared to recall, mean recall [5, 47] better characterizes the performance on categories with fewer samples. These results suggest that our model is better at capturing those tail categories.

4.4. Open-set Scene Graph Generation

Moving forward, we consider a challenging open-set setting for SGG, where the categories of target concepts (objects and predicates) are unknown during training. We believe this is the first result for open-set SGG.

Setup. In this experiment, our model is trained on COCO Caption and evaluated on VG. During training, we parsed concept categories from captions, remove the low-frequency categories, and formed a vocabulary of 4273 objects and 677 predicates. This vocabulary was then used to

Model	#Objects	#Predicates	#Triplets	#Images	SGDet	
					R@50	R@100
Ours	143	56	154K	64K	3.8	4.5
Ours	4273	677	758K	105K	4.1	4.8

Table 8. Results of open-set SGG. Evaluation is performed on VG with the vocabulary and model learned from COCO.

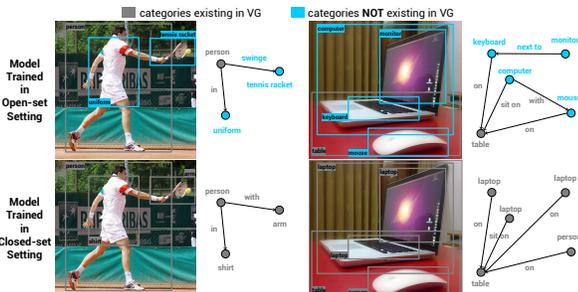


Figure 4. Qualitative results of our models (trained in open-set and closed-set settings) on VG test set for SGG.

train our model. At inference time, we first generated scene graphs using our vocabulary, and then matched the detected categories in our vocabulary to target concepts on VG (150 objects and 50 predicates) for evaluation.

Results. Table 8 compares the results of our models trained in closed-set and open-set settings using the same COCO caption dataset. The model trained in open-set setting has slightly better recall (4.8 vs. 4.5 R@100). Our open-set results are also comparable to VSPNet (supervised by unlocalized scene graphs on VG in a closed-set setting). We hypothesize that the open-set setting allows the model to learn from more concepts and thus benefits SGG. To verify this hypothesis, we plot the output scene graph from our models trained on closed-set and open-set settings in Fig. 4. Compared to our closed-set model, our open-set model detects more concepts outside VG (e.g. “swinge”, “mouse”, “keyboard”). Our results suggest an exciting avenue of large-scale training of open-set SGG using image captioning dataset such as CC.

5. Conclusion

We proposed one of the first methods of learning to generate scene graphs from image-sentence pairs. Our key idea is to use off-the-shelf object detectors, so as to match detected object tags to parsed tokens from captions, thus creating “pseudo” labels for training. Further, we designed a Transformer-based model and demonstrated strong results across different levels of supervision. Our model learned from only image-sentence pairs, outperformed a state-of-the-art weakly supervised model trained by human-annotated unlocalized scene graphs. More importantly, we presented the first result for open-set scene graph generation. We hope our work points to exciting avenues of learning structured visual representation from

natural language.

Limitation and Future Work. A main limitation of our method is the dependency on an object detector covering a wide range of concepts. We anticipate that our method will benefit from the development of open-vocabulary detectors [59, 62]. Moreover, there is fundamental ambiguity during the label assignment step in learning from image-sentence pairs when there are multiple object instances of the same category. We hypothesize that contextual cues such as surrounding objects might help and will leave this as future work. Finally, we also plan to explore using scene graphs learned from image-sentence pairs for vision and language tasks (e.g. VQA).

Acknowledgement: YZ and YL acknowledge the support provided by the UW-Madison OVCGRG with funding from WARF. JS and CX were supported by the National Science Foundation (NSF) under Grant RI:1813709. The article solely reflects the opinions and conclusions of its authors but not the funding agency.

References

- [1] Iro Armeni, Zhi-Yang He, JunYoung Gwak, Amir R Zamir, Martin Fischer, Jitendra Malik, and Silvio Savarese. 3D scene graph: A structure for unified semantics, 3d space, and camera. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5664–5673, 2019. 1
- [2] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. 4, 16
- [3] Federico Baldassarre, Kevin Smith, Josephine Sullivan, and Hossein Azizpour. Explanation-based weakly-supervised learning of visual relations with graph networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 612–630. Springer, 2020. 2
- [4] Long Chen, Hanwang Zhang, Jun Xiao, Xiangnan He, Shiliang Pu, and Shih-Fu Chang. Counterfactual critic multi-agent training for scene graph generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4613–4623, 2019. 1, 2
- [5] 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 (CVPR)*, pages 6163–6171, 2019. 6, 8
- [6] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO captions:

- Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015. 2, 5, 13
- [7] Xinlei Chen and Abhinav Gupta. Webly supervised learning of convolutional networks. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 1431–1439, 2015. 1, 2
- [8] Xinlei Chen, Abhinav Shrivastava, and Abhinav Gupta. NEIL: Extracting visual knowledge from web data. In *Proceedings of the IEEE international conference on computer vision (ICCV)*, pages 1409–1416, 2013. 1, 2
- [9] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. UNITER: Universal image-text representation learning. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 104–120. Springer, 2020. 2, 3, 6
- [10] 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 (CVPR)*, pages 3076–3086, 2017. 2
- [11] Karan Desai and Justin Johnson. VirTex: Learning visual representations from textual annotations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 1, 2
- [12] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. 4
- [13] Helisa Dhama, Azade Farshad, Iro Laina, Nassir Navab, Gregory D. Hager, Federico Tombari, and Christian Rupprecht. Semantic image manipulation using scene graphs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 1
- [14] Santosh K Divvala, Ali Farhadi, and Carlos Guestrin. Learning everything about anything: Webly-supervised visual concept learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3270–3277, 2014. 1, 2
- [15] Ali Farhadi, Ian Endres, Derek Hoiem, and David Forsyth. Describing objects by their attributes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1778–1785. IEEE, 2009. 2
- [16] Tanmay Gupta, Arash Vahdat, Gal Chechik, Xiaodong Yang, Jan Kautz, and Derek Hoiem. Contrastive learning for weakly supervised phrase grounding. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2019. 1
- [17] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016. 16
- [18] Drew Hudson and Christopher D Manning. Learning by abstraction: The neural state machine. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 32. Curran Associates, Inc., 2019. 1
- [19] Achiya Jerbi, Roei Herzig, Jonathan Berant, Gal Chechik, and Amir Globerson. Learning object detection from captions via textual scene attributes. *arXiv preprint arXiv:2009.14558*, 2020. 2
- [20] Jingwei Ji, Ranjay Krishna, Li Fei-Fei, and Juan Carlos Niebles. Action Genome: Actions as compositions of spatio-temporal scene graphs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10236–10247, 2020. 1
- [21] Mao Jiayuan and Kasai Seito. Scene graph parser. <https://github.com/vacancy/SceneGraphParser>, 2018. 5
- [22] 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 (CVPR)*, June 2018. 1
- [23] Armand Joulin, Laurens Van Der Maaten, Allan Jabri, and Nicolas Vasilache. Learning visual features from large weakly supervised data. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 67–84. Springer, 2016. 1, 2
- [24] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yanis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual Genome: Connecting language and vision using crowdsourced dense image annotations. *International Journal of Computer Vision (IJCV)*, 123(1):32–73, 2017. 2, 5, 13
- [25] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Alexander Kolesnikov, et al. The open images dataset v4. *International Journal of Computer Vision (IJCV)*, pages 1–26, 2020. 6
- [26] Christoph H Lampert, Hannes Nickisch, and Stefan Harmeling. Learning to detect unseen object classes by between-class attribute transfer. In *Proceedings of*

- the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 951–958. IEEE, 2009. 2
- [27] Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 121–137. Springer, 2020. 2, 3
- [28] 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. 1, 2
- [29] 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 (ICCV)*, pages 1261–1270, 2017. 1, 2
- [30] Cewu Lu, Ranjay Krishna, Michael Bernstein, and Li Fei-Fei. Visual relationship detection with language priors. In *Proceedings of the European conference on computer vision (ECCV)*, pages 852–869. Springer, 2016. 2, 6
- [31] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. ViLBERT: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 13–23, 2019. 2, 3
- [32] Jiasen Lu, Vedanuj Goswami, Marcus Rohrbach, Devi Parikh, and Stefan Lee. 12-in-1: Multi-task vision and language representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 2, 3
- [33] George A Miller. WordNet: a lexical database for english. *Communications of the ACM*, 38(11):39–41, 1995. 5
- [34] Jeffrey Pennington, Richard Socher, and Christopher D Manning. GloVe: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, 2014. 6
- [35] Julia Peyre, Josef Sivic, Ivan Laptev, and Cordelia Schmid. Weakly-supervised learning of visual relations. In *Proceedings of the IEEE international conference on computer vision (ICCV)*, pages 5179–5188, 2017. 2
- [36] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning (ICML)*, 2021. 1, 2
- [37] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems (NeurIPS)*, volume 28. Curran Associates, Inc., 2015. 2, 6
- [38] Mert Bulent Sariyildiz, Julien Perez, and Diane Larlus. Learning visual representations with caption annotations. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020. 1, 2
- [39] Sebastian Schuster, Ranjay Krishna, Angel Chang, Li Fei-Fei, and Christopher D. Manning. Generating semantically precise scene graphs from textual descriptions for improved image retrieval. In *Proceedings of the Fourth Workshop on Vision and Language*, pages 70–80, Lisbon, Portugal, Sept. 2015. Association for Computational Linguistics (ACL). 5
- [40] Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. Objects365: A large-scale, high-quality dataset for object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019. 7
- [41] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 2556–2565, 2018. 1, 2, 5, 13
- [42] Jiaxin Shi, Hanwang Zhang, and Juanzi Li. Explainable and explicit visual reasoning over scene graphs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8376–8384, 2019. 1
- [43] Jing Shi, Yiwu Zhong, Ning Xu, Yin Li, and Chenliang Xu. A simple baseline for weakly-supervised scene graph generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021. 2, 3
- [44] Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. VL-BERT: Pre-training of generic visual-linguistic representations. In *International Conference on Learning Representations (ICLR)*, 2020. 2, 3
- [45] Hao Tan and Mohit Bansal. LXMERT: Learning cross-modality encoder representations from transformers. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2019. 2, 3

- [46] Kaihua Tang, Yulei Niu, Jianqiang Huang, Jiaxin Shi, and Hanwang Zhang. Unbiased scene graph generation from biased training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3716–3725, 2020. [1](#), [2](#), [6](#), [8](#)
- [47] 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 (CVPR)*, pages 6619–6628, 2019. [1](#), [2](#), [3](#), [6](#), [8](#)
- [48] 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 (CVPR)*, pages 1–9, 2017. [1](#)
- [49] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 30, 2017. [4](#), [13](#), [16](#)
- [50] Johanna Wald, Helisa Dharmo, Nassir Navab, and Federico Tombari. Learning 3D semantic scene graphs from 3d indoor reconstructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. [1](#)
- [51] Josiah Wang, Katja Markert, and Mark Everingham. Learning models for object recognition from natural language descriptions. In *The British Machine Vision Conference (BMVC)*, volume 1, page 2, 2009. [1](#), [2](#)
- [52] Wenbin Wang, Ruiping Wang, Shiguang Shan, and Xilin Chen. Sketching image gist: Human-mimetic hierarchical scene graph generation. In *European conference on computer vision (ECCV)*. Springer, 2020. [2](#)
- [53] Danfei Xu, Yuke Zhu, Christopher B Choy, and Li Fei-Fei. Scene graph generation by iterative message passing. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 5410–5419, 2017. [1](#), [2](#), [3](#), [5](#), [6](#), [8](#)
- [54] 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](#)
- [55] Jianwei Yang, Jiasen Lu, Stefan Lee, Dhruv Batra, and Devi Parikh. Visual curiosity: Learning to ask questions to learn visual recognition. In *Conference on Robot Learning (CoRL)*, 2018. [1](#), [2](#)
- [56] Xu Yang, Kaihua Tang, Hanwang Zhang, and Jianfei Cai. Auto-encoding scene graphs for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10685–10694, 2019. [1](#)
- [57] Ting Yao, Yingwei Pan, Yehao Li, and Tao Mei. Exploring visual relationship for image captioning. In *Proceedings of the European conference on computer vision (ECCV)*, pages 684–699, 2018. [1](#)
- [58] Keren Ye and Adriana Kovashka. Linguistic structures as weak supervision for visual scene graph generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8289–8299, 2021. [2](#), [3](#), [6](#), [7](#)
- [59] Keren Ye, Mingda Zhang, Adriana Kovashka, Wei Li, Danfeng Qin, and Jesse Berent. Cap2Det: Learning to amplify weak caption supervision for object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9686–9695, 2019. [1](#), [2](#), [9](#)
- [60] Ruichi Yu, Ang Li, Vlad I Morariu, and Larry S Davis. Visual relationship detection with internal and external linguistic knowledge distillation. In *Proceedings of the IEEE international conference on computer vision (ICCV)*, pages 1974–1982, 2017. [2](#)
- [61] Alireza Zareian, Svebor Karaman, and Shih-Fu Chang. Weakly supervised visual semantic parsing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3736–3745, 2020. [1](#), [2](#), [3](#), [5](#), [6](#), [7](#)
- [62] Alireza Zareian, Kevin Dela Rosa, Derek Hao Hu, and Shih-Fu Chang. Open-vocabulary object detection using captions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. [1](#), [2](#), [9](#)
- [63] 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 (CVPR)*, pages 5831–5840, 2018. [1](#), [2](#), [3](#), [5](#), [6](#), [7](#), [8](#)
- [64] Hanwang Zhang, Zawlin Kyaw, Shih-Fu Chang, and Tat-Seng Chua. Visual translation embedding network for visual relation detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 5532–5540, 2017. [2](#), [3](#), [8](#)
- [65] Hanwang Zhang, Zawlin Kyaw, Jinyang Yu, and Shih-Fu Chang. PPR-FCN: Weakly supervised visual relation detection via parallel pairwise R-FCN. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 4233–4241, 2017. [2](#), [3](#)
- [66] Ji Zhang, Kevin J Shih, Ahmed Elgammal, Andrew Tao, and Bryan Catanzaro. Graphical contrastive losses for scene graph parsing. In *Proceedings of the*

IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 11535–11543, 2019. [1](#), [2](#)

- [67] Yiwu Zhong, Liwei Wang, Jianshu Chen, Dong Yu, and Yin Li. Comprehensive image captioning via scene graph decomposition. In *European Conference on Computer Vision (ECCV)*, pages 211–229. Springer, 2020. [1](#)
- [68] Luowei Zhou, Hamid Palangi, Lei Zhang, Houdong Hu, Jason Corso, and Jianfeng Gao. Unified vision-language pre-training for image captioning and vqa. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 34, pages 13041–13049, 2020. [2](#), [3](#)

Appendices

In appendices, we provide additional implementation details, describe the domain gap among different datasets, and present detailed results of scene graph generation. We hope that this document will complement our main paper.

A. Implementation Details

We now present implementation details on triplet filtering (Sec. 3.2) and our model architecture (Sec. 4).

Triplet Filtering. We describe details for triplet filtering during an open-set setting. In the open-set setting, the subject, predicate and object categories in the target dataset are unknown during training. There is thus no limits on these categories. We simply removed the parsed concepts with low frequency. The frequency threshold was set to 3 and 10 for objects/subjects and predicates, respectively. At the end, we kept 4273-677 subject/object-predicate categories for COCO Caption dataset [6]. These categories correspond to 758K triplet instances and 105K images. These images and triplet instances were used for model training in the open-set setting.

Model Architecture. Table 10 lists the architecture of our Triplet Transformer. Our model consists of the visual embedder, the textual embedder, the multi-layer Transformer encoder, and the classification heads. In the Transformer encoder, each self-attention layer [49] takes the token embeddings as inputs, and outputs a contextualized embedding for each token. We used an input/output dimension of 768-D with 12 attention heads each with 64 dimensions.

B. Domain Gap among Different Datasets

Further, we study the domain gap among the datasets considered in the paper (Conceptual Caption [41], COCO Caption [6], Visual Genome [24]). Our main results are obtained by training our model on Conceptual Caption and COCO, and evaluating the trained model on Visual Genome. Our observation is that the gap between Conceptual Caption and COCO Caption vs. Visual Genome is large and thus our task is very challenging.

To verify our observation, we plot the distribution of the most frequent predicate categories (a) and noun (subject/object) categories (b) across three datasets in Figure 5. These categories are parsed from image captions and matched to a target dictionary defined on Visual Genome. The top 15 most frequent noun/predicate categories are displayed. The domain gap between different datasets can be clearly observed in Figure 5. For example, the predicate category “wearing” is very common in Visual Genome dataset, yet quite rare in COCO Caption and Conceptual Caption. As described in Sec. 3.2, we applied the weighted loss dur-

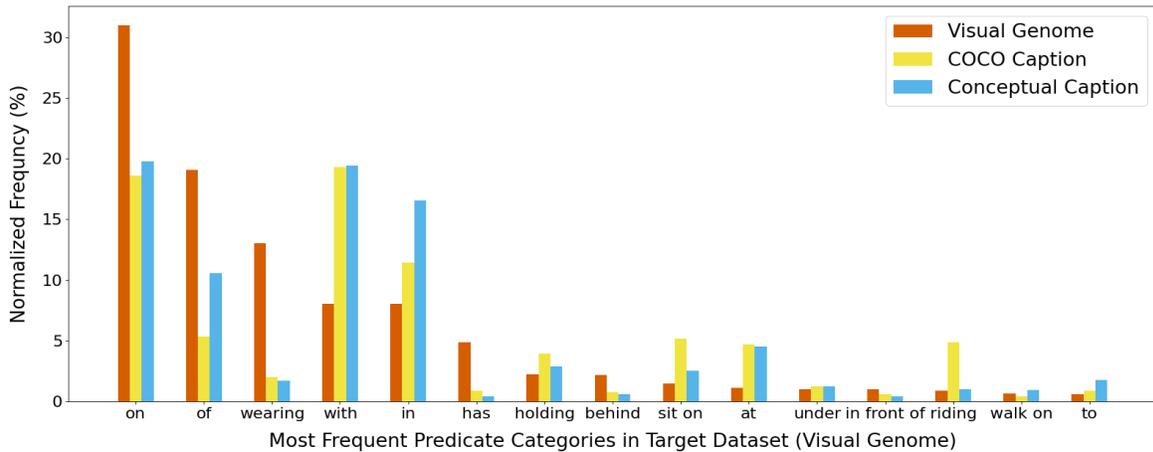
ing training to bridge the domain gap. According to our experiments in Sec. 4.2, the weighted loss can help to improve the performance of scene graph generation.

C. Further Analysis of Language Supervised Scene Graph Generation

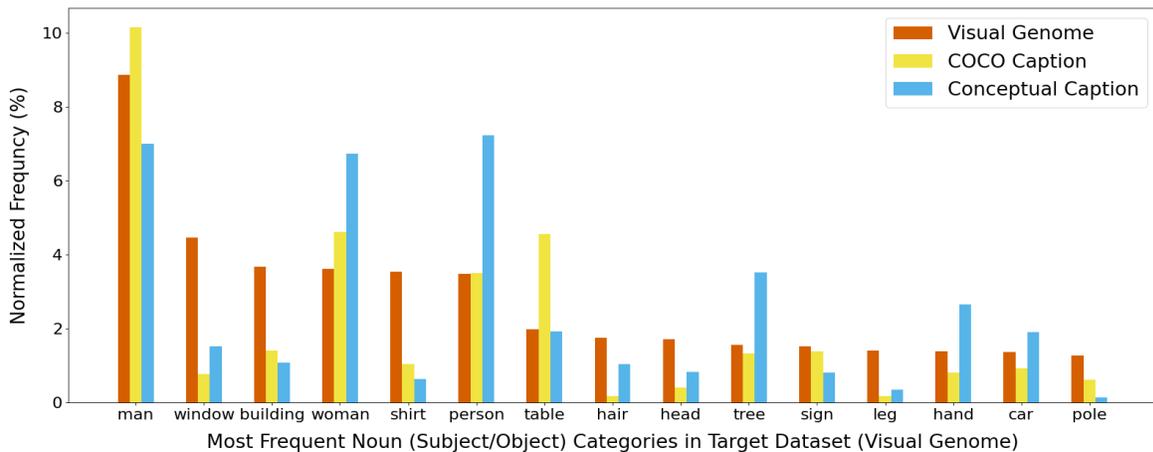
We present additional details for our main results on language supervised scene graph generation, shown in Table 2 of our main paper. Specifically, we provide a breakdown of our results, and compare our models obtained by different levels of supervision.

Setup. We present additional details in the setting of SGM@100, shown in Table 2 of our main paper. We show the per-category recalls for 15 most common predicate categories. These categories are selected by their frequency in the scene graph annotation of Visual Genome. The per-category recalls are calculated in the same way as the mean Recall (mR) except that we present the individual recall before averaging over all categories.

Results. Table 9 shows the per-category recalls for our models trained by different supervisory signals, including language supervised, weakly supervised and fully supervised. Each row in Table 9 corresponds to one of the models in Table 2 of our main paper (*e.g.*, CC+COCO represents our model trained by image descriptions from Conceptual Caption and COCO Caption). The per-category recalls vary drastically among different supervision settings. For example, our model trained by the localized scene graph labels (fully supervised) achieves the highest recall (16.8) for the category “on”, while the model trained by image captions (language supervised) obtains the lowest recall (7.4). On the other hand, our language supervised model largely outperforms fully supervised and weakly supervised models on several categories, such as “with” (7.4 vs. 0.9 and 1.0) and “eating” (17.9 vs. 10.6 and 8.9). We conjecture that the performance difference is again produced by the domain gap of the supervisory signals (*i.e.* datasets). For example, “eating” is more frequent on Conceptual Caption and COCO Caption than on Visual Genome. It will be an exciting direction to investigate how to further bridge this domain gap.



(a). Predicate category distribution across different datasets.



(b). Noun (object/subject) category distribution across different datasets.

Figure 5. The distribution of (a) the predicate categories and (b) the noun (subject/object) categories parsed from different image caption datasets (Visual Genome, COCO Caption and Conceptual Caption). Y-axis represents the category frequency normalized in respective dataset. X-axis indexes the most common categories in the target dataset (Visual Genome). These three datasets have very different distributions with a large domain gap.

Models	Supervision	R@100	on	has	in	of	wearing	near	with	above	holding	behind	sit on	in front of	at	riding	eating
Ours	CC+COCO	7.0	7.4	1.0	3.8	6.2	24.4	8.4	7.4	0.4	5.3	0.4	3.5	2.0	8.4	3.7	17.9
Ours	VG-Weak	11.5	12.3	14.5	4.5	6.9	27.9	6.8	0.9	0.7	19.1	12.7	7.9	3.5	18.9	18.4	10.6
Ours	VG-Full	15.3	16.8	15.8	7.6	7.5	38.4	10.5	1.0	1.4	23.2	13.6	6.8	2.6	9.5	22.8	8.9

Table 9. Per-category recall@100 for scene graph generation. Each row corresponds to one of the models in Table 2 of our main paper. CC+COCO represents the model trained by the image captions from Conceptual Caption (CC) and COCO Caption (COCO). VG-Weak represents the model trained by the unlocalized scene graphs in Visual Genome. VG-Full denotes the model trained by the localized scene graphs in Visual Genome. We list the overall recall as well as the recalls for most frequent predicates.

ID	Module	Input ID	Layer Type	Weights Size	Output Size	Comments
1	Input	-	Region Features	-	(N, 1536)	Obtained from the off-the-shelf object detector
2		-	Region Boxes	-	(N, 7)	
3		-	Region Tags	-	(N)	
4	Visual Embedder	1	FC	(1536, 768)	(N, 768)	
5		4	LN	(768)	(N, 768)	
6		2	FC	(7, 768)	(N, 768)	
7		6	LN	(768)	(N, 768)	
8		-	Type Embedding	(3, 768)	(N, 768)	Region type: subject vs. object vs. context
9		5+7+8	LN	(768)	(N, 768)	
10		9	Dropout	-	(N, 768)	p=0.1
11	Textual Embedder	3	Word Embedding	(604, 200)	(4, 200)	Text tokens: [subject tag, MASK, object tag, SEP]
12		11	FC	(200, 768)	(4, 768)	
13		-	Position Embedding	(4, 768)	(4, 768)	Embedding for text token position
14		12+13	LN	(768)	(4, 768)	
15		14	Dropout	-	(4, 768)	p=0.1
16	Transformer Encoder Layer 1	[10, 15]	Concatenation	-	(N+4, 768)	Concatenate the visual and textual features
17		16	Self-attention	(768, 64, 3, 12)	(N+4, 768)	Multi-head self-attention [49]
18		17	FC	(768, 768)	(N+4, 768)	
19		18	Dropout	-	(N+4, 768)	p=0.1
20		16+19	LN	(768)	(N+4, 768)	Residual connection followed by LN
21		20	FC	(768, 3072)	(N+4, 3072)	
22		21	GELU	-	(N+4, 3072)	
23		22	FC	(3072, 768)	(N+4, 768)	
24		23	Dropout	-	(N+4, 768)	p=0.1
25		20+24	LN	(768)	(N+4, 768)	Residual connection followed by LN
26	Transformer Encoder Layer 2	25	Self-attention	(768, 64, 3, 12)	(N+4, 768)	Multi-head self-attention [49]
27		26	FC	(768, 768)	(N+4, 768)	
28		27	Dropout	-	(N+4, 768)	p=0.1
29		25+28	LN	(768)	(N+4, 768)	Residual connection followed by LN
30		29	FC	(768, 3072)	(N+4, 3072)	
31		30	GELU	-	(N+4, 3072)	
32		31	FC	(3072, 768)	(N+4, 768)	
33		32	Dropout	-	(N+4, 768)	p=0.1
34		29+33	LN	(768)	(N+4, 768)	Residual connection followed by LN
35	Classification Heads	34	Indexing	-	(1, 768)	Subject visual embedding
36		34	Indexing	-	(1, 768)	Object visual embedding
37		34	Indexing	-	(1, 768)	Subject textual embedding
38		34	Indexing	-	(1, 768)	Object textual embedding
39		34	Indexing	-	(1, 768)	Predicate embedding
40		11	Indexing	-	(1, 200)	Subject word embedding
41		11	Indexing	-	(1, 200)	Object word embedding
42		40	FC	(200, 768)	(1, 768)	
43		41			(1, 768)	
44		35+42	FC, ReLU, FC	(768,768), -, (768,151)	(1, 151)	Subject logits
45		36+43			(1, 151)	Object logits
46		44	Softmax	-	(1, 151)	Used for subject cross-entropy loss
47		45	Softmax	-	(1, 151)	Used for object cross-entropy loss
48		35	FC	(768, 768)	(1, 768)	
49		36	FC	(768, 768)	(1, 768)	
50		37	FC	(768, 768)	(1, 768)	
51		38	FC	(768, 768)	(1, 768)	
52		39+48+49+50+51	FC, ReLU, FC	(768,768), -, (768,151)	(1, 51)	Predicate logits
53		52	Softmax	-	(1, 51)	Used for predicate cross-entropy loss

Table 10. Model architecture of our proposed Triplet Transformer. For each input image, $N = 36$ detected regions from the off-the-shelf object detector are used as our model inputs. Within the Input ID column, “+” means summing up the input features, “[]” means concatenating the input features. Within the Layer Type column, “FC” denotes the fully-connected layer, “LN” denotes the Layer Normalization [2], “GELU” represents the GELU activation function [17], and “Indexing” means extracting the features corresponding to a particular slot (e.g., visual or textual embeddings of subject, predicate and object.).