

From Pixels to Graphs: Open-Vocabulary Scene Graph Generation with Vision-Language Models

Rongjie Li^{1*} Songyang Zhang^{2*} Dahua Lin² Kai Chen^{2‡} Xuming He^{1,3‡}

¹School of Information Science and Technology, ShanghaiTech University

²Shanghai AI Laboratory

³Shanghai Engineering Research Center of Intelligent Vision and Imaging

{lirj2, hexm}@shanghaitech.edu.cn, {zhangsongyang, lindahua, chen kai}@pjlab.org.cn,

Abstract

Scene graph generation (SGG) aims to parse a visual scene into an intermediate graph representation for downstream reasoning tasks. Despite recent advancements, existing methods struggle to generate scene graphs with novel visual relation concepts. To address this challenge, we introduce a new open-vocabulary SGG framework based on sequence generation. Our framework leverages vision-language pre-trained models (VLM) by incorporating an image-to-graph generation paradigm. Specifically, we generate scene graph sequences via image-to-text generation with VLM and then construct scene graphs from these sequences. By doing so, we harness the strong capabilities of VLM for open-vocabulary SGG and seamlessly integrate explicit relational modeling for enhancing the VL tasks. Experimental results demonstrate that our design not only achieves superior performance with an open vocabulary but also enhances downstream vision-language task performance through explicit relation modeling knowledge.

1. Introduction

The main objective of scene graph generation (SGG) is to parse an image into a graph representation that describes visual scenes in terms of object entities and their relationship. Such a generated scene graph can serve as a structural and interpretable representation of visual scenes, facilitating connections between visual perception and reasoning [37]. In particular, it has been widely used in various vision-language (VL) tasks, including visual question answering [8, 10, 11, 30, 33], image captioning [42, 43], referring expressions [41] and image retrieval [12].

Most previous works have focused on addressing the

*Equal contribution and this work was partially done when Li was a research intern at the Shanghai AI Laboratory. ‡ denotes corresponding author. Code is available: https://github.com/SHTUPLUS/Pix2Grp_CVPR2024

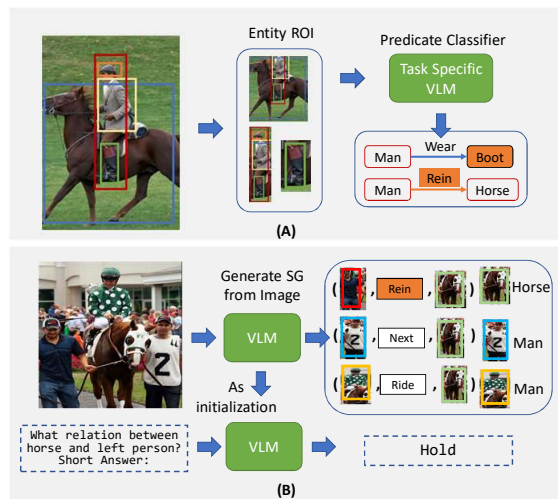


Figure 1. **An illustration of open-vocabulary SGG paradigm comparison.** (A) previous work adopt the task-specific VLM as predicate classifiers from given entity proposals; (B) Our framework offers a unified framework for generating scene graph with novel predicates from images directly and conducting VL tasks.

SGG problems in a close-world setting. While much effort has been made in tackling long-tailed data bias [6, 32, 53] and insufficient labeling [3, 19, 51], the typical label space of those methods only covers a limited subset of the diverse visual relationships in the real world. This can result in incomplete scene representations and also domain gaps when used as intermediate representations in downstream vision-language tasks [38]. To address this limitation, recent works [7, 48, 52] start to tackle the SGG problem under various open-vocabulary settings by exploiting the image-text matching capability of pre-trained vision-language models (VLM). Nevertheless, these attempts typically focus on a simplified setting of the open-vocabulary SGG task, such as allowing only novel entities [7, 48] (Ov-SGG), or a subtask of SGG, such as classifying open-set predicates with given entity pairs [7, 52], as shown in part

A of Fig. 1. It remains open to building an end-to-end SGG model in a general open-vocabulary setting. Moreover, those methods often employ an additional pre-training step to enhance the representation power of VLM on relation modeling, which induces a high training cost for large-scale models.

In this work, we aim to address the open-vocabulary SGG problem in a more general setting, i.e., generating scene graphs with both known and novel visual relation triplets from pixels. To this end, we propose an efficient end-to-end framework that leverages the image-to-text generation paradigm of pre-trained VLMs, as illustrated in part B of Fig. 1. This framework, dubbed Pixels to Scene Graph Generation with Generative VLM (PGSG), in particular, we formulate the SGG as an image-to-sequence problem and introduce a fine-tuning strategy based on generative VLMs. In contrast to previous methods relying on image-to-text matching VLMs, our generative framework provides a more efficient way to utilize the rich visual-linguistic knowledge of pre-trained VLMs for relation-aware representation without requiring additional pre-training of VLMs. In addition, by converting SGG into a sequence generation problem, our method unifies the SGG task with a diverse set of VL tasks under the generative framework, which enables us to transfer visual relation knowledge to other VL tasks in a seamless manner (e.g., via model initialization).

Specifically, we develop an image-to-graph generation framework consisting of three main components. First, we introduce *scene graph prompts*, which transform scene graphs into a sequence representation with *relation-aware tokens*. Given those graph prompts, we then employ a pre-trained VLM to generate a corresponding scene graph sequence for each input image. Finally, we design a plug-and-play *relationship construction module* to extract locations and categories of relation triplets, which produces the output scene graph. Our generation pipeline only requires fine-tuning on target SGG datasets. Such a framework enables us to generate scene graphs with diverse predicate concepts thanks to the large capacity of pre-trained VLMs. Moreover, our unified framework supports a seamless adaptation strategy, employing our fine-tuned VLM as an initialization for image-to-text generative models in various VL tasks.

We validate our framework on three SGG benchmarks: Panoptic Scene Graph [40], OpenImages-V6 [16], and Visual Genome [15], achieving state-of-the-art performance in the general open-vocabulary setting. Furthermore, we apply our SGG-based VLM to multiple VL tasks and obtain consistent improvements, highlighting our effective relational knowledge transfer paradigm.

In summary, the contribution of our work is threefold:

- We propose a novel framework for a general open-vocabulary SGG problem based on the image-to-text generative VLM, which can be seamlessly integrated with

other VL tasks.

- Our method introduces a scene graph prompt and a plug-and-play relation-aware converter module on top of generative VLM, allowing more efficient model learning.
- Our framework achieves superior performance on the general open-vocabulary SGG benchmarks and consistent improvements on downstream VL tasks.

2. Related Work

Scene Graph Generation Scene graph generation (SGG) aims to localize and classify entities and visualize the relations between them in images. The concept of SGG was first introduced in [15]. Initially, it was used as an auxiliary representation to enhance image retrieval [12, 29]. The diverse semantic concepts and compositional structures of visual relations create a large concept space, making this task difficult. Researchers have approached SGG from various perspectives, including addressing intrinsic long-tailed data bias [2, 18, 20, 22, 32], incorporating scene context to model diverse visual concepts [24, 31, 49, 50], and reducing labeling costs through weakly or semi-supervised training [19, 45, 52, 54]. However, most SGG work focuses on limited categories, while few address SGG models' generalization capabilities. Recent studies, such as [7], have systematically investigated this problem by designing VL-pretrained models to classify predicates with unseen object categories, and [52] have investigated visual relations with unseen objects. We generate scene graphs with unseen predicate categories to explore a more realistic setting.

Scene Graph for VL Tasks Scene graphs have been extensively studied for improving vision-language (VL) tasks. Early works used scene graphs to improve image retrieval. Several studies showed that scene graphs can be used for image captioning, visual question answering (VQA), and visual grounding. Adding explicit scene graphs to VL models for downstream tasks is difficult. As discussed in [27, 35], generated scene graph noise can harm VL models. Second, SGG models trained on public datasets with few classes often have a large semantic domain gap compared to visual concepts needed for downstream VL tasks [38]. Some works incorporate scene graph representations within model training to address these issues [38, 46]. Our work aligns with this line: we formulate a unified framework that allows joint optimization between explicit scene graph generation modeling and VL tasks.

Open-vocabulary Scene Graph Generation Previous research has predominantly focused on enhancing the generalization of new entity combinations in visual relationships [7, 13, 15, 32, 52]. Notably, He et al. [7] and Yu et al. [48] have recently addressed open-vocabulary predicate classification through visual-language pre-training. Nevertheless, these methods struggle to effectively detect visual relations involving unseen predicates in real-world scenarios. Ad-

ditionally, some recent works [23, 28] have attempted to handle the challenge of unseen verbs in human-object interaction detection by employing CLIP as a teacher model for feature representation. However, these approaches often neglect to explore the synergies between open-vocabulary SGG and VL tasks. In this study, we introduce a unified framework leveraging VLMs to tackle open-vocabulary SGG and enhance the reasoning capabilities of VLMs.

3. Preliminary

3.1. Scene Graph Generation

The goal of scene graph generation is to generate a scene graph $\mathcal{G}_{sg} = \{\mathcal{V}_e, \mathcal{R}\}$ from an image, which consists of visual relationships $\mathcal{R} = \{\mathbf{r}_{ij}\}_{i \neq j}$ and N^v entities $\mathcal{V} = \{\mathbf{v}_i\}_{i=1}^{N^v}$. The relation triplet $\mathbf{r}_{ij} = (\mathbf{v}_i, c_{ij}^e, \mathbf{v}_j)$ represents the relationship between the i -th and j -th entities, with a predicate category denoted as c_{ij}^e . The entity $\mathbf{v}_i = \{c_i^v, \mathbf{b}_i\}$ consists of a category label c_i^v in entity category space \mathcal{O}^v and a bounding box \mathbf{b}_i , indicating its location in the image. The predicate category c_{ij}^e belongs to a category space \mathcal{O}^e . In this work, we mainly focus on open-vocabulary predicate SGG settings. For open-vocabulary predicate SGG, the predicate category space \mathcal{O}^e is divided into: seen base classes \mathcal{O}_b^e , and novel unseen classes \mathcal{O}_n^e . Unlike previous works [7, 52] tackle detecting relations with entity category \mathcal{O}^v with unseen open-vocabulary set, or classifying novel predicates [7, 48] from given GT entities, we explore a more realistic and challenging setting that generates scene graphs with unseen predicate concepts.

3.2. Vision-language Models

In this work, we use detector-free BLIP [17] and instruct-BLIP as base VLM, which achieves state-of-the-art performance on many VL tasks. It performs the *image to sequence generation* task with a vision encoder-text decoder architecture. The vision encoder, denoted as $\mathcal{F}_{v.enc}$, maps the input image \mathbf{I} to vision features $\mathbf{Z}^v \in \mathbb{R}^{M \times d}$. For sequence generation, the text decoder with token predictor $\mathcal{F}_{t.dec}$ outputs hidden states of sequence $\mathbf{H}^s = [\mathbf{h}_0, \dots, \mathbf{h}_{T'}] \in \mathbb{R}^{T' \times d}$ from encoder features \mathbf{Z}^v , by auto-regressive generation methods such as beam search or nucleus sampling. The token sequence $\mathbf{t}^s = [t_0, \dots, t_{T'}] \in \mathbb{D}^{T'}$ with token classification score $\mathbf{P}^s = [\mathbf{p}_0, \dots, \mathbf{p}_{T'}] \in \mathbb{R}^{T' \times |\mathcal{C}_{voc}|}$ in vocabulary space \mathcal{C}_{voc} .

4. Our Approach

4.1. Method Overview

Our proposed PGSG framework formulates open-vocabulary SGG as an image-to-sequence generation task, which is achieved by a network \mathcal{F}_{sg} with parameters Θ_{sg} . Moreover, with a unified formulation, the PGSG framework enables a seamless transfer of explicit scene graph representation knowledge by Θ_{sg} to VL tasks.

Specifically, our framework is composed of three components: scene graph sequence generation, relation triplet construction, and adaptation to downstream VL tasks, as shown in Fig. 2. First, we generate the scene graph sequence \mathbf{s}_{sg} from the input image \mathbf{I} by image-to-text generation of VLM with the scene graph sequence prompts (Sec. 4.2). Second, the relationship construction module extracts \mathcal{R} from the scene graph sequence \mathbf{s}_{sg} to construct the scene graph \mathcal{G}_{sg} (Sec. 4.3). In addition, we introduce the learning and inference pipeline of our SGG framework in Sec. 4.4. Finally, we adapt the VLM with Θ_{sg} for transferring the explicit scene graph representation knowledge to VL tasks (Sec. 4.5).

4.2. Scene Graph Sequence Generation

The scene graph sequence \mathbf{s}_{sg} is composed of a set of relation triplets \mathcal{R} from \mathcal{G}_{sg} . To achieve this, we propose a scene graph sequence prompt:

“Generate the scene graph of [triplet sequence] and [triplet sequence] ...”

It consists of two primary components: the prefix instruction “Generate the scene graph of” and K relation triplet sequences, separated by “and” or “,”. Specifically, each triplet sequence is transformed from the \mathbf{r}_{ij} by the natural language grammar of subject-predicate-object:

“ \mathbf{t}_i^v [ENT] \mathbf{t}_{ij}^e [REL] \mathbf{t}_j^v [ENT]”

The token sequences \mathbf{t}_i^v , \mathbf{t}_{ij}^e , and \mathbf{t}_j^v represent the tokenized category names of the subject, object entities c_i^v, c_j^v , and predicate c_{ij}^e . We also introduce specified relation-aware tokens, [ENT] and [REL], to represent the compositional structure of relationships and the position of entities.

The text decoder of VLM takes \mathbf{Z}^v and the prefix instruction as input and generates \mathbf{s}_{sg} . This procedure is similar to standard image captioning. By using the vocabulary space of natural language, we can use generalized semantic representation to generate the visual relationship. The following modules extract both the spatial and categorical information and construct the instance-aware relation triplet from the sequence using relation-aware tokens.

4.3. Relationship Construction

As referred to in Sec. 3.1, the standard visual relation triplet r_{ij} contains the predicate category label c_{ij}^e , entities with position $\mathbf{b}_i, \mathbf{b}_j$, and category c_i^v, c_j^v . The relation construction aims to extract spatial and category labels from \mathbf{s}_{sg} to construct the relation triplets. This relationship construction has two submodules: 1) *Entity Grounding Module*, which outputs entity positions, and 2) *Category Convert Module*, which converts language prediction from vocabulary space to category space.

4.3.1 Entity Grounding Module

The *entity grounding module* predicts the bounding box of an entity to ground the generated relations within the scene

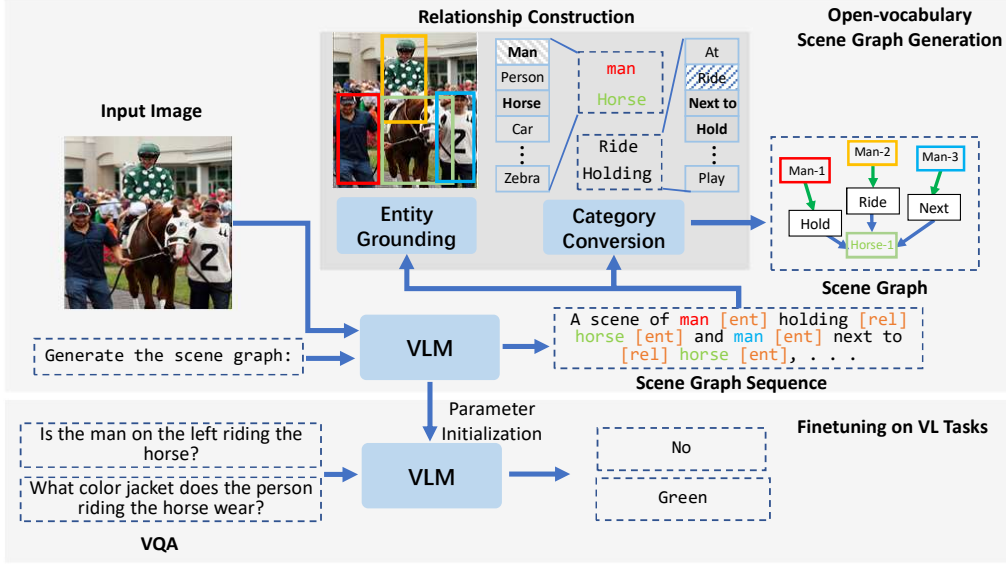


Figure 2. **Illustration of overall pipeline of our PGSG.** We generate scene graph sequences from the images using VLM. Then, the relation construction module grounds the entities and converts categorical labels from the sequence. For VL tasks, the SGG training provides parameters as initialization for VLM in fine-tuning.

graph sequence. Unlike existing multi-output VLMs [36, 44, 46] that output sequences comprising a mixture of coordinates and words, the entity grounding module predicts bounding boxes $\mathbf{B} \in \mathbb{R}^{2N \times 4}$ of entity token sequence “ \mathbf{t}_i^v [ENT]”, “ \mathbf{t}_j^v [ENT]” from N relation triplets of \mathbf{s}_{sg} . This process significantly enhances the quality of modeling spatial coordinates for these entities.

As shown in Fig. 3, for each token sequence “ \mathbf{t}_i^v [ENT]”, we extract the \mathbf{b}_i by hidden states $\mathbf{H}_i = [\mathbf{h}_1^{t_i^v}, \dots, \mathbf{h}_{T_v}^{[ENT]}] \in \mathbb{R}^{T_v \times d}$ of , where T_v is length of \mathbf{t}_i^v . We utilize the token hidden states as queries \mathbf{Q} to re-attend the image features \mathbf{Z}_v through the attention mechanism to more accurately locate the objects. We first transform $[\mathbf{h}_1^{t_i^v}, \dots, \mathbf{h}_{T_v}^{[ENT]}]$ into query vector $\mathbf{q}_i^{ent} \in \mathbb{R}^{d_q}$ by average pooling and linear projection:

$$\mathbf{q}_i^{ent} = \frac{1}{T_v} \sum_{k=1}^{T_v} \mathbf{h}_k^{t_i^v} \cdot \mathbf{W}_q^T. \quad (1)$$

Subsequently, we decode the \mathbf{b}_i by cross-attention between queries of $2N$ entity token sequences: $\mathbf{Q} = \{\mathbf{q}_1^{ent}, \dots, \mathbf{q}_{2N}^{ent}\} \in \mathbb{R}^{2N \times d}$ and \mathbf{Z}_v using transformer encoder $\mathcal{F}_{enc}(\cdot)$ and decoder $\mathcal{F}_{dec}(\cdot)$. Finally, the \mathbf{B} is predicted by the feed-forward network $\text{FFN}(\cdot)$.

$$\hat{\mathbf{Q}} = \mathcal{F}_{dec}(\mathbf{Q}, \mathcal{F}_{enc}(\mathbf{Z}_v)), \quad (2)$$

$$\mathbf{B} = \text{FFN}(\hat{\mathbf{Q}}). \quad (3)$$

By introducing this dedicated module, we are able to generate standard scene graphs with instance-level spatial prediction from an image-level scene graph sequence.

4.3.2 Category Conversion Module

During the sequence generation, the category of entities and predicates is represented by a token from the language vocabulary space \mathcal{C}_{voc} . However, the SGG benchmark category concept space is only part of the language vocabulary. To align the open-vocabulary predictions with the category space of the target SGG benchmark, we need to convert the vocabulary token into a category label. As illustrated in Fig. 3, the *Category Conversion Module* maps the tokens prediction score from the vocabulary space \mathcal{C}_{voc} to the entity and predicate category space $\mathcal{O}^v, \mathcal{O}^e$ in the target SGG benchmark, respectively.

The category conversion module is parameter-free, which transforms the vocabulary score $\mathbf{P}_i^s, \mathbf{P}_{ij}^s \in \mathbb{R}^{T_v \times |\mathcal{C}_{voc}|}$ of each entity and the predicate token sequences \mathbf{t}_i^v and \mathbf{t}_{ij}^e into the categorical prediction scores $\mathbf{p}_i^v \in \mathbb{R}^{|\mathcal{O}^v|}$ and $\mathbf{p}_{ij}^e \in \mathbb{R}^{|\mathcal{O}^e|}$ respectively. To achieve this transfer, we first tokenize target category sets, include entity \mathcal{O}^v , and predicate \mathcal{O}^e into the token sequences \mathbf{t}_c^e and \mathbf{t}_c^v by tokenizer of VLM.

$$\mathbf{t}_c^e = [t_1^1, t_2^1, t_1^2, \dots, t_{T_{C_e}}^{C_e}], \quad \mathbf{t}_c^v = [t_1^1, t_2^1, t_1^2, \dots, t_{T_{C_v}}^{C_v}]. \quad (4)$$

Here, the superscript c denotes the category index, and T is the sequence length for each tokenized category name. Subsequently, we compose \mathbf{p}_r^e and \mathbf{p}_i^v from \mathbf{p}^s using the corresponding token indices of \mathbf{t}_c^v and \mathbf{t}_c^e :

$$\mathbf{p}_{ij}^e = \left\{ \frac{\beta_1}{2} \sum_{i=1}^2 \mathbf{P}_{ij}^s[t_i^1]; \dots; \frac{\beta_{C_e}}{T_{C_e}} \sum_{i=1}^{T_{C_e}} \mathbf{P}_{ij}^s[t_i^{C_e}] \right\}. \quad (5)$$

Here, $[\cdot]$ represents the indexing operation, and $\{\cdot; \cdot\}$ denotes the concatenation operation. The same procedure is

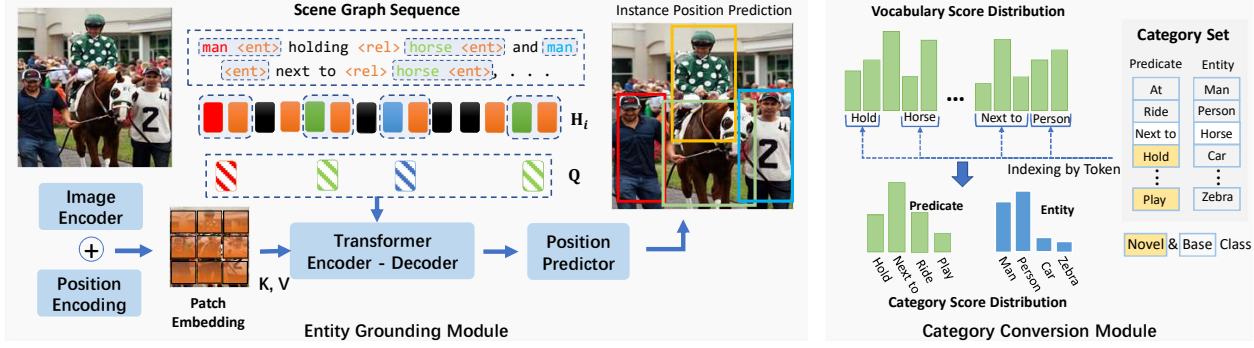


Figure 3. **Illustration of entity grounding module and category conversion module.** Left): The entity grounding module localizes the entities within scene graph sequences by predicting their bounding boxes. Right): The category conversion module maps the vocabulary sequence prediction into categorical prediction.

applied to the scores for entity classification \mathbf{p}_i^v . Additionally, if the generated tokens t_i^v and $t_{i_j}^e$ exactly match the given category name, we amplify the score of this category by β_i . To this end, the category conversion module enables the evaluation and analysis of SGG performance with respect to the defined categories.

4.4. Learning and Inference

4.4.1 Learning

For SGG training, we use a multi-task loss consisting of two parts: the loss for standard next token prediction language modeling \mathcal{L}^{lm} and the loss for the entity grounding module \mathcal{L}^{pos} . The \mathcal{L}^{lm} loss follows the same definition as in the VLM backbone, which is a standard self-regressive language modeling loss. It predicts the next token t_i based on previously generated tokens ($\mathbf{t} = [t_1, t_2, \dots, t_{i-1}]$) and visual features \mathbf{Z}^v . This prediction is represented as: $t_i = \mathcal{F}_{vlm}(\mathbf{Z}^v, \mathbf{t}; \Theta_{sg})$. To train this VLM, we optimize its parameters (Θ_{sg}) by maximizing the likelihood of correctly predicting each subsequent token in the sequence of length K : $\sum_{i=1}^K \log P(t_i | \mathbf{Z}, t_1, t_2, \dots, t_{i-1}; \Theta_{vlm})$. Compared to previous SGG methods, our self-regression-based approach effectively captures the semantic relationships between predicates and entities. The loss \mathcal{L}^{pos} is for the bounding box regression of entity grounding module. For each box prediction \mathbf{B} of scene graph sequence, we have GIOU and L1 loss for calculating the distance \mathbf{B} and GT \mathbf{B}_{gt} : $\mathcal{L}^{pos} = \text{GIOU}(\mathbf{B}, \mathbf{B}_{gt}) + \|\mathbf{B} - \mathbf{B}_{gt}\|_1$.

4.4.2 Inference

During inference, we first generate the scene graph sequence using VLMs. The following relationship construction module extracts category labels and spatial positions for constructing visual relations from the sequences.

Scene Graph Sequence Generation The scene graph sequence is generated using the image-to-text generation of VLMs. Our goal is to generate a diverse set of relationships to create a more comprehensive scene graph representation. To achieve this, we employ the nucleus sampling strat-

egy [9] with multiple round sequence generation. Specifically, for each image, we generate M sequences, whose maximum length is L . This strategy encourages the model to generate diversely, allowing scene graph predictions to capture various visual entities and their relationships.

Relationship Triplet Construction To construct relationship triplets $\hat{\mathbf{r}}_{ij} = (\mathbf{v}_i, c_{ij}^e, \mathbf{v}_j)$ from the scene graph sequence, we initially extract the token sequence of subject and object entities t_i^v , t_j^v , and predicate t_{ij}^e . Specifically, we employ a heuristic rule to match patterns such as “*subject [ENT] predicate [REL] object [ENT]*”, allowing us to extract these triplets. Subsequently, utilizing a relationship construction module, we extract $\mathbf{b}_i, \mathbf{p}_i^v$ of entities \mathbf{v}_i , and $\mathbf{b}_j, \mathbf{p}_j^v$ for \mathbf{v}_j , with predicate scores \mathbf{p}_{ij}^v . To obtain the category labels of triplets c_i^v, c_j^v, c_{ij}^e , we select the top 3 categories according to prediction scores: $\mathbf{p}_i^v, \mathbf{p}_j^v$, and \mathbf{p}_{ij}^v .

Post-Processing We refine the initial triplets to obtain the final relationship triplets through a process of filtering and ranking. Firstly, we remove self-connected relationships where the subject and object entities are the same. Then, we apply a non-maximum suppression (NMS) strategy to eliminate redundant relationships, resulting in a more concise predicted scene graph. Finally, we rank all relationships based on the score of each component triplet. This score, denoted as S^t , is calculated as the product of prediction scores for the involved entities s_i^v, s_j^v and the predicate s_p^e , extracted from category score distributions $\mathbf{p}_i^v, \mathbf{p}_j^v, \mathbf{p}_{ij}^e$ using their respective category indices.

4.5. Adaptation to Downstream VL Task

In our framework, depicted in Fig. 2, we employ a unified image-to-text generation approach, facilitating seamless knowledge transfer for relation modeling. We initialize the visual encoder Θ_{v_enc} and text decoder Θ_{t_dec} parameters of the Scene Graph Generation (SGG) model, denoted as Θ_{sg} , for VLMs fine-tuning on VL tasks. Additionally, we employ initial pre-trained weights for the token predictors of text decoders to maintain vocabulary prediction quality during fine-tuning. Modules of the VLM not utilized

D	S	M	Novel+base		Novel
			mR50/100	R50/100	mR50/100
VG	PCIs	CaCao	10.3 / 12.6	-	-
		SVRP	8.3 / 10.8	33.5 / 35.9	-
		PGSG*	10.8 / 13.9	26.9 / 33.9	5.2 / 7.7
	SGCls	SVRP	3.2 / 4.5	19.1 / 21.5	-
		PGSG*	8.4 / 11.0	22.6 / 27.2	4.8 / 6.0
		VS3	5.1 / 5.7	11.0 / 12.8	0.0 / 0.0
PSG	SGDet	SGTR [†]	6.4 / 8.4	14.2 / 18.2	0.0 / 0.0
		PGSG	13.5 / 16.4	18.0 / 20.2	7.4 / 11.3
		SVRP	-	-	-
Oiv6		SGTR [†]	11.0 / 16.7	36.1 / 38.4	0.0 / 0.0
		PGSG	20.8 / 23.0	41.3 / 43.3	3.8 / 8.9

Table 1. **The open-vocabulary scene graph generation on VG, PSG, and OIV6 datasets** † denotes the method is reproduced with BLIP visual encoder ViT as the backbone; D denotes the dataset name; S is the SGG setting; M represents the model; * denotes that we use the grounding truth assignment on predictions.

during SGG training, such as the text encoder, retain their initial pre-trained weights. This adaptive knowledge sharing from SGG parameters during fine-tuning enhances the performance of VL tasks.

5. Experiments

5.1. Experiments Configuration

We evaluate our method on both SGG and downstream VL tasks. For SGG benchmarks, we evaluate PGSG on three benchmarks: Panoptic Scene Graph Generation (PSG) [40], Visual Genome (VG) [15], and OpenImage V6 (Oiv6) [16]. We randomly select 50% predicate categories as novel classes for open-vocabulary predicate SGG.

We assess our method’s performance in the SGG task using three evaluation protocols: SGDet, PCIs, and SGCls, as proposed by [39]. We calculate overall recall (R@K) and class-balance metric mean recall (mR@K) on the top K predictions. Furthermore, we present mR@K specifically for novel classes to illustrate the model’s performance in generalizing to unseen predicates.

For VL tasks, we inspect our model on visual grounding on RefCOCO/+g [26, 47], visual question answering on GQA [11] and image captioning on COCO [4].

We use the BLIP [17] as our VLM backbone, which employs ViT-B/16 as the visual backbone and BERT_{base} as the text decoder. For downstream VL tasks, we initialize the parameters with the pre-trained SGG model. More implementation details can be found in the supplementary materials.

5.2. Open-vocabulary SGG

Setup We evaluate our design on the VG, PSG, and Oiv6 datasets by comparing it with OV-SGG methods: CaCao[48], SVRP[7], and VS3[52] in Tab 1. The CaCao and SVRP have the capability for open-vocabulary

predicate classification with GT entities. In contrast, the VS3 achieves SGDet on close-set predicate, where relation triplets have unseen entity categories. Additionally, we adapt a one-stage SGG method (SGTR [21]) for predictive open-vocabulary SGG since it has the most balanced performance in a close-vocabulary SGG setting. To apply close-set SGG methods to predicate open vocabulary, we train models on base categories only and set the classifier weights of novel categories to zero at inference time.

VG Our method demonstrates strong performance compared to previous methods on PCIs, SGCls, and SGDet settings on VG. Since the VLM cannot take a bounding box as input, we simulate entity-based relations in PCIs and SGCls settings by directly replacing entity prediction with GT entities. Our method improves base+novel class mR@100 and R@100 by 6.5 and 5.7 over SVRP in SGCls. We outperform SVRP’s task-specific pretrain model with a general pretrain VLM in an image-to-text generation fashion.

For the SGDet setting, we compare with two baselines: VS3 and SGTR. Our method outperforms both methods, achieving 2.9 and 0.9 on mR@100 and R@100 for all categories. It’s worth noting that VS3 and SGTR are unable to generalize to novel classes and can only predict relationships for base categories.

PSG Our PGSG method achieves significant performance gains for both seen base and unseen novel categories on the PSG dataset. We observe improvements of 9.4 and 14.3 on mR@100 and R@100 compared with SGTR equipped with a pre-trained ViT backbone in the whole category set, respectively. For novel predicate categories, PGSG outperforms the strong baseline SGTR method by 4.5 on mR@100.

OpenImage Our approach excels on the OpenImage V6 benchmark (see Tab. 1). Comparing our design to the SGTR baseline, we achieve better mR@100 and R@100 with margins of 5.3 and 4.9.

5.3. Close-vocabulary SGG

Our method is also tested in a closed-vocabulary SGG. Our method generates relationships using VLMs without explicitly pairing entity proposals, so we compare it to one-stage SGG methods with top-down architecture, such as PSGTR [40], PSGFormer [40], SGTR [21], ReITR [5], ISG [14], SSRCNN [34]. In zero-shot triplet SGG, our method shows good compositional generalization. To ensure a fair comparison, we substituted the category conversion module for a close-set classifier, aligning with previous SGG methods. As shown in Tab. 2, the BGSg with the close-set classifier achieves competitive or superior performance when compared to SGG SOTA methods, SGTR, in a standard SGG setting.

Zero-shot Triplets In the middle part of Tab. 2, we evaluate our method’s generalization ability through zero-shot triplet SGG. On the PSG dataset, our PGSG method en-

D	B	M	Zs-SGG	Close-SGG	
			zR50/100	mR50/100	R50/100
PSG	R101	PSGTR	3.1/6.4	20.3/21.5	32.1/35.3
		PSGFormer	2.2/4.9	19.3/19.7	20.4/20.7
		SGTR	4.1/5.8	24.3/27.2	33.1/36.3
	ViT-B*	SGTR	3.1/3.6	22.1/23.5	30.1/32.0
		PGSG	5.1/7.5	14.9/18.1	26.0/28.9
		PGSG-c	6.8/8.9	20.9/22.1	32.7/33.4
X101-FPN	SSRCNN	3.1/4.5	18.6/22.5	23.7/27.3	
R50	SVRP	-	10.5/12.8	31.8/35.8	
Swin-L	VS3	-	-	34.5/39.2	
VG	R101	RelTR	2.6/3.4	6.8/10.8	21.2/27.5
		ISG	2.7/3.8	8.0/8.8	29.7/32.1
		SGTR	2.5/5.8	12.0/15.2	24.6/28.4
	ViT-B*	SGTR	2.1/4.9	8.7/10.1	18.1/20.4
		PGSG	6.2/8.5	8.9/11.5	16.7/21.2
		PGSG-c	4.8/7.6	10.5/12.7	20.3/23.6

Table 2. **The close-vocabulary SGGDet performance on VG and PSG dataset.** D: dataset name; B: visual backbone; M: Method; ViT-B*: the use BLIP ViT visual encoder as backbone. -c: denotes use the close-set classifier.

hances zR@100 by 4.1 over the ViT SGTR baseline and matches the ResNet-101 PSGTR model, despite lower image input. PGSG performed well on the VG dataset, outperforming SSRCNN by 4.0 on zR@100 and achieving a 3.6 margin over the ViT SGTR baseline.

Discussion We also evaluate our method against previous SGG approaches under standard close-vocabulary SGG settings in the Tab. 2. While our method performs similarly to prior works using the same ViT backbone in closed-vocabulary settings, there remains a gap between our results and previous SOTA performances.

- Limited input resolution: Transformer-based VLMs like BLIP use smaller images (384x384), while traditional SGG methods (e.g., ResNet-101 FPN, 800x1333) may miss or mislocalize small or ambiguous objects, which have been shown to pose challenges in SGG tasks [19, 21].
- Single-stage training: In PGSG, the BLIP is used for SGG training directly, as opposed to the conventional two-stage training. However, this approach imposes limitations on the model’s ability to accurately detect small objects.

5.4. Ablation Study

We conduct an ablation study to investigate the impact of different model structures and the effectiveness of SGG training. Specifically, we compare the design choices regarding the entity grounding module with the scene graph prompt on the PSG dataset, which is presented in Tab. 3.

Entity Grounding Module First, we evaluate the effectiveness of *Entity Grounding with Relation-aware Tokens* for predicting entity location in the lower part of Tab. 3. Our approach demonstrates notable improvements in entity detection quality compared to previous methods such as UniTab and OFA, which rely on sequence generation by VLMs in-

Layer of Entity Grounding Module						
L	mR20	mR50	mR100	R20	R50	R100
0	2.0	3.8	10.3	7.1	13.2	20.1
3	3.2	8.1	15.9	11.0	18.7	27.9
6	3.0	8.4	17.0	11.1	19.2	28.4
12	3.2	8.1	16.7	12.1	18.9	28.2
Relation-aware Token for Entity Grounding						
C	mR20	mR50	mR100	R20	R50	R100
mixture	1.2	3.1	8.5	6.4	11.8	16.3
tokens	3.0	8.4	17.0	11.1	19.2	28.4

Table 3. **Ablation study for entity grounding module on PSG dataset.** L: indicates the layer of transformer within entity grounding module. C: denotes the different design choices for entity grounding.

SL	Open Vocab SGG			ZS Trp. R50/100	Time Sec
	Novel+Base R50/100	mR50/100	Novel mR50/100		
1024	16.9/18.8	22.9/25.0	6.7/9.9	3.6/6.4	6.9
768	16.8/18.0	23.6/24.9	7.9/10.0	8.5/8.1	4.8
512	15.3/15.6	22.0/22.5	7.4/8.1	3.5/4.6	2.2
256	12.8/12.8	18.3/18.3	5.3/5.3	2.4/2.4	1.8

Table 4. **The SGG performance and inference time with different length of output sequences.** SL: output sequence lengths; Sec: inference time of second per image.

corporating both coordinates and text. Moreover, we investigate the impact of the number of transformer layers (L) in the entity grounding module. As seen in the upper part of Tab. 3, we vary L from 0 to 12. Here, $L = 0$ signifies direct prediction of the bounding box from the grounding entity query Q , with a performance plateau at $L = 6$.

Effectiveness of SGG Supervision To evaluate the impact of scene graph supervision on downstream VL tasks, we fine-tune the BLIP model after training it on entity-only sequence prompts with the same format: “A photo of entity₁ [ENT] and entity₂ [ENT]...”. Fine-tuning is performed on the GQA dataset. Results from Tab. 7 show that positional supervision does not lead to significant improvements over initial pre-training. In contrast, scene graph supervision enhances the performance of GQA, especially for relation and attribute questions.

5.5. Model Analysis

In this section, we analyze our sequence generation-based SGG framework. We explore the impact of different prefix prompts on scene graph generation and evaluate the diversity, quality, and time consumption of generated scene graph sequences. We also present the visualization of generated sequences in supplementary.

Quality of Scene Graph Sequence We use a heuristic rule to parse the triplet sequence based on the scene graph prompt. To showcase our approach’s effectiveness, we track [REL] token occurrences and the number of relation triplets across different sequence lengths, as shown in Tab. 5. The

SL.	#Trip.	#Uni.Trip.	#[REL]	% Valid
1024	87.2	53.4	95.1	95.3%
768	76.1	40.2	79.2	96.0%
512	45.0	34.3	46.2	97.4%
256	22.4	20.6	23.3	96.1%
128	17.4	11.3	18.7	93.0%

Table 5. **Analysis for quality of generated Scene graph sequences.** SL.: Generated sequence Length; # Trip.: Average of number of relation triplets; # U. Trip.: Average of number of unique relation triplets; # [REL]: Average of number of occurrence of [REL] tokens; % Valid: percentage of valid relation triplets.

Model	RC			RC+			RCg
	val	testA	testB	val	testA	testB	val
init	83.1	86.0	77.0	77.1	83.1	69.8	74.3
PGSG	86.0	88.9	82.4	79.8	84.3	72.3	77.8

Table 6. **Performance of visual grounding task.** RC, RC+ and RCg represents the RefCOCO, RefCOCO+ and RefCOCOg datasets respectively.

results demonstrate that VLM generates the formatted scene graph sequence, with the heuristic rule efficiently retrieving most relation triplets.

Time Consumption of Sequence Generation Our framework also achieves a good trade-off between performance and computational cost by varying the lengths of the generated sequences. As shown in Tab. 4, PGSG maintains comparable performance when the output length is reduced by 50%, with comparable inference times to 2-stage SGG methods (more comparison is in supplementary).

5.6. Downstream VL-Tasks

We apply PGSG to VL tasks, including visual grounding, visual question answering, and image captioning. Through comparisons with the initial pre-trained model, we assess how explicit scene representation modeling influences vision-language reasoning.

Visual Question Answering On the GQA benchmark, PGSG outperforms the initial pre-trained model by 1.7 in overall accuracy, with the largest gains in relation (1.9) and object (1.3) questions. Additionally, our unified SGG training also boosts BLIPv2 VLM’s zero-shot performance on GQA from 32.3 to 33.9, with improvements across all question types.

Image Captioning We conduct the evolution of image captioning on COCO [4]. The results demonstrate that our framework achieves comparable performance to the initial pre-trained VLMs, including BLIP and BLIPv2, especially the SPICE metric [1], which indicates the scene description quality of the generated caption.

Visual Grounding Our method surpassed initial pre-trained model by 2.9 and 5.4 on the validation and test A/B splits of RefCOCO and by 2.7, 1.2, and 2.5 for RefCOCO+ and RefCOCOg, respectively (see Tab. 6). These results underscore the effectiveness of knowledge transfer in SGG mod-

Model	Question Type					
	Relation	Attribute	Object	Category	Global	Overall
BLIP	Finetuned					
Init	54.9	64.5	86.6	62.2	67.6	62.5
PGSG-E	55.8	64.3	87.5	62.6	67.9	63.1
PGSG	56.8	66.3	87.9	62.8	68.1	64.2
BLIPv2	Zeroshot					
Init	29.8	28.8	53.6	31.1	28.2	32.3
PGSG	31.3	29.6	53.9	34.3	26.1	33.9

Table 7. **Performance comparison between initial pre-trained VLM and PGSG on GQA.** We report the answering accuracy respectively according to the question type proposed by benchmark. E: denotes the entity-only scene graph sequence for ablation;

		Bleu4	CIDEr	SPICE
BLIP	Init	0.44	132.1	23.6
	Det-T	0.43	132.0	23.5
	PGSG	0.41	131.2	23.9
BLIP v2	Init	0.42	145.1	25.6
	PGSG	0.43	145.2	26.0

Table 8. **Performance comparison on COCO image captioning.** eling, enhancing reasoning and localization following unified SGG training for visual grounding.

6. Conclusion

In this paper, we present a VLM-based open-vocabulary scene graph generation model that formulates the SGG through an image-to-graph translation paradigm, thereby addressing the difficult problem of generating open-vocabulary predicate SGG. The holistic architecture unifies the SGG and subsequent VL tasks, thereby enabling explicit relation modeling to benefit the VL reasoning tasks. Extensive experimental results validate the validity of the proposed model.

Discussion of Limitations:

Performance on Close-vocabulary SGG. Due to the inherent weakness of the visual backbone and labeling noise, our PGSG performs suboptimal on VG and PSG datasets in a standard setting. Improving the perception performance of VLMs with high resolution input is a challenging and important issue for future research.

Study on More VL Models and Tasks. Furthermore, our study primarily focuses on the BLIP applied to the SGG and BLIP and BLIPv2 for VL tasks: VQA, visual grounding, and image captioning. However, our approach exhibits potential for extension to other VL models and tasks. Exploring these possibilities remains for future research.

Acknowledgement: This work was supported by NSFC under Grant 62350610269, National Key R&D Program of China 2022ZD0161600, Shanghai Frontiers Science Center of Human-centered Artificial Intelligence, MoE Key Lab of Intelligent Perception and Human-Machine Collaboration (ShanghaiTech University), and Shanghai Postdoctoral Excellence Program 2022235.

References

- [1] Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. Spice: Semantic propositional image caption evaluation. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part V 14*, pages 382–398. Springer, 2016. [8](#)
- [2] Tianshui Chen, Weihao Yu, Riquan Chen, and Liang Lin. Knowledge-embedded routing network for scene graph generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6163–6171, 2019. [2](#)
- [3] Vincent S Chen, Paroma Varma, Ranjay Krishna, Michael Bernstein, Christopher Re, and Li Fei-Fei. Scene graph prediction with limited labels. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2580–2590, 2019. [1](#)
- [4] 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. [6](#), [8](#), [11](#)
- [5] Yuren Cong, Michael Ying Yang, and Bodo Rosenhahn. Reltr: Relation transformer for scene graph generation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023. [6](#)
- [6] Yuyu Guo, Lianli Gao, Xuanhan Wang, Yuxuan Hu, Xing Xu, Xu Lu, Heng Tao Shen, and Jingkuan Song. From general to specific: Informative scene graph generation via balance adjustment. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16383–16392, 2021. [1](#)
- [7] Tao He, Lianli Gao, Jingkuan Song, and Yuan-Fang Li. Towards open-vocabulary scene graph generation with prompt-based finetuning. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXVIII*, pages 56–73. Springer, 2022. [1](#), [2](#), [3](#), [6](#)
- [8] Marcel Hildebrandt, Hang Li, Rajat Koner, Volker Tresp, and Stephan Günnemann. Scene graph reasoning for visual question answering. *arXiv preprint arXiv:2007.01072*, 2020. [1](#)
- [9] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. *arXiv preprint arXiv:1904.09751*, 2019. [5](#)
- [10] Drew Hudson and Christopher D Manning. Learning by abstraction: The neural state machine. *Advances in Neural Information Processing Systems*, 32, 2019. [1](#)
- [11] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6700–6709, 2019. [1](#), [6](#), [11](#)
- [12] 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](#), [2](#)
- [13] Xuan Kan, Hejie Cui, and Carl Yang. Zero-shot scene graph relation prediction through commonsense knowledge integration. In *Machine Learning and Knowledge Discovery in Databases. Research Track: European Conference, ECML PKDD 2021, Bilbao, Spain, September 13–17, 2021, Proceedings, Part II 21*, pages 466–482. Springer, 2021. [2](#)
- [14] Siddhesh Khandelwal and Leonid Sigal. Iterative scene graph generation. *arXiv preprint arXiv:2207.13440*, 2022. [6](#)
- [15] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123(1):32–73, 2017. [2](#), [6](#), [11](#)
- [16] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Alexander Kolesnikov, Tom Duerig, and Vittorio Ferrari. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *International Journal of Computer Vision(IJCV)*, 2020. [2](#), [6](#)
- [17] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR, 2022. [3](#), [6](#), [11](#)
- [18] Lin Li, Long Chen, Yifeng Huang, Zhimeng Zhang, Songyang Zhang, and Jun Xiao. The devil is in the labels: Noisy label correction for robust scene graph generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18869–18878, 2022. [2](#)
- [19] Lin Li, Long Chen, Hanrong Shi, Hanwang Zhang, Yi Yang, Wei Liu, and Jun Xiao. Nicest: Noisy label correction and training for robust scene graph generation. *arXiv preprint arXiv:2207.13316*, 2022. [1](#), [2](#), [7](#)
- [20] 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. [2](#), [11](#)
- [21] Rongjie Li, Songyang Zhang, and Xuming He. Sgtr: End-to-end scene graph generation with transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19486–19496, 2022. [6](#), [7](#)
- [22] Wei Li, Haiwei Zhang, Qijie Bai, Guoqing Zhao, Ning Jiang, and Xiaojie Yuan. Ppdl: Predicate probability distribution based loss for unbiased scene graph generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19447–19456, 2022. [2](#)
- [23] Yue Liao, Aixi Zhang, Miao Lu, Yongliang Wang, Xiaobo Li, and Si Liu. Gen-vlkt: Simplify association and enhance interaction understanding for hoi detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20123–20132, 2022. [3](#)
- [24] 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. [2](#)

- [25] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. **11**
- [26] Junhua Mao, Jonathan Huang, Alexander Toshev, Oana Camburu, Alan L Yuille, and Kevin Murphy. Generation and comprehension of unambiguous object descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 11–20, 2016. **6, 11**
- [27] Kien Nguyen, Subarna Tripathi, Bang Du, Tanaya Guha, and Truong Q Nguyen. In defense of scene graphs for image captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1407–1416, 2021. **2**
- [28] Shan Ning, Longtian Qiu, Yongfei Liu, and Xuming He. Hoiclip: Efficient knowledge transfer for hoi detection with vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23507–23517, 2023. **3**
- [29] 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, 2015. **2**
- [30] Jiaxin Shi, Hanwang Zhang, and Juanzi Li. Explainable and explicit visual reasoning over scene graphs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8376–8384, 2019. **1**
- [31] 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 Conference on Computer Vision and Pattern Recognition*, pages 6619–6628, 2019. **2**
- [32] 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*, pages 3716–3725, 2020. **1, 2**
- [33] 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**
- [34] Yao Teng and Limin Wang. Structured sparse r-cnn for direct scene graph generation. *arXiv preprint arXiv:2106.10815*, 2021. **6**
- [35] Dalin Wang, Daniel Beck, and Trevor Cohn. On the role of scene graphs in image captioning. In *Proceedings of the Beyond Vision and Language: Integrating Real-world Knowledge (LANTERN)*, pages 29–34, 2019. **2**
- [36] Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In *International Conference on Machine Learning*, pages 23318–23340. PMLR, 2022. **4**
- [37] Yu-Siang Wang, Chenxi Liu, Xiaohui Zeng, and Alan Yuille. Scene graph parsing as dependency parsing. *arXiv preprint arXiv:1803.09189*, 2018. **1**
- [38] Zhecan Wang, Haoxuan You, Liunian Harold Li, Alireza Zareian, Suji Park, Yiqing Liang, Kai-Wei Chang, and Shih-Fu Chang. Sgeitl: Scene graph enhanced image-text learning for visual commonsense reasoning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 5914–5922, 2022. **1, 2**
- [39] Danfei Xu, Yuke Zhu, Christopher B Choy, and Li Fei-Fei. Scene graph generation by iterative message passing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5410–5419, 2017. **6, 11**
- [40] Jingkang Yang, Yi Zhe Ang, Zujin Guo, Kaiyang Zhou, Wayne Zhang, and Ziwei Liu. Panoptic scene graph generation. *arXiv preprint arXiv:2207.11247*, 2022. **2, 6, 11**
- [41] Sibe Yang, Guanbin Li, and Yizhou Yu. Cross-modal relationship inference for grounding referring expressions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4145–4154, 2019. **1**
- [42] Xu Yang, Kaihua Tang, Hanwang Zhang, and Jianfei Cai. Auto-encoding scene graphs for image captioning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 10685–10694, 2019. **1**
- [43] Xuewen Yang, Yingru Liu, and Xin Wang. Reformer: The relational transformer for image captioning. *arXiv preprint arXiv:2107.14178*, 2021. **1**
- [44] Zhengyuan Yang, Zhe Gan, Jianfeng Wang, Xiaowei Hu, Faisal Ahmed, Zicheng Liu, Yumao Lu, and Lijuan Wang. Unitab: Unifying text and box outputs for grounded vision-language modeling. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXVI*, pages 521–539. Springer, 2022. **4**
- [45] Yuan Yao, Ao Zhang, Xu Han, Mengdi Li, Cornelius Weber, Zhiyuan Liu, Stefan Wermt, and Maosong Sun. Visual distant supervision for scene graph generation. *arXiv preprint arXiv:2103.15365*, 2021. **2**
- [46] Yuan Yao, Qianyu Chen, Ao Zhang, Wei Ji, Zhiyuan Liu, Tat-Seng Chua, and Maosong Sun. Pevl: Position-enhanced pre-training and prompt tuning for vision-language models. *arXiv preprint arXiv:2205.11169*, 2022. **2, 4**
- [47] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context in referring expressions. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part II 14*, pages 69–85. Springer, 2016. **6, 11**
- [48] Qifan Yu, Juncheng Li, Yu Wu, Siliang Tang, Wei Ji, and Yueting Zhuang. Visually-prompted language model for fine-grained scene graph generation in an open world. *arXiv preprint arXiv:2303.13233*, 2023. **1, 2, 3, 6**
- [49] Alireza Zareian, Svebor Karaman, and Shih-Fu Chang. Bridging knowledge graphs to generate scene graphs. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020. **2**
- [50] 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. **2, 11**
- [51] 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 *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXVII*, pages 409–424. Springer, 2022. [1](#)

- [52] Yong Zhang, Yingwei Pan, Ting Yao, Rui Huang, Tao Mei, and Chang-Wen Chen. Learning to generate language-supervised and open-vocabulary scene graph using pre-trained visual-semantic space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2915–2924, 2023. [1](#), [2](#), [3](#), [6](#)
- [53] Chaofan Zheng, Xinyu Lyu, Lianli Gao, Bo Dai, and Jingkuan Song. Prototype-based embedding network for scene graph generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22783–22792, 2023. [1](#)
- [54] Yiwu Zhong, Jing Shi, Jianwei Yang, Chenliang Xu, and Yin Li. Learning to generate scene graph from natural language supervision. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1823–1834, 2021. [2](#)

A. More Experiment Configurations

Datasets and tasks We evaluate our method on both SGG task and downstream VL-tasks. For SGG task, we use two large-scale SGG benchmarks: Panoptic Scene Graph Generation (PSG) [40], Visual Genome (VG) [15]. We mainly adopt the data splits from the previous work [20, 40, 50]. For Visual Genome [15] dataset, we take the same split protocol as [39, 50] where 62,723 images are used for training, 26,446 for test, and 5,000 images sampled from the training set for validation. The most frequent 150 object categories and 50 predicates are adopted for evaluation. The Panoptic Scene Graph Generation [40] has 4,4967 images are used for training, 1,000 for test, and 3,000 images sampled from the training set for validation. There 133 object categories and 56 predicates categories in total. We use it bound box annotation for SGG task rather than segmentation masks.

For the open-vocabulary predicate SGG settings, we randomly select 30% predicate categories as novel class. For VL tasks, we inspect our model on VL-task which potentially need the visual scene representation, such as visual grounding on RefCOCO+/g [26, 47], visual question answering on GQA [11], and image captioning on COCO image caption [4].

Implementation Details We initialize our PGSG by using the BLIP [17] model with ViT-B/16 as the visual backbone and BERT_{base} as the text decoder. For scene graph training process, we use the image size of 384 times 384, an AdamW [25] optimizer with lr = 1e-5, weight decay of 0.02 with a cosine scheduler. We increase the learning rate of position adaptors to 1e-4 for faster convergence. We train our model on 4 A100 GPUs with 50 epoch. For downstream tasks fine-tuning, we following the training setup of BLIP [17]. We use the image encoder and text decoder of PGSG model, and the text encoder and word embedding remain the same as in the pre-trained BLIP model. During the scene graph sequence generation, we generate M=32 number of sequences which length is L=24. For category amplifier β_i , we set this hyper-parameter as 5.0 for entity categories and 1.0 for predicate classification.

B. More Experimental Results

In this section, we propose mre Experimental results, includes quantitative and quantitative analysis of our method.

C. Quantitative Results

C.1. Close-vocabulary SGG on OpenImage V6

In Tab. 9, we present the close-vocabulary SG-Det performance on OpenImage V6 across various visual backbones and zero-shot triplet (**Zs Trp.**) scenarios. With the same

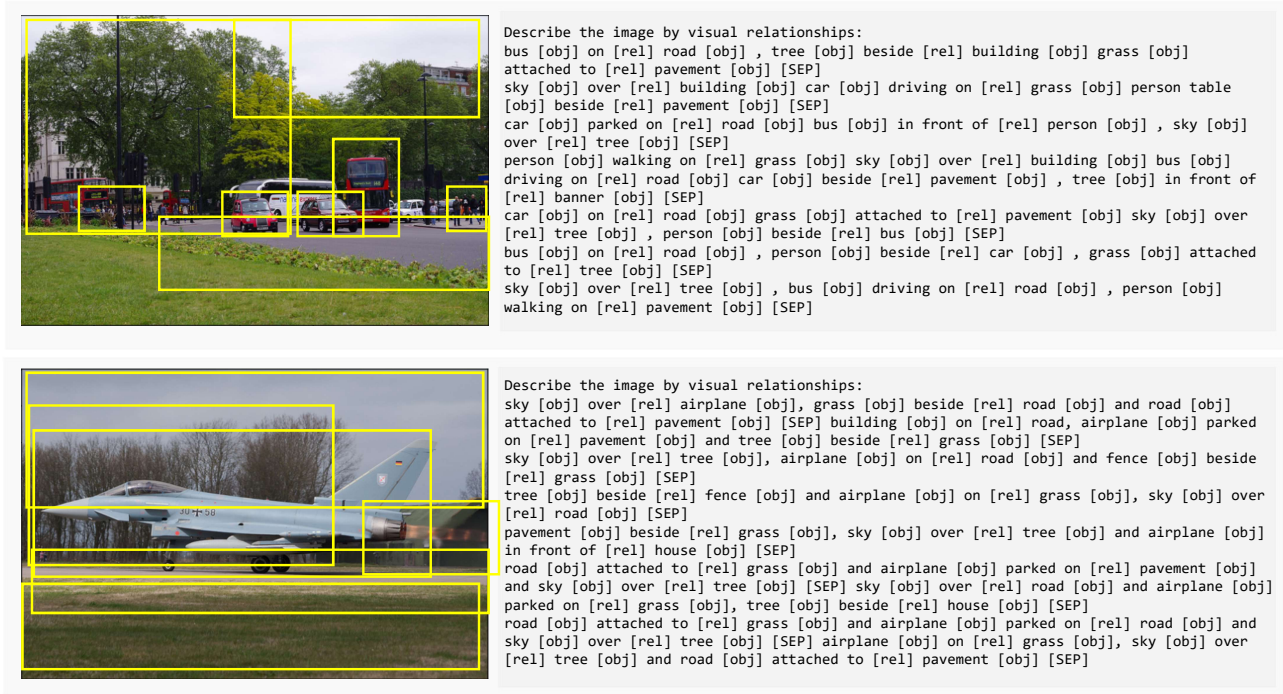


Figure 4. The visualization of scene graph sequence prediction of PGSG.

B	M	Zs Trp. R@50/100	Standard SGG				
			mR50	R50	wmAP phr	rel	score_wtd
R101	RelDN	-	39.7	72.1	28.7	29.1	38.6
	HOTR	-	36.8	52.6	21.5	19.4	26.8
	SGTR	-	38.6	59.1	36.9	38.7	42.8
ViT-B*	SGTR	19.4/31.6	30.5	52.6	28.0	22.7	30.8
	PGSG	23.1/38.6	40.7	62.0	27.8	19.7	28.7

Table 9. The close-vocabulary SG-Det performance on Open-Image V6.

Prompt	Open Vocab SGG			ZS Trp. R@50/100
	Novel+Base		Novel	
	mR@50/100	R@50/100	mR@50/100	
A	8.2/10.5	14.5/16.4	2.3/7.0	3.6/6.4
B	9.3/11.7	17.7/20.4	3.7/8.6	4.4/7.6
C	9.1/10.1	16.3/19.4	4.1/9.0	4.1/6.6

Table 10. Ablation study on different prompt for SGG task on PSG dataset. A: "A visual scene of: "; B: "Describe the image by relationships: "; C: "A picture of: "

backbone as BLIP ViT-B, our PGSG achieves comparable performance with baseline SGTR in a standard close-vocabulary setting and reasonable performance with the previous one-stage SGG method, which has a larger in-

put resolution ResNet-101 backbone. For compositional generalization setting, zero-shot triplet SGG, our method achieves a remarkable 7.0 improvement over the SGTR baseline.

C.2. Sensitivity of different Prefix Prompts

We also study the different prompt structures for generating the scene graph, as shown in Tab. 10. We experiment with the PGSG framework with different prefix instructions for the scene graph generation task. The results show that more specific instructions yield a slight improvement in performance, which indicates that our method has robustness for different instructions.

C.3. Time Complexity Comparison With Previous Method

Despite potential inference time increases due to the self-regression generation with a large model, we have effectively mitigated this issue by reducing output size. We achieve a boost in inference speed reasonable open-vocabulary SGG performance, in the Tab. 4 of the main paper. We also compare the inference times with other SGG methods, as shown in Tab. 11. The results demonstrate that PGSG attains comparable time efficiency while maintaining its competitive open-vocabulary SGG performance.

M	VCTree	GPS-Net	BGNN	PGSG	PGSG*
Time	1.69	1.02	1.32	4.8	1.8

Table 11. **the inference speed (Second per image) comparison with previous two-stage SGG methods**

C.4. Time Complexity

Despite potential inference time increases due to the self-regression generation with a large model, we have effectively mitigated this issue by reducing output size. We achieve a boost in inference speed reasonable open-vocabulary SGG performance, in the Tab. 4 of the main paper. We also compare the inference times with other SGG methods, as shown in Tab. 11. The results demonstrate that PGSG attains comparable time efficiency while maintaining its competitive open-vocabulary SGG performance.

D. Qualitative Results

We also present the qualitative analysis for the PGSG framework to take a close look at the sequence generation-based SGG framework. In Fig. 4, we show a few examples of generated sequences from our validation set of the PSG dataset. At inference time, the VLM generates scene graph sequences with entity-aware tokens as indicators by using several short token sequences with nucleus sampling, which are able to obtain diverse visual relations. The following entity grounding module extracts the boxes for each entity within the sequences.