

A Survey of Scene Graph: Generation and Application

Pengfei Xu, Xiaojun Chang, Ling Guo, Poyao Huang, Xiaojiang Chen, and Alexander G. Hauptmann

Abstract—**Scene Graph** is a data structure, which is mainly used to describe the objects, attributes and object relationships in a scene. Scene Graph is a deep representation of a scene, and is very conducive to many visual tasks, such as image retrieval, image/video captions, VQA, and even to image generation and specific relationship detection. At present, numbers of research works about scene graph are proposed, including the scene graph generation methods and the related applications. These proposed methods based on scene graph have great improvements in relative performances compared with the corresponding traditional methods, which also proves the effectiveness of scene graph in the visual understanding of a scene. Therefore, In this paper, we provide a systematic review of the existing techniques of scene graph generation and application, including not only the state-of-the arts but also those with latest trends. Particularly, we discuss the scene graph generation methods according to the inference models for visual relationship detection, and the applications of scene graph are stated according to the specific visual tasks. Finally, we point out several problems in the current scene graph generation methods, related applications and the future research directions of scene graph.

Index Terms—Scene Graph, Object Detection, Visual feature extraction, Prior Information, Visual Relationship Recognition.

1 INTRODUCTION

A scene graph is first proposed as a data structure that describes the object instances in a scene and the relationships between the objects [1]. As shown in Fig.1, a complete scene graph can represent the detailed semantics of a dataset of scenes, but not a single image or a video; and it has powerful representations that encode 2D/3D images [1], [2] and videos [3], [4] into their abstract semantic elements without any restriction on the types and attributes of objects and the relationships between objects. Fig. 1 (a) illustrates a scene graph, and we can see that a scene graph G is a data structure of directed graph, which can be defined as a tuple $G = (O, E)$, where $O = O_1, \dots, O_n$ is a set of objects detected in the images. Each object has the form $o_i = (c_i, A_i)$, where c_i and A_i are the category and attributes of the object respectively. While $E \subseteq O \times R \times O$ is a set of directed edges to represent the relationships between objects. At present, a scene graph is commonly associated to an image dataset, but not to only one image; So it can be considered as a visual understanding to relevant images. While, a part of scene graph is grounded to an image by associating the objects to the corresponding regions in an image, as shown in Fig. 1 (b). Scene graph has a powerful representations for semantic features about the scene, and is beneficial for a wide range of visual tasks.

There are some similarities of scene graphs with the commonsense knowledge graph, such as their graphical structures and constituent elements. However, scene graph is a different type of knowledge graph, which is mainly reflected in the following aspects: (a) Each node in scene graph is associated with an image region, and these nodes come in pairs, namely a subject and an object; while each node in knowledge graph is the general concept

of its semantic label. (b) In a scene graph, the directed edges represent the relationships between pairs of objects; while each edge in knowledge graph encodes a relational fact involving a pair of concepts [5].

The idea of using the visual features of different objects in the image and the relationships between them have been proposed for achieving the visual tasks of action recognition [6], image captioning [7] and other relevant computer vision tasks [8] as early as 2015. Then, Johnson et al. proposed the concept of scene graph [1], and gave the corresponding notation representations. In [1], scene graph is generated manually from a dataset of real-world scene graphs, so as to capture the detailed semantics of a scene. Since then, the research on scene graph has received extensive attentions. Subsequently, several scene graph datasets are introduced [9], [10], [11], [12]. Based on these datasets, many scene graph generation (SGG) methods are proposed, and these methods can be divided into SGG methods with facts alone (such as the methods based in pixel-level object detection [13], GAN [14], and so on) as well as introducing prior information. At present, these SGG methods pay more attention to the methods with fact alone, including CRF-based (conditional random field) SGG [1], [15], [16], TransE-based (visual translation embedding) SGG [17], [18], [19], CNN-based SGG [20], [21], [22], RNN/LSTM-based SGG [23], [24], [25], GNN [26], [27], [28], and other SGG methods with fact alone [14], [29], [30]. In addition, different types of prior information are introduced for SGG, such as Language Priors [9], visual contextual information [31], [24], Knowledge priors [32], [33], visual cue [34], and so on. Scene graph has the powerful representations for the semantic features of a scene, thus, it has widely applied to related visual tasks, such as image retrieval [1], [35], image generation [36], [37], image/video captioning [38], [39], [40], Visual question answering (VQA) [41], [42], Human object interaction (HOIs) [43], [44], [45], Image understanding and reasoning [46], [47], [48], 3D scene graph [2], [49], and so on. Therefore, we can see that scene graph has become a hot research topic in computer vision, and it will still receive

- P. Xu, L. Guo and X. Chen are with the School of Information Science & Technology, Northwest University.
- X. Chang is with the Faculty of Information Technology, Monash University. Email: cxj273@gmail.com.
- P. Huang and A. Hauptmann is with School of Computer Science, Carnegie Mellon University. Email: alex@cs.cmu.edu.

continuous attention in the future.

Since the concept of Scene graph was proposed in 2015 and first applied to image retrieval, then the relevant researches on Scene graph have increased significantly, especially in 2019 (As shown in figure 2). In these research results, we mainly focus on the scene graph generation (SGG) methods and the applications of scene graph. Fig.3 (a) shows the relevant works on SGG, and it can be seen that more researches are focused on SGG by using GNN models and introducing relevant prior information. While the applications of Scene graph mainly refer to image generation, image/video captioning and image semantic understanding and reasoning, etc. as shown in Fig.3 (b). There also are a few applications on VQA and image retrieval. In addition, several works utilized 3D scene graph for 3D object detection and recognition. With the increasing researches on scene graph, the scene graph databases related to specific tasks are constantly updated and established, which enable reliable data for the further researches on scene graph in the future.

At present, the researches on scene graph mainly try to solve the following three problems: (1) How to generate a more accurate and complete scene graph; (2) How to simplify the computational complexity of SGG; (3) How to apply scene graph to more tasks in a more appropriate and extensive way. Although there have been many related methods proposed for solving these problems, there still need deep researches on the solution of these problems. Moreover, there are still other problems that need to be further solved. For example, the unbiased scene graph data has always been a problem in scene graph generation, and will be a problem to be solved in the later research. In addition, the descriptions of the relationships between objects in datasets are rough and inaccurate. Therefore, we need to further optimize the annotations in related scene graph datasets.

In this paper, we mainly discuss the generation and application of scene graph relevant to computer vision in this paper. In section 2, we first introduce several existing datasets that are commonly used for scene graph, as well as the performance evaluation of scene graph generation models. Section 3 briefly introduces basic notations of scene graph, and then provide a thorough review of current available scene graph generation techniques, including those work with facts alone, as well as using different types of prior information. Meanwhile, We describe the overall frameworks of models, model training, as well as pros and cons of such techniques. In section 4, we further explores the applications of scene graph to a wide variety of computer vision tasks. Furthermore, Section 4 f Section 5 will discuss the main problems in the generation and application of scene graph at present and the future researches of scene graph. Finally, we present our concluding remarks in Section 6.

2 DATASETS FOR SCENE GRAPHS

A long-standing goal of computer vision is to develop models that can understand the visual information in scenes, and further reason some unseen visual events from the current scenes. While in terms of current AI technologies, the performance of the relevant network models is still largely dependent on the knowledge learned from the existing datasets. If these models are transferred from their original datasets to other datasets with relatively unfamiliar scenes, the performance of the models is likely to decline dramatically or even fail to work. Therefore, large scale visual datasets for specific tasks are critical to the computer

vision network models. In this section, We discuss several existing datasets that have been released for scene graph generation and applications of relevant downstream tasks. We briefly state the basic data structure of these main scene graph datasets, and make a further comparative analysis on these data sets.

Real-World Scene Graphs Dataset. In 2015, Johnson proposed the notion of scene graph, as well as Real-World Scene Graphs Dataset (RW-SGD) [1], which may be the first dataset explicitly created for scene graph generation and application (image retrieval). RW-SGD is built by manually selecting 5,000 images from YFCC100m [50] and Microsoft COCO datasets [51], and then Amazon’s Mechanical Turk (AMT) is used to produce a human-generated scene graph from these selected images. The Final RW-SGD contains over 93,832 object instances, 110,021 attribute instances, and 112,707 relationship instances.

Visual Relationship Dataset (VRD) [9] is constructed for the task of visual relationship prediction. VRD has 100 object classes detected from 5000 images, and also contains 37,993 relationships. However, the distribution of the visual relationships has the common problem of the long tail of infrequent relationships in scene graph datasets.

Visual Genome Dataset (VGD) [10] is a large scale visual dataset, and consists the components of objects, attributes, relationships, question answer pairs, and so on. At present, VGD has widely applied to scene graph generation and application for its large number of images, objects, relationships, and so on. In addition, another scene graph dataset (Visually-Relevant Relationships Dataset (VrR-VG)) [21] is constructed based on VGD.

UnRel Dataset (UnRel-D) [11] is a new challenging dataset of unusual relations, and contains more than 1000 images, which can be queried with 76 triplet queries.

HCVRD Dataset [12] has 52,855 images with 1,824 object categories and 927 predicates, and also contains 28,323 relationships types. Similar to VRD, HCVRD also has the long-tail distribution of infrequent relationships.

3 SCENE GRAPHS GENERATION

The concept of scene graph is first proposed by Johnson in [1], and manually established the corresponding scene graph on a real-time World scene Graph dataset. A scene graph is a topological representation of a scene, which mainly encodes object and their relationships. The task of scene graph generation (SGG) is to construct a graph structure that best associates its nodes and edges with the objects and their relationships in a scene. While the key challenge task is to detect/recognize the relationships of the objects.

Currently, there are two main scene graph generation approaches [27]. The first approach has the two stages, that is object detection and pair-wise relationship recognition [15], [54], [9], [55]. The other approach is to jointly detect and recognize the objects and their relationships [22], [26], [53]. The subsequent SGG methods are proposed to generate a complete scene graph with facts alone or by introducing additional prior information. In this section, we will review the SGG methods using only facts observed in the given images/videos; and further discusses the techniques that incorporate other priors.

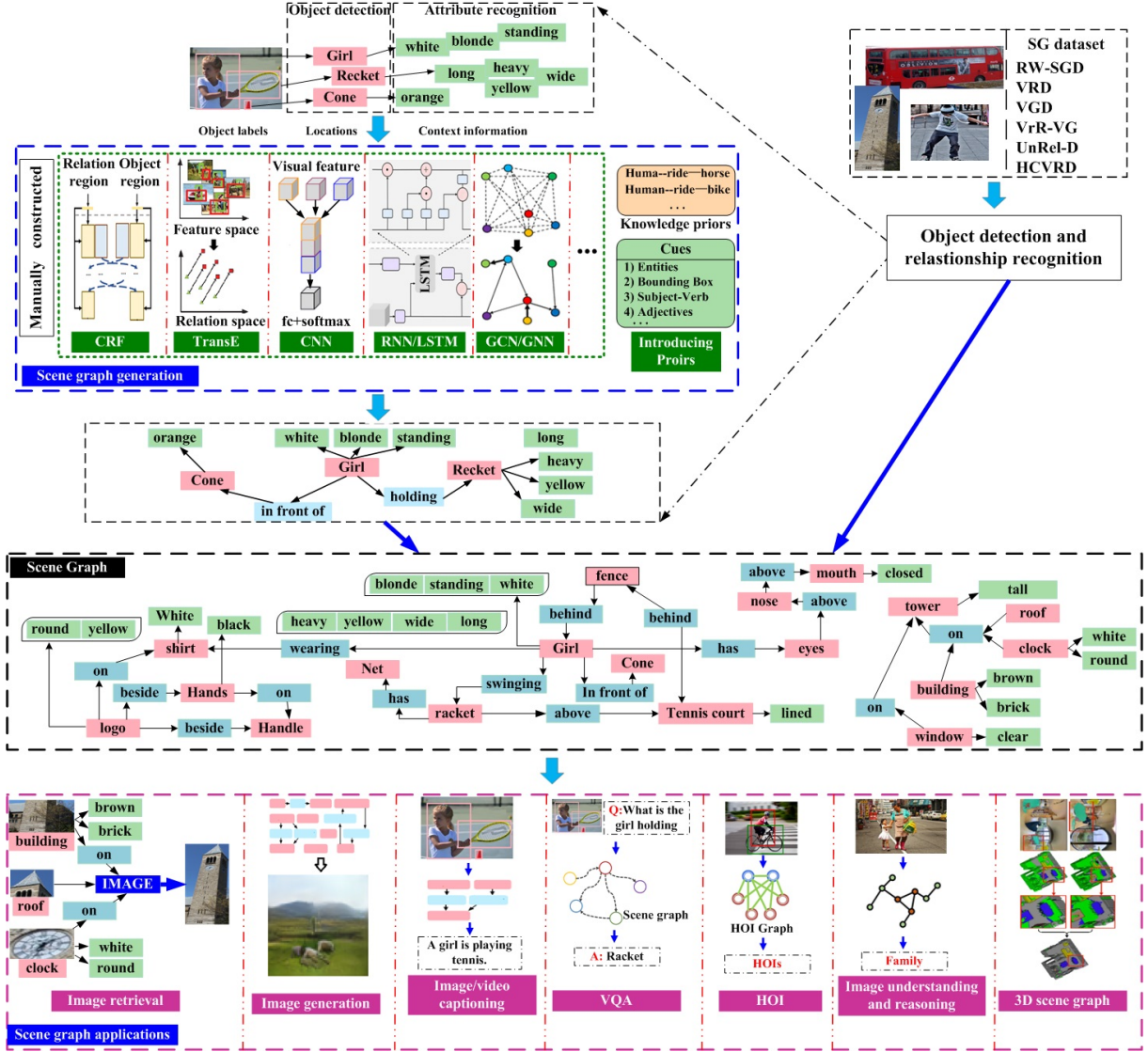


Fig. 1: (a): An example of a scene graph encodes objects, attributes, and relationships in a scene. (b): A grounding of the scene graph associates the object to the regions in an image.

Dataset	images/videos	Obj. instances	Obj. classes	Att. instances	Att. types	Rel. Instances	Rel. types	Pre. per Obj. Category	Pre.
COCO [51]	124,828	886,284	80	-	-	-	-	-	-
YFCC100m [50]	-845735	534,309	200	-	-	-	-	-	-
RW-SGD[1]	5000	93,832	6745	110,021	3743	112,707	1310	3.3	-
VRD [9]	5000	-	100	-	-	37,993	6,672	24.25	-
VGD [10]	100k	33,877	3,843,636	-	-	-	40,480	-	-
UnRel [11]	1000	-	-	-	-	76	-	-	-
HCVRD [12]	52,855	-	1824	-	-	256,550	28,323	10.63	927
VrR-VG [21]	58,983	282,460	1600	-	-	203,375	117	-	-
Visual Phrase [52]	2,769	3,271	8	-	-	2040	13	120	-
VG150[53]	87,670	738,945	150	-	-	413,269	50	-	-

TABLE 1: Aggregate statistics for scene graph datasets.

3.1 Scene Graphs Generation with Facts Alone

3.1.1 CRF-based SGG

Johnson et al. proposed the concept of scene graph, and give the corresponding formulations. While they used Amazon’s Mechanical Turk (AMT) to construct a scene graph manually on RW-SGD [1]. Furthermore, conditional random field (CRF) is construct

for image retrieval using the generated scene graph. However, it takes much cost for generating a scene graph manually, and it has the influence of subjective factors of understanding a scene. subsequently, Schuster et al. [8] proposed a method of scene graph generation automatically using two parsers: a rule-based parser and a classifier-based parser, which map dependency syntax

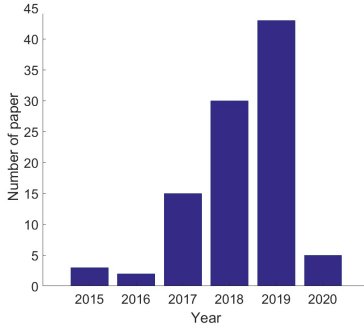


Fig. 2: Classification and statistics of the researches on scene graph from 2015 to 2020.

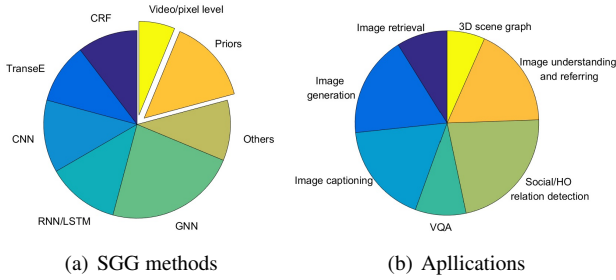


Fig. 3: The Classification and statistics of SGG methods and Applications

representations to scene graphs. Based on the constructed scene graph, they also achieved the image retrieval task via CRF. These may be the two early methods that involved the construction and applications of scene graph.

Formally, for a given scene graph $G = (O, E)$, there are many possible ways of grounding the scene graph to an image I . The inference tasks and other visual tasks at the high level are to recognize objects, predict the objects' coordinates, and detect/recognize pairwise relationship predicates between objects [56]. Therefore, the first stage of identifying the categories and attributes of the detected objects is achieved mainly using RPN or Faster RCNN [57]. Furthermore, most of the works focus on the key challenge of reasoning the visual relationship. In [15], the CRF model of scene graph also has two unary potentials that associate the objects with their visual features and relationships. While, for relational modeling, Deep Relational Network (DR-Net) are explored to detect the relationships.

In [16], SG-CRF is proposed to improve the accuracy of SGG. In this methods, Semantic Compatibility Network (SCN) are used to learn the semantic compatibility of nodes in the scene graph, and approximates scene graph inference by mean-field approximation algorithm, which can be expressed as $Q^t = \text{MeanField}(\psi_u, L_e, Q^{t-1})$. Then the pairwise potential ψ_p of each node is calculated according to the label word embeddings of its 1-hop neighbors. Q^T is the final output of last mean-field iteration. Assume that I is the given input image, and SG denotes the final generated scene graph. Then the objective for SG-CRF can be formulated as maximizing the following probability

function [16]:

$$P(SG|I) = \prod_{o_i \in O} P(o_i, o_i^{bbox}|I) \prod_{r_{i \rightarrow j} \in E} P(r_{i \rightarrow j}|I) \quad (1)$$

Where, the term $P(o_i, o_i^{bbox}|I)$ is a unary potential, which models the agreements of the visual features of the box o_i^{bbox} with the category and attributes of the object o_i . Then, CRFs for SGG can be formulated to find the optimal solution of $x^* = \arg \max_x P(X)$, which obeys the Gibbs distribution. Where,

$$P(X) = \frac{1}{Z(X)} \exp(-\sum_i \psi_u(x_i) - \sum_{j \neq i} \psi_p(x_i, x_j)) \quad (2)$$

Similar to Eq.1, the unary potential $\psi_u(x_i)$ is the measurement for assigning the node x_i , and the pairwise potential $\psi_p(x_i, x_j)$ is the cost of assigning x_i to x_j .

3.1.2 TransE-based SGG

There are similarities between scene graph and knowledge graph in terms of object relationship reasoning. Therefore, inspired by the advances of Translation Embedding networks (TransE) in relational representation learning of knowledge bases and object detection networks, relevant models and methods based on TransE are explored for visual relationship detection/recognition to build a scene graph. [17], [18], [19], [56]. These TransE-based SGG methods place the visual relationships between objects in a low-dimensional relation space, where the relations are modeled as simple translation vectors.

VTransE [17] extends TransE networks [58] for modeling visual relations and the predicates. In VTransE, the detected subjects and objects are mapped into a low-dimensional relation space, and their relationships are formulated as translation vectors for scene graph generation. Similar to other SGG methods, object detection needs to carry out first, and VTransE networks can be married to any object detection networks such as Faster-RCNN [57], SSD [59] and YOLO [60], which are used to locate the objects and recognize their categories for the following task of relation recognition.

VTransE represents any valid relation $\langle \text{subject}, \text{predicate}, \text{object} \rangle$ in vectors s , p and o respectively. If the relation holds, then this relation can be represented as a translation: $s + p \approx o$ in the embedding space, otherwise $s + p \not\approx o$. Besides VTransE learns the relation translation vector $t_p \in R^r$ as in TransE, and it learns two projection matrices W_s, W_o by $s = W_s x_s$ and $o = W_o x_o$:

$$W_s x_s + t_p \approx W_o x_o \quad (3)$$

Where x_s, x_o are the visual features of subjects and objects, respectively. Furthermore, a prediction loss is proposed to solve the problem of problematic sampling negative triplets due to the incomplete relation annotation:

$$l_{rel} = \sum_{(s,p,o) \in R} -\log \text{softmax}(t_p^T (W_o x_o) - W_s x_s) \quad (4)$$

Finally, the relation detection scores are obtained by summing the scores of subject/object detection and relation predicate prediction.

UVTransE [18] is proposed to improve generalization for the rare or unseen relations based on VTransE. There are lots of obvious object relations in scenes, but also exist many unseen relations. Therefore, the relation detection models also need to recognize the hidden relations. Inspired by VTransE [17], UVTransE introduces the union of subject and object, and a context-augmented translation embedding model is proposed to capture both common and rare relations in scenes. Similar to [17], UVTransE needs to learn three projection matrices W_s, W_o and W_u by minimizing the following multi-class cross-entropy loss function:

$$L_{vis} = \sum_{(s,p,o) \in T} -\log \frac{\exp(p^\top \hat{p})}{\sum_{q \in P} \exp(q^\top \hat{p})} + C(\|W_s s\|_2^2 - 1)_+ + \|W_o o\|_2^2 - 1)_+ + \|W_u u\|_2^2 - 1)_+ \quad (5)$$

Where, T and P are the set of all relationship triplets and the set of all predicate labels. $\hat{p} = W_u u - W_s s - W_o o$, $[x]_+ = \max(0, x)$. C is a hyper-parameter, which is used to determine the importance of the soft constraints. Eq.(5) is different from VTransE [17] in terms of the introduced contextual union feature. Finally, similarly to [17], the score of the entire triplet is also to use the sum of the subject/object detection score and the predicate recognition score.

MATransE (Multimodal Attentional Translation Embeddings) [19] is proposed to satisfy $s + p = o$ by guiding the features' projection with attention and Deep Supervision. Similar to [17], MATransE needs to learn the projection matrices W_s, W_p , and W_o by employing a Spatio-Linguistic Attention module (SLA-M) [15]. As shown in Eq.(3), a two-branch architecture is designed in MATransE: one branch is to drive the predicate features into scores $t_p = W_p x_p$ (P-branch), and another branch is used to classify the object-subject features $W_o x_o - W_s x_s$ (OS-branch).

Finally, both the scores of P-branch and OS-branches are connected into a single vector, which would be used to train a meta-classifier to predict the categories of the predicates. Thus, with $W = (W_s, W_p, W_o)$, the total loss can be formulated as:

$$L(W) = \lambda_f L_f(W) + \lambda_p L_p(W_p) + \lambda_{os} L_{os}(W_o, W_s) \quad (6)$$

Where λ is used to balance the importance of each term.

RLSV [61] is proposed to solve the problem of the incomplete scene graph, and the formulation of RLSV is to predict the missing relations between the objects. RLSV is staged by three modules: visual feature extraction, hierarchical projection and train objective module. By combining location and visual information of entities, the visual feature extraction model embeds the inputting image as visual projection vectors $v_{p_h}, v_{p_r}, v_{p_t}$ for head h , relation r and tail t respectively. Based on $v_{p_h}, v_{p_r}, v_{p_t}$, the hierarchical projection module projects a given visual triple (h, r, t) onto attribute space, relation space and visual space, resulting in a new presentation $(h_\perp, r_\perp, t_\perp)$. Then followed by TransE, the score function can be defined as:

$$E_I(h, r, t) = \|h_\perp + r_\perp - t_\perp\|_{L_1/L_2} \quad (7)$$

Finally, a max-margin function with negative sampling is formulated as the training objective:

$$L = \sum_{I \in \mathcal{I}} \sum_{(h,r,t) \in \mathcal{T}_I} \sum_{(h',r',t') \in \mathcal{T}'_I} [E_I(h, r, t) - E_I(h', r', t') + \gamma]_+ \quad (8)$$

where γ is a marginal hyperparameter, \mathcal{T}'_I is the negative sampled visual triple set generated from positive visual triple set \mathcal{T}_I .

3.1.3 CNN-based SGG

DR-Net [15] is a framework which formulates a triplet in the form of $\langle \text{subject}, \text{predicate}, \text{object} \rangle$ as its prediction output, and jointly predicts the category label of the triplet by exploiting the spatial configuration and statistical dependency among the elements of the triplet. The whole framework has three stages: object detection, filtering pairs of objects and joint recognition. In the object detection stage, Fast RCNN is used to detect a set of candidate objects in the images, and each detected object comes with a bounding box and its visual features. The next stage is to filter out some pairs of detected objects by a low-cost neural network. Then the retained pairs of objects are fed to the joint recognition module. Finally, The joint recognition module would produce a triplet as the output by considering the visual features of each object, the spatial configurations between any paired objects and the statistical dependency between the relationship predicates. Specifically, to represent the spatial configurations, dual spatial masks are designed by deriving from the bounding boxes, which may overlap with each other. To exploit the statistical relations, DR-Net is developed to incorporate statistical relational modeling into a deep neural network framework.

SIN (Structure Inference Network) [62] is a detector which is designed to infer the object category labels by improving Faster R-CNN with a graphical model. The visual features of the objects, the relationships between objects in a single image and the scene contextual information are exploited in SIN for improving the performances of object detection. The framework of SIN is as follows. ROIs are derived from the input images, and then each ROI is pooled into a feature map with fixed-size f_i^v , which is considered as a node in graph modeling. Meanwhile, a scene of an image is generated from its global feature f^s in the same way. The scenes and nodes are put into the SIN as Scene GRUs. Afterwards, both the spatial and visual features of the nodes v_i and v_j are jointly combined to form a directed edge $e_{j \rightarrow i}$ from v_j to v_i , which represents the influence of v_j on v_i . All edges will be passed into SIN as Edge GRUs. In SIN, the state of each GRU is updated at each iterative, and the final integrated node representations are used to predict the object category and bounding box offsets.

Rel-PN [20]. Relationship Proposal Networks (Rel-PN) first detect all meaningful proposals of object, subject and relationship by running 3-branch RPN in Faster RCNN [57] respectively. Although the object instances and subject instances belong to the same category space, their distribution is inconsistent, so they are extracted separately.

The relationship branch is to reduce the number of the pairs of objects, otherwise there would have $object \times subject$ pairs of relations. In [20], 9 kinds of relationship proposals are selected according to several conditions. Then two branches of visual compatibility and spatial compatibility modules are used to output the visual and spatial scores. For visual compatibility module, three visual features are connected to obtain a (5x5x512) vector, and then the module output visual score s_v . Moreover, three groups of spatial difference features are connected to obtain a (64x64) vector, and then output spatial score s_s . Finally, p_v and p_s are integrated into a final score.

$$p = \alpha p_v + (1 - \alpha) p_s \quad (9)$$

where, α is the ratio of visual compatibility.

Based on the model in [20], the model in [63] considers three types of features: visual, spatial and semantic features using three corresponding models, and these features are then fused for the final relationship identification. Different from [20], the model in [63] used an additional semantic module to learn the semantic features, and achieved better performances.

ViP-CNN [22] has the capacity to jointly learn the specific visual features for the interaction and to consider the visual dependency. ViP-CNN has four branches for triplet proposal and phrase recognition. Likely to other SGG models, Faster R-CNN with VGG-Net [64] as backbone is used to detect the objects and locate the corresponding bounding boxes, so as to provide the triplet proposals. For the triplet proposal branches, the extracted CNN features by VGG-Net are used for proposing regions of interest (ROIs), and then triplet proposals are obtained by grouping these ROIs. Furthermore, triplet non-maximum suppression (triplet NMS) is proposed to solve the problem of the sparsity of relationship annotations, so as to reduce the redundancy information. While, the remained triplets are used for the branch of phrase recognition.

BAR-Net[65] uses the standard object detection methods to detect pair-wise relationships, which is achieved by decomposing the relation detection task into two tasks of retentive object detection. In BAR-Net, one detector (Such as faster RCNN) is used to detect all objects in the image, and then the other detector was used to detect the objects, which have interactions with each object. The bounding boxes obtained by the first detector are used as the inputs for the second detector, and the the joint probability can be represent by simpler conditional probabilities:

$$pro(s, p, o|I) = pro(s|I)pro(p, o|s, I) \quad (10)$$

The second probability item $pro(p, o|s, I)$ models the probability that an object o in the image is related to the subject S , which is called Box Attention.

LinkNet [56] is proposed to improve scene graph generation by explicitly modeling inter-dependency among all related objects, rather than an object in isolation. Linknet mainly has three modules: 1). A relational embedding module is used to classify the objects and their relationships. Given an image, objects' proposals and labels are extracted by a object detection method, such as Faster R-CNN [57]. 2). A global context encoding module is used to extract global information, which contains as much as possible all proposal information in the image, and is used to assist the classification of object relations. 3). A geometrical layout encoding module is used to assist in the classification of object relations using the spatial information between the object proposals. Finally, the two categories can be used to generate the scene graph, and the loss function of whole network is the weighted sum of the losses for predicting the bounding boxes and categories of the detected objects, and even the relationship categories between the objects.

3.1.4 RNN/LSTM-based SGG

Iterative Message Passing [53]. As many previous works focused on building a scene graph given an image, surrounding context in it is ignored. However, scene graph prediction based on context information could resolve vagueness since local predictions is isolated. Motivated by this observation, Xu et.al. proposed an iterative message passing based model to fulfill scene graphs generation.

Given an image, their model first generates object proposals B_I by Region Proposal Network (RPN). With these object proposals, three tasks is to be fulfilled: object class label inferring, bounding box offsets computation and predicate prediction. Thus, the scene graph generation problem is formulated as optimizing

$$x^* = \operatorname{argmax}_x Pr(x|B_I, I) \quad (11)$$

,where

$$Pr(x|B_I, I) = \prod_{i \in V} \prod_{j \neq i} Pr(x_i^{cls}, x_i^{bbox}, x_{i \rightarrow j} | B_I, I), \quad (12)$$

x_i^{cls} is the i^{th} object proposal, x_i^{bbox} is the i^{th} proposal's bounding box offsets, $x_{i \rightarrow j}$ is the predicate that i^{th} object proposal applies to j^{th} object proposal.

To cut down the cost of inference on a dense graph, a general RNN unit, namely GRU, is used to seek out the hidden states in the training framework. Based on GRU, $Pr(x|B_I, I)$ is transformed as follows.

$$Pr(x|B_I, I) = \prod_{i \in V} Q(x_i^{cls}, x_i^{bbox} | h_i) Q(h_i | f_i^v) \prod_{j \neq i} Q(x_{i \rightarrow j} | h_{i \rightarrow j}) Q(h_{i \rightarrow j} | f_{i \rightarrow j}^e) \quad (13)$$

where f_i^v is the i^{th} feature node and $f_{i \rightarrow j}^e$ is the edge feature from node i to node j . To prompt inferring efficiency, a prime dual message passing scheme between node GRU graph and edge GRU graph is used with the bipartition of a scene graph.

PANet [23]. Many previous works focus on context and scene information for relationship prediction, which ignores internal connection among predicates. Therefore, this paper proposed Predicate Association Network (PANet) to effectively acquiring contexts and relationship between predicates.

PANet is a two-stage network. In the first stage, Faster-RCNN is used to generate object proposals, of which each b_i represents three kinds of object features: category embedding (E_{b_i}), spatial information (S_{b_i}) and visual feature (F_{b_i}).

$$V_{b_i} = \sigma(W_b(E_{b_i} \circ F_{b_i} \circ S_{b_i}) + b_b) \quad (14)$$

Based on these object proposals, instance-level context and scene-level context are obtained via an RNN. For each object pair $\langle s_i, o_i \rangle$, their class probability $P(s_i|I)$ and $P(o_i|I)$ are acquired by these contexts.

In the second stage, the connections of predicates are found by another RNN, where matching and attention principle are used. Each predicate p_i is represented by a word embedding E_{p_i} . The combination of instance-level and scene-level contexts of each object pair $\langle s, o \rangle$ is denoted as $\langle G_s, G_o \rangle$. Feature maps of their union bounding box $F_{s,o}$ are used to represent the state of the corresponding union region. $F_{s,o}$ is then threw into a complete layer to reduce dimension. The converged feature vector $U_{s,o}$ of two objects is:

$$U_{s,o} = (G_s * G_o) \circ \sigma(W_u F_{s,o} + b_u) \quad (15)$$

where $G_s * G_o$ is used to find out each object pair's contexts.

Alignment feature is obtained by refining predicate label E_{p_i} with $P_{s,o}$. Then these features R_{p_i} are feed into an RNN module to extract predicate sets $\gamma_{p_i}^{(2)}$ of the i th predicate:

$$\gamma_{p_i}^{(2)}, h_{p_i}^{(2)} = RNN(R_{p_i}, h_{p_i-1}^{(2)}) \quad (16)$$

where, $h_{pi}^{(2)}$ is the hidden state in step i of the RNN, and $\gamma_{pi}^{(2)}$ is contextual information of predicate p_i . Then the final weighted contexts γ_{att} are computed for each pair $\langle s, o \rangle$ as:

$$\gamma_{att} = \sum_{i=1}^m w_i \gamma_{pi}^{(2)} \text{ s.t. } 0 \leq w_i \leq 1 \quad (17)$$

The predicate label is determined as with the highest probability:

$$P(p_i | I, s_j, o_j) = \max(W_r \gamma_{att} + b_r) \quad (18)$$

where W_r and b_r are weights and bias respectively.

CMNs [66]. Two issues exist in previous works on referential expressions, where one is that referential expressions were treated uniformly, thus failing to uncover the consistence between textual components and visual units in an image, the other is that relationship categories of most previous works are fixed.

To solve these two problem, this paper focus on referential expressions via internal relationships, which are denoted as triplets $\langle \text{subject}, \text{relationship}, \text{object} \rangle$ and proposed Compositional Modular Networks (CMNs) [66]. It is an end-to-end architecture that models the language structure of referential expressions and their groundings and supports parsing of any language.

CMNs decomposes a referential expression into a triplet $\langle \text{subject}, \text{relationship}, \text{object} \rangle$ by three distinguished attention maps and arranges the extracted textual notations based on image regions by a modular neural framework. Two modules work in CMNs, one for seeking textual components by generating scores for each component, the other for addressing the relationship between two bounding-box pairs by pairwise score of region-region pairs. LSTM here is used for expression parsing with attention in CMNs.

VCTREE [25]. Prior layout structures, like chain, fully connected graphs, are reliable for visual context encoding. Such prior layout structures are not perfect for the following two reasons. First, linear structures are too simple and might only obtain some spatial information or co-occurrence bias, and complete graphs are not discriminative for hierarchical relations since dense connections could result in saturated message passing during a sub-sequential context encoding. Second, object layouts should vary as contents or questions change. Therefore, fixed chains and complete graphs are inadequate for dynamic visual contexts.

In this paper, a model named VCTREE, which imposes dynamic "trees" on encoding object-level visual contexts for visual inferring tasks, is proposed. VCTREE model has the following four steps [25]. 1) Faster-RCNN is implemented for object proposal detection. The visual feature of proposal i is denoted as x_i , combining a RoIAlign feature $v_i \in R^{2048}$ with a spatial feature $b_i \in R^8$ with 8 variables, the bounding box coordinates (x_1, y_1, x_2, y_2) , center $(\frac{x_1+x_2}{2}, \frac{y_1+y_2}{2})$ and size (x_2-x_1, y_2-y_1) . Segmented features, like instance segmentation or panoptic segmentation, could also be options since the visual feature x_i is not just bounding box. 2) A matrix is learned to construct VCTREE. Since the VCTREE construction could be separate and the score matrix is non-differentiable from the loss of end-task, a hybrid learning strategy is developed. 3) Bidirectional Tree LSTM is employed to encode the contextual information through VCTREE. 4) The encoded contexts will be decoded for each specific end-task.

AHRNN [67]. Most existing approaches have two limitations to generate a good scene graph effectively. One is that futile object

bounding boxes or unnecessary relationship pairs are generated by object detection based approaches. The other is that linguistic superfluous relationships would be build by probability ranking method. Motivated by these two observations, the authors proposed an architecture by first focusing on regions of interest and recognizing them without extra object detection and ranking the score of relationship triplets automatically based on features.

Its architecture is constituted by a CNN model, an Attention-based Hierarchical Triplet RNN (AHRNN) and a scene graph construction algorithm. The CNN model is of "encoder-decoder" framework, of which the goal is to extract a set of global visual feature vectors of an input image. AHRNN is proposed to automatically transform the above vectors into a set of triplets. AHRNN is made up of two parts, an Attention-based Triplet RNN (ATRN) to turn image features into a topic vector successively by roughly leveraging them to compose relationship triplets, and an Attention-based Word RNN (AWRNN) to recognize each target word in a $\langle \text{subject}, \text{predicate}, \text{object} \rangle$ triplet based on topic vector. Finally, an automatic scene graph construction algorithm is performed to finalize the components in a scene graph based on the predicted relationship triplets.

MOTIFNET [24]. Visual scenes have strong structural regularities. Based on this motivation, the authors extract many structural repetitions in scene graphs from Visual Genome dataset, which comes from COCO and contains millions of objects and relationships with annotation for 100k images. Their analysis reveals two important results. First, the object categories of objects might determine the distribution of the relationships, but not vice versa. This implies that strong regularities exist in a local sub-graph. Second, structured modules appear even in large sub-graphs and regular substructures in scene graphs appear repeatedly over half of images, which are defined as motifs.

Based on the above analysis, a baseline is induced: given object categories, forecast the most frequent relationship between object pairs. This baseline improves by 1.4% on the parameter of average recall, which implies that an effective scene graph generation framework should consider both the asymmetric connection between objects and their relationships and frequent appearing substructure in scene graphs.

Thereafter, a neural network architecture called Stacked Motif Network (MOTIFNET) is proposed, which contains three stages: bounding box detection, labels prediction for regions and relationship forecasting.

Global context is extracted by different bidirectional LSTMs and transmitted between each stage for different tasks. In the first stage, a detector results in bounding boxes and object context between bounding boxes is extracted. The global context in this stage is to label these bounding boxes. Given bounding boxes and labels, edge context is prepared to predict relationships. Finally, category labels are given to relationships through an outer product, which combines contextualized head, tail, and union bounding region information.

MSDN [26]. The authors explored understanding images by a single neural network model for three tasks: object detection, scene graph generation and image caption [26]. Since features for these tasks are of high correlation and complementary for each other, the authors proposed an end-to-end Multi-level Scene Description Network (MSDN) to detect objects, predict their relationships and caption images simultaneously, which sufficiently used semantic annotations of three levels and their connections between them.

The framework of MSDN is as follows. 1) Region proposal,

including proposals for objects, phases and region captions. Object region proposals are extracted via Region Proposal Network (RPN). Phase region proposals are generated as the full combinations between all the object region proposals. And caption region proposals are generated via another RPN based on grounding bounding boxes. 2) Feature specialization. Different feature is specialized for different task, which is realized by leveraging different FC layers for different branches. 3) Dynamic graph construction. The goal of this stage is to dynamically build a graph to reveal connections between feature nodes based on semantic and spatial relations among ROIs for different branches. 4) Feature refining. Three parallel feature refining modules are performed iteratively via message transmission along the graph for different tasks. 5) Final prediction. The above features are then used to label objects, forecast predicates and generate captions. The scene graph is generated based on objects and their relationships.

The key procedure of MSDN is dynamic graph construction where region features, phrase features and object features are extracted from the original features separately and delivered for image caption, phrase detection and object detection respectively after feature refining.

3.1.5 GNN-based SGG

Factorizable Net: An Efficient Subgraph-based Framework for Scene Graph Generation [27]. Most approaches to generate scene graphs follow two popular frameworks. One performs object detection first and then applies pair-wise relationships recognition, the other first generates object region proposals and then separately predicts object labels and their relationships based on these proposals [27]. Both frameworks will build a quadratical number of objects, which is time-consuming.

A different framework called Factorizable Net is proposed in [27]. The object region proposals are detected by RPN first and then they are paired to construct a complete directed graph. Thereafter, a preciser graph is generated by merging edges corresponding to similar union regions into a sub-graph. Based on the above objects and sub-graphs, 2-D feature maps for sub-graph as well as feature vectors for objects are both generated. These two kinds of features are then refined through Spatial-weighted Message Passing (SMP) structure. The refined features would be passed to Spatial-sensitive Relation Inference (SRI) module for predicate recognition. Here, SMP, a GNN approach, is used for better feature representation. A novel clustering method is used to decompose the entire graph into sub-graphs during inference, where each sub-graph has several objects and their relationships as edges. By substituting the original scene graph with these sub-graphs, the generation efficiency of scene graph is much promoted.

Predicate Prior Model. Similar to the idea that generating a scene graph by using language prior of relationship triples [9], Hwang et.al. generated it by jointly combining the prior of predicate distribution [68]. The framework of [68] is similar to that of [53]. The framework first extracts node context information from visual features and edge context information from the relation among bounding box offsets. Then, mean field is used to perform approximate inference by Gated Recurrent Units (GRU), which is an iterative message passing scheme to label objects, locate their bounding box offsets, and predict relationship predicates between object pairs. The difference between [68] and [53] is that a pre-trained tensor-based relational module was added as a dense relational prior in [68] to refine the relationship estimation

during the iterative message passing period, which is also a fine-tuning of the learning process of the scene graph module [53]. Here, an iterative message passing scheme with GRUs is used as a GNN way to promote the scene graph generation performance with better feature representation.

Graph R-CNN for Scene Graph Generation [69]. This work takes regularities of object relationships and proposes the Graph R-CNN framework. Graph R-CNN is made up of three stages: 1) detect object region via RPN, 2) trim relationship edge between the above object regions by Relation Proposal Network (RePN), 3) graph context integration by an attentional GCN. In the first stage, local object regions are obtained by a popular object detection module. Based on realistic regularities of object relations, a relation proposal network (RePN) effectively computes relation scores between object pairs to explore rare scene graph connections. Thereafter, an attentional graph convolution network (AGCN) is applied to transmit higher-level context information throughout the graph by iteratively updating objects and relationships. Different from existing work, node-edge attentions are predicted to regulate information flow between unreliable or unlikely edges.

A scene graph generator should reveal the connection between objects and relations in order to improve prediction precision [70]. Motivated by this observation, this paper demonstrated that the architecture of a neural network should stay invariant to a particular type of input permutation. Formally, given the same features, the same result is supposed to be obtained by a framework or a function \mathcal{F} even given a permutation of the input. For example, consider a category space with three variables y_1, y_2, y_3 , and assume that \mathcal{F} takes as input $z = (z_1, z_2, z_3, z_{12}, z_{13}, z_{23}) = (f_1, f_2, f_3, f_{12}, f_{13}, f_{23})$, and outputs a label vector $y = (y_1^*, y_2^*, y_3^*)$. Given a permuted input like $z' = (f_2, f_1, f_3, f_{21}, f_{23}, f_{13})$ to function \mathcal{F} , the output should still be $y = (y_1^*, y_2^*, y_3^*)$.

The authors proved this property based on the fact that such architecture or framework can gather information from the holistic graph in a permutation invariant way. Based on this property, they suggested several common architectural structures like attention neural networks and RNNs, which was used in their scene graph module.

It is a challenge to deal with the long-tailed distribution of relationships for most approaches since they are mostly trained on those predicates with sufficient labels. Motivated by this observation, the authors attempt to construct a scene graph by few-shot learning of predicates, which can scale to new predicates.

The pipeline of Few-Shot Learning is as follows. 1) Fully train Graph Convolution model and spatial and semantic shift functions on relationships with abundant data. 2) Define shift functions for new rare relationships with few examples using fully trained shift functions. 3) Fine-tune new shift functions with few training examples

The novelty of this model is that predicates are defined as functions such that object notations are useful for few-shot predicate forecasting, including a forward function that turns subject notations into objects and a corresponding function that changes the object representation back into subjects [71].

ARN [28] The Attentive Relational Network (ARN) is mainly proposed to address the following two problems. One is that most existing papers ignore the semantical relation between visual features and natural languages as well as internal connections among the triplets. The other is that most works leverage a

stepwise method to represent nodes and edges, which neglects the global structure and information of an image.

The Attentive Relational Network mainly has four parts: 1) a object detection module to extract visual features and relation information about entities with bounding boxes. A softmax function is used to score the initial classification of each entity and relation; 2) a semantic transformation module to produce semantical embeddings by transforming label embeddings and visual features into the same space; 3) a graph self-attention module to leverage a self-attention framework to embed entities; 4) a relation inference module to predict entity category and relationship as the final scene graph result.

Graphical Contrastive Losses [72]. Since a subject or an object relates to many instances of one category, most models fail to discriminate one target instance from another and obtain the exact pairs since an image contains similar subject-object interacting. These two obstacles result in two errors respectively, Entity Instance Confusion and Proximal Relationship Ambiguity. In this paper, a set of contrastive losses between graphs is proposed to tackle these two issues. The losses in form of offset based triplet loss are specifically constructed to solve the aforementioned errors. It also adds supervision in the form of negatives specific to Entity Instance Confusion and Proximal Relationship Ambiguity.

The proposed Relationship Detection Network (ReIDN) [72] has two stages. The first stage generates all bounding box regions. In the second stage, it derives three types of features for semantic, visual and spatial relationship proposal. Each feature outputs a set of class logits, which are combined via entity-wise addition and applied by a softmax normalization to obtain a distribution of predicate labels.

CMAT [73]. Existing SGG mechanisms do not effectively grasp the relation among visual context for the following reasons. objects and relationships in previous SGGs are independent from each other since the goal of training with cross-entropy is not graph-coherent. Besides, this goal should change with the update of even a single node.

Based on the aforementioned analysis, this paper proposes Counterfactual critic Multi-Agent Training (CMAT), which is a new framework to simultaneously fit the requirements of graph-coherent and local-sensitive. Its framework is as follows. It first leverages RPN to create some object regions according to an image. Thereafter, visual context is encoded based on the communication between these objects. Specifically, objects are considered as synergic agents and these agents relate to each other based on pairwise visual features to figure out their object labels. During communication, category scores are obtained for all objects. Based on the scores, object labels are assigned and visual relationship for each object pair is inferred. Thus, a entire scene graph is constructed. During training, a counterfactual critic is used to compute the individual contribution.

Differentiable Scene Graphs [74]. Different from most existing works, this paper attempts to figure out corresponding bounding boxes of subjects and objects given an image and a triplet query $\langle subject, relation, object \rangle$.

Leveraging scene graphs to solve this problem is a natural way but several challenges are still there. Scene graphs are not easy to be learned from a up-bottom task since they are discrete. Besides, even pre-training SG generators separately can be an alternative, this method is not feasible since numerous manual annotations are required and the coverage of pre-trained SG way is low.

This work designs Differentiable Scene-Graphs (DSG) to solve the above obstacles with following architecture. Thereafter, object features are extracted from the backbone based on these boxes. Parallely, each box pair is used to obtain a union box and features are derived from these box pairs in the same way. As inputs to a Differentiable Scene-Graph Generator Module, these features will be used to create the Differential Scene Graph (DSG), which is a set of new node and edge features. The DSG could refine the original boxes and be viewed as a Referring Relationships Classifier to identify a bounding box as either a Subject, an Object, Other or Background. Other means that the corresponding box is involved in another query relationship over this image. Otherwise the label will be Background. The novelty of this architecture exists that the components of a scene graph is decomposed so that each element in a triplet could be represented by a dense descriptor. Given an image and triplet query, bounding box proposals are created by a detector.

Triplet-Aware Scene Graph Embeddings [75]. This paper attempts to solve the layout prediction problem which forecasts scene layout masks and object localization according to a scene graph [36]. The authors propose the following layout prediction framework. Given a scene graph, the framework first uses a GCN to construct corresponding embeddings of object nodes. These embeddings are propagated to the layout prediction network to create a sequence of $\langle subject, predicate, object \rangle$ triplet embeddings. Besides individual class labels, this network classifies objects as either subject or object, inducing a new ordering and combination among objects. The triplet embeddings are also given to a triplet superbox regression network to jointly localize subject and object bounding boxes. Ultimately, all of the outputs would be merged to generate a scene layout mask with object localization. In this work, triplet embeddings with supervisory signals are used to improve scene layout prediction. Besides, data augmentation technique is performed to maximize triplet number during training. These two methods, namely, additional supervision and data augmentation, will enhance the embedding representation.

3.1.6 Other SGG methods with facts alone

SG-GAN [14]. Most of previous SGG methods use detectors to detect all the objects, and then generate the whole scene graph. Therefore, these methods have limitations of bounding boxes being available and without using the objects' attributes. The method first generates small sub-graphs, which can describe a specific region of the input image about a scene. Then, all of the generated sub-graphs are used to construct the complete scene graph. In this method, the images and noise information are first fed to a generator, then a CNN is used to extract the image features, and a dynamic image representation and attention vector are obtained using an attention mechanism. Finally, the image representations are used to produce triples by LSTM. Inspired by GAN, the triple generator is trained adversarially. While the trained triple Generator would resolve all the triples into a graph.

VRL [30] may be the first SGG method by using reinforcement learning [76]. This method is to gradually generate the scene graph, and the relationships between subjects and objects are generated in each step, so that the final complete scene graph will be gradually formed like a tree. For the whole model framework, the input states of reinforcement learning is parts of state features, including image features, subject features, object features and history phrase information. Then there are three branches of output

actions, which are to determine the properties of the subjects, the relationships between the current subjects and objects and the categories of the next objects. Variation-structured reinforcement learning actually refers to that the action space of the model varies according to the state in each step, so as to reduce the action selecting space and improve the accuracy. To this end, Directed Semantic Action Graph is constructed by the training set, which is actually the statistical information of relations and attributes in the data set relative to the object categories. Finally, three reward functions is defined to reflect the detection accuracy of taking action in a specific state.

CMAT [73]. To improve the quality of scene graph, the most important thing is to improve the performances of relationship recognition. Therefore, CMAT combines objects recognition and relation recognition to effectively improve the quality of scene graph, and each object in the images is regarded as an agent. The existing algorithms use the cross entropy as the loss function of object detection and recognition, but there is a problem that the importance of each object is different. To this end, graph-level metrics (such as Recall @k [9] and SPICE [77]) are used to evaluate the detection results, and used as a supervisory signal for model training. Then, The final multi-agent policy-gradient is used to maximize the graph-level metrics.

Analogies Transfer [78]. During generating the scene graph, there are many unseen relationships of the individual entities in the dataset. In order to generate a complete scene graph, Peyre et al. proposed to use analogy transformations to detect the unseen relationships that involve similar objects for the model. The whole network model has two stages. In the first stage, all the subjects and objects are detected, and the module of visual phrase embedding is used to learn the features of subjects, objects, predicates and visual phrases by optimizing the joint loss $L_{joint} = L_s + L_o + L_p + L_{vp}$. Then if we need to identify a unseen triplet, the model can utilize analogy transformation to compute the similarity between the unseen triplet and its similar triplets to estimate this unseen relationship.

GB-NET [5]. Due to a unified formulation of the two constructs of Knowledge Graph and Scene Graph, Graph Bridging Network (GB-NET) is proposed to incorporate the combination of the rich visual and commonsense information. In GB-NET, the scene graph and entity bridges are first initialized using Faster R-CNN. Then a variant of GGNN [79] is used to propagate the messages throughout the graph to update node representations, which establishes the bridge between the instance-level visual knowledge and commonsense knowledge, and a scene graph would be generated.

In addition, A simple and effective SSG method was proposed in [80] by jointly embedding the images and scene graphs. This method try to generate a scene graph from images by investigating several existing methods based on bag-of-words, sub-path representations and GNNs.

3.2 SGG by introducing additional information

To generate a scene graph faster and more accurately, scene graph generation models pay more attention to introducing multiple types of prior information, such as language priors, visual priors, knowledge priors, contexts, and so on. In this section, we discuss the related works of SGG by introducing additional information.

Phrase Cues [34]. Plummer et al. proposed a model framework for localizing or grounding the phrases in the images by

introducing linguistic and visual cues, which are constructed from the captions. Then the single phrase cues (SPCs) and the phrase pair cues (PPCs) are used to combine with Canonical Correlation Analysis (CCA) [81] to detect visual relationships. Therefore, in [34], the introduced priors are a list of the cues with corresponding phrases from the sentence, and these cues are extracted from the captions.

Language Prior Model [9]. Given the relationship annotations between the objects, which are detected in a set of fully supervised images. Lu et.al. proposed to train a visual appearance module and a language module individually, and later these two modules are combined together through a objective function, so as to improve the final performance to generate a scene graph. Compared to the method of Visual Phrase, in which a separate detector is designed for each single relationship, Language Prior Model uses the visual appearance module to learn the individual visual features of its comprising objects and predicates. In computational complexity, for the N objects and K predicates, Visual Phrases are used to train $O(N^2K)$ unique detectors, while only $O(N + K)$ detectors need through visual appearance module in Language Prior Model. In addition, the language module in Language Prior Model is very novel to project relationships into a word embedding space for SGG. In language module, the similar relationships are optimized to be more close together based on the semantic priors of relationships. In this way, the rare relationships can also be predicted so as to solve the problem of the long tail of infrequent relationships to a certain extent.

LK Distillation [55]. In most previous SGG methods, the visual relationship between two entities are generated. While in [55], Yu et al. try to model the three entities in a scene jointly, which can more accurately reflect these entities' relationships compared to modeling them independently. However, to reduce the complexities of model learning, the prior knowledge of linguistic statistics is used to regularize the visual feature learning. The useful linguistic knowledge can be extracted from the training relation annotations (internal knowledge) and other publicly information, such as Wikipedia. While in the teacher-student knowledge distillation framework [82], the distilled linguistic knowledge is used to predict the predicates using the visual features.

CDDN [31]. Cui et al. proposed a context-dependent diffusion network (CDDN) framework to identify the visual relationships. Before carrying out CDDN, object detectors are used to obtain the locations, labels and confidence scores of all the detected objects, which would be used as the input for CDDN. Then two different types of global context information (semantic priors and spatial scenes) are used for visual relationship detection. Semantic priors are learned by a word semantic graph from language priors, and spatial scenes are obtained by a visual scene graph to extract the visual features. Then these two types of global context information are adaptively aggregated by a diffusion network to estimated the predicates.

CISC [83] is another SGG method by introducing the context information. Besides significative visual pattern is als be explored for SGG. In Relationship Context-InterSeCtion Region (CISC) method, the context for relationships is constructed to benefit the relationship recognition from their association, and the proposed intersection region are used to discover the effective visual pattern for relationship recognition.

Knowledge-embedded routing network. In the real world, the distribution of the relationships is unbalanced, which leads to the poor performance of the existing methods in recognizing

the relationships with the low frequency. To solve this problem, the SGG model based on knowledge-embedded routing network is proposed by Chen et al. in [32]. In this method, a series of object regions are generated using Faster RCNN. Then, a graph network is used to propagate the features of nodes on the graph to learn the more contextualized features, so as to predict the labels in each object pair. Moreover, another graph is used to correlate the pairs of objects, and their relationships are predicted by a GNN model. The same process is repeated for all the pairs of objects to recognize their relationships, and the final scene graph is generated. Therefore, the statistical correlations between pairs of objects and their given relationships are used as the introduced priors for SGG in Knowledge-embedded routing network.

KB-GAN [33]. Since the existing scene graph datasets have the problem of the long tail in the distribution of object and relationship labels. Commonsense knowledge extracted from the external knowledge bases (KB) is used to refine object and phrase features for SGG, and an auxiliary image reconstruction path based on GAN is introduced to regularize the whole SGG network (KB-GAN). Therefore, in fact KB-GAN is also an application of scene graph on image generation.

3.3 videos and pixels-level for SGG

SGFB. In [84], a new data structure: Action Genome is introduced as a representation of spatio-temporal scene graphs. To generate the spatio-temporal scene graphs, Scene Graph Feature Banks (SGFB) is proposed, and the spatio-temporal scene graphs are further incorporated into a sequence of scene graph features as the final representation $F_{SG} = [f_1, f_2, \dots, f_T]$, which is used to predict action labels by 3D CNNs. With Action Genome, the action recognition task has achieved better performance on the Charades dataset.

Ontology graph is proposed in [3] to describe objects, parts, actions and attributes in a scene. Ontology graph has several similarities with scene graph, for example, these two types of graph structures have objects, attributes and relationships, and both of them also have their sub-graphs. In [3], ontology graph is used for scene-centric joint-parsing of cross-view videos, and the tasks of object detection, multi-object tracking, action recognition and human attributes recognition are used to evaluate the proposed scene-centric joint-parsing framework.

Pixels2Graph [13]. the existing relationship detection methods usually have two steps: object detection and relationship recognition, while Pixels2Graph is to directly get objects and relationships from the pixels in the original images. In the method of Pixels2Graph, The elements of the scene graph, including nodes and edges, are detected first, actually that is the objects and their bounding boxes of the relations on the graph are detected. Then these elements are combined with associative embedding to form the relationships of the objects.

4 APPLICATIONS OF SCENE GRAPH

Scene graph can describe the objects in a scene and the relationships between the objects, which provides better visual representations for relevant visual tasks, and can greatly improve the model performance of these visual tasks. In this section, we stated the applications of scene graph to different types of visual tasks.

4.1 Image Retrieval

Image retrieval is a classic visual task in computer vision. For retrieving the target images, the query could be the content of an image or the text describing the image. Commonly, most content-based image retrieval methods use low-level visual features. Recently, there has more and more interest in the models of jointly reasoning about the visual and textual features. However, these models have their limitations in terms of expressiveness. While text-based image retrieval methods have the problem of the inherent referential uncertainty of textual representations. Scene graphs are a structured representation of visual scenes, and a scene graph can explicitly represents the objects, attributes and relationships in images. Therefore, scene graph-based image retrieval has broad development prospects.

In 2015, J.Johnson et al.[1] proposed the concept of scene graph, and design a conditional random field model for image retrieval by utilizing the scene graph, which is constructed manually. In [85], a new text-based image retrieval framework is proposed based on binary representations and semantic graphs, and mainly consists of four parts: cross-modal binary representation, cross-modal semantic graph, the joint objective function and online updating. In 2019, Ramnath et al. [86] proposed a neural-symbolic approach for one-shot image retrieval based on the caption descriptions. This method constructs the catalogs and captions as scene-graphs and achieves the retrieval task by solving the problem of learnable graph matching.

In 2020, Wang et al. [87] proposed to represent image and text with two kinds of scene graphs: visual scene graph (VSG) and textual scene graph(TSG),and the image-text retrieval task is then naturally formulated as cross-modal scene graph matching. Given a query in one modality (a sentence query or an image query), the goal of the image-text cross-modal retrieval task is to find the most similar sample from the database in another modality. Therefore, their Scene Graph Matching (SGM) model aims to evaluate the similarity of the image-text pairs by dissecting the input image and text sentence into scene graphs. The framework of SGM is illustrated in Figure, which consists of two branches of networks. In the visual branch, the input image is represented into a visual scene graph (VSG) and then encoded into the visual feature graph (VFG). Simultaneously, the sentence is parsed into a textual scene graph (TSG) and then encoded into the textual feature graph (TFG) in the textual branch. Finally, the model collects object features and relationship features from the VFG and TFG and calculates the similarity score at the object-level and relationship-level, respectively.

4.2 Image Generation

Image generation for the complex scenes with multiple objects and desired layouts is a hot topic in computer vision research. Despite image generation based on computer vision technology has significant recent progress, it is still a difficult problem to generate the images with multi-objects and complex scenes.

Johnson et al. [36] attempted to generate a realistic image given the corresponding scene graph with object labels and their relationships by Image Generation Network (IG-Net). This problem is a rebuilding work which meets the following three challenges: how to process the graph-structured input, how to guarantee the uniformity between the generated images and their corresponding scene graphs, and how to ensure the authenticity of the synthesized images. These challenges are settled as follows. A

scene graph, which specifies the objects and relationships, is used as the input in IG-Net. IG-Net passes the information along the edges to compute the objects' feature vectors, which thereafter are used to predict the locations of the objects' bounding boxes and segmentation masks. Then, a scene layout can be formed based on these bounding boxes and segmentation masks, and a rough image \hat{I} is generated by using cascaded refinement network (CRN). Subsequently, the authenticity of \hat{I} is solved by adversarially training IG-Net against a pair of discriminator networks D_{img} and D_{obj} . In this process, \hat{I} is encouraged to appear realistic, and also to contain the realistic and recognizable objects. Finally, the generated images can be obtained.

To generate images, Zhao et al. proposed an end-to-end approach (Layout2Im) [88] to generate images from the layouts. In Layout2Im, the representation of each object is broken down into specified and unspecified parts. For the specified parts, the category is coded by word embedding, and for unspecified parts, the visual features are extracted as a low-dimensional vector. Then, individual object representations are grouped together using convolution LSTM to obtain the encoding of the complete layout, which is then decoded into an image. During the process of image generation, several loss terms are introduced to improve the performance of image generation.

Since the previous image generation methods cannot introduce new additional information to the existing description, and are limited to generating images at one time. Therefore, Mittal et al. [89] proposed a recursive network architecture that preserves the image content generated in previous steps and modifies the accumulated images based on newly provided scene information. This method allows to preserve the context in sequentially generated images by subjecting certain information to subsequent image generation conditions.

To solve the problem that it needs to ensure whether the generated image conforms to the scene graph, Tripathi et al. [90] propose an image generation method by harnessing scene graph context to improve image generation. In this method, a scene graph context network is introduced to pool the context features, which are generated by a GCN network. Then these pooled context features are passed to a fully-connected layer, where embeddings are generated for both the generator and the discriminator networks during training. The scene context network encourages the generated images to appear realistic, but also to respect the scene graph relationships.

In [91], a semi-parametric method (PasteGAN) is proposed by Yikang et al. for image generation based on the scene graph and the object crops. The scene graph defines the spatial arrangements of the objects and their relationships, while the given object crops are used to determine the object appearances. Then, two branches of networks are trained simultaneously: One branch focuses on diverse image generation with the object crops, which are retrieved from the external memory; While in the other branch, the original crops are used to reconstruct the ground-truth images.

To improve the quality of generated images, several previous methods are proposed for mapping Scene Graph to images, which is invariant to a set of logical equivalences. Tripathi et al. [92] proposed a new image generation method based scene graph. In this method, the scene graph representations are first enhanced with heuristic-based relations, which increases the minimal storage overhead. Then, the extreme points representations are used to supervise the scene composition network learning.

It is a challenging task of generating the realistic images with

complex visual scenes, especially when we want to control the layouts of the generated images. To this end, Herzig et al. [37] present a novel model to inherently learn the canonical graph representations. In the proposed model, similar predictions can be semantically obtained from similar scene graphs, and the model networks can capture the representation independently of the objects.

To generate a narrative collage for given images, Fang et al. [93] introduced a layer graph and a scene graph to represent the relative depth order and semantic relationships between the objects, so as to generate the main textual descriptions of the images.

4.3 Visual-textual transformer

Since the scene graphs contain the structured semantic information in a visual scene, and the semantic information is mainly reflected in the representations of the objects, attributes and pairwise relationships in the images. Thus, the scene graph can provide beneficial priors for the vision tasks of Image/video Captioning and Visual Question Answering (VQA).

4.3.1 Image/video Captioning

Different from the traditional image captioning methods, a method with scene-graph based semantic representation for image captioning is proposed in [94]. To embed scene graph as an intermediate state, the task of image captioning is divided into two phases: concept cognition and sentence construction respectively. In this method, a CNN-RNN-SVM framework is proposed to generate the scene-graph-based sequence, which is then transformed into a bit vector, as the input of RNN in the next phase for generating the captions.

Accurately grounding text descriptions to the visual relations is critical to most language-and-vision tasks. In [95], Neural Scene Graph Generators are proposed for tackling two language-and-vision tasks of image-text matching and image captioning. Since the Scene Graph Generators can learn the visual relation features more effectively, which facilitates to ground language to visual relations. Subsequently, these two language-and-vision applications have better performance.

Since the graphical representations with conceptual positional binding can improve image captioning, a novel approach, which is derived from regional visual features of the images, is proposed for generating the image captions, and this approach is called Tensor Product Scene-Graph-Triplet Representation (TPsgtR) in [96]. In TPsgtR, the technique of neuro-symbolic embedding is introduced to embed the identified relationships among different image regions into concrete forms. These neural symbolic representations are beneficial to better definition of neural symbolic space for neuro-symbolic attention, and can be transformed to better captions for the images.

The language inductive bias is exploited as language priors in [38], and Scene Graph Auto-Encoder (SGAE) is proposed to incorporate these inductive bias into the encoder-decoder models for image captioning, which is expected to help this encoder-decoder model have less overfitting to the dataset bias. Specifically, in the textual domain, SGAE is used to learn a dictionary (D) for sentence reconstruction. While in the vision-language domain, D can be shared to guide the encoder-decoder models for image captioning. The scene graph has the powerful representation of the complex structural layouts of images and sentences as well as

the shared dictionary D , then the language inductive bias can be transferred across the textual and visual domains.

In [40], a new scene graph-based approach is present for image captioning, and this whole framework comprises two generators (image scene graph generator and sentence scene graph generator) and a pair of encoder and decoder (scene graph encoder and a sentence decoder). Specifically, the pair of encoder and decoder are trained on textual modality. Moreover, an unsupervised feature extraction method is proposed to learn the scene graph features by mapping from the visual features of the images to the textual features of the sentences.

The image captioning framework based on scene graphs is proposed in [39] to solve the problem that the entities in images are considered individually in most previous works, which results in lacking the structured information for generating the sentences. Therefore, the scene graphs are structured by leveraging both visual features and semantic knowledge. The visual features are extracted from the object entities by CNN models, and the semantic relationship features are learned from triples. Based on these obtained features, a hierarchical-attention-based module is designed to learn the discriminative features for generating the sentences.

In [97], the Scene Graph Captioner (SGC) framework is proposed for image captioning. In this framework, the comprehensive structural semantic features are extracted by explicitly modeling the objects, attributes and relationships in the visual scenes, and the LSTM-based models translate these semantic features into the final text descriptions.

Storytelling from an image stream. In [4], the scene graph is used to generate the story from an image stream. The proposed SGVST models visual relations in one image and cross-images, which is conducive to image description. Experimental results show that this method can significantly improve the quality of story generation. The Scene Graph Parser converts an image into a Scene Graph G . Then, the scene Graph is input multi-modal Graph ConvNet, and the nodes in the scene graph are enhanced by Graph convolutional neural network (GCN). In order to model the interaction between images, the temporal convolutional neural network (TCN) is used to further optimize the visual representations of images. Finally, the features of relation aware, which is a set of internal relation and cross-image relation, are obtained and input to Hierarchical Decoder to generate stories.

4.3.2 Visual Question Answering

In [98], an alternative approach is investigated based scene graphs for Visual Question Answering (VQA) inspired by traditional QA systems on knowledge graphs. Specifically, The graph networks (GN) can be applied for encoding the scene graph and performing reasoning according to the given questions. Since scene graphs can capture the essential information of images by the form of graph structures, which is benefit for QA methods to outperform the traditional VQA algorithms.

A VQA method with visual attention is proposed based scene graphs [42]. In this method, natural language explanations comprising of evidences are generated for answering the questions, which are asked to images using two sources of information: the entity annotations generated from the scene graphs and the attention map generated by a VQA model.

For achieving the tasks of visual question answering and visual relationship detection, a new multimodal fusion model

BLOCK is proposed based on the block-superdiagonal tensor decomposition [41] to represent the fine interactions between multi-modalities, while the powerful mono-modal representations are also maintained. Moreover, the end-to-end learnable architectures are designed for representing the relevant interactions between modalities.

In [99], a Scene Graph Convolutional Network (Scene GCN) is designed to jointly reason the object properties and relational semantics for VQA task. In this method, to effectively represent visual relational semantics, a visual relationship encoder is built to yield discriminative and type-aware visual relationship embeddings constrained by both the visual context and language priors. Moreover, SceneGCN is proposed to reason about the visual clues for the correct answer under the guidance of the question.

4.4 Visual social and Human object interaction (HOI) recognition

Visual social relationship recognition and visual human-object relationship recognition are two important computer vision tasks. In this section, we will discuss the existing methods based scene graph for these two tasks.

Social relationships are mainly reflected in the relationship between people, and is the concrete expression of human social structure. At present, We always want to build a intelligent machine that can make better interactions with humans in different types of social environments, then it is critical to develop models to understand and predict the social relationships from images/videos. In [43], Li et al. proposed a Dual-Glance model for social relationship recognition. In this method, the persons are detected, and attention mechanisms are used to exploit contextual cues. Furthermore, to solve the problem that visually identifying social relationship contains uncertainty, an Adaptive Focal Loss is designed by leveraging the ambiguous annotations for the models to more effectively learn the relationship features.

In [100], a Multi-Granularity Reasoning (MGR) framework is designed by Zhang et al. for recognizing the social relationships from images. In MGR, two branches are used to extract the global features and mid-level details, respectively. One branch is constructed to extract the global features about the scenes, and the other branch is mainly focused on the regional cues and fine interactions. In addition, two graphs of pose-guided Person-Object Graph and Person-Pose Graph are designed to model the interactions of human-objects and human-human. Based on these two graphs, social relation reasoning is carried out by GCN. Finally, the human social relations can be predicted by integrating the global CNN features and the reasoning feature obtained by GCNs.

Human object interaction (HOI) recognition is an important basis of distinguishing different types of human actions happened in real world, and HO-RCNN is proposed to detect HOIs in two steps [101]. First, human and object detectors are used to provide the proposals of paired human-object regions, which are then passed into a ConvNet to output the HOI classification scores for the final HOI recognition. Specifically, for a human-object proposal, HO-RCNN classifies its HOIs using a multi-stream network to extract the features from the detected humans, objects and the spatial relations between them.

Furthermore, a multi-task approach based on Zero-Shot Learning is proposed in [102] to scale HOIs. This approach address the challenges of scaling HOI recognition to the long tail of categories

by introducing a zero-shot learning approach, which reasons on the decomposition of HOIs as verbs and objects. In this approach, a factorized mode, which consists of shared neural network layers and independent verb/object networks, is introduced for HOI recognition. The whole model is trained jointly in a multi-task fashion, and the final scores of all verb-object prediction pairs are calculated for the final HOI predictions.

In [103], Transferable inter-activeness Prior (TIAP) is explored for HOI detection by indicating whether human and object have the interactions with each other. The inter-activeness prior of HOI can be learned from the HOI datasets, regardless of HOI category settings. Therefore, the core idea of learning the inter-activeness prior is to design an Inter-activeness Network to learn the general inter-activeness prior from HOI datasets, and to perform Non-Interaction Suppression before HOI classification.

In [104], InteractNet is proposed for HOI detection and recognition. Inspired by a human-centric approach, this network model would be used to achieve the task of detecting $\langle human, verb, object \rangle$ triplets in the challenging photos. For InteractNet, there is a hypothesis that the visual features of the detected persons has powerful cues for localizing the objects they are interacting with, so that the model learns to predict the action-specific density over the object locations based on the visual features of the detected persons. Moreover, InteractNet also can detect the persons of interest and objects jointly, and efficiently infer the interaction triplets in a jointly trained end-to-end system.

Graph Parsing Neural Network (GPNN) is proposed by Qi et al. in [105] for addressing the task of HOIs detection and recognition in images and videos. GPNN is a new framework that can learn the structural knowledge from the images. Given a scene, GPNN can infer the HOI graph structure represented by an adjacency matrix and node labels. During information passing, GPNN iteratively computes the adjacency matrix and node labels.

4.5 Image understanding and reasoning

Compared with unstructured textual descriptions, scene graphs help us make a prediction of the detected objects and their relationships, and even to infer the possible layouts of the images. The detection and recognition of the triplet $\langle subject, relation, object \rangle$ are the key to image understanding and reasoning [46], [47], [106]. Fully understanding an image needs to detect and recognize different visual components, to make a combination of visual modules, reasoning modules and priors, so as to infer the higher-level events and activities.

Wang et al.[48] proposed a deep CNN (IDW-CNN) to improve the image segmentation accuracy by learning from Image Descriptions. IDW-CNN consists of three key parts: ResNet-101 for feature extraction, a network stream for image segmentation label-map prediction and the final network stream for estimating the object interactions. IDW-CNN jointly trains IDW on the existing image segmentation datasets, and fully explores the prior knowledge from other relevant datasets, so as to improve the segmentation performance. As only weak labels are used, so IDW-CNN can also be used for the tasks of Semi- and Weakly-supervised Image Segmentation.

In [107], Yang et al. proposed a novel approach for predicting the precise support relations, and introduced a framework for constructing semantic scene graphs and assessing the quality. In this method, a CNN-based method is used to detect the objects from the given images first. Then, the precise support relations

between objects are predicted by considering the auxiliary information in indoor environments. Finally, a semantic scene graph is constructed to describe the contextual relations between the indoor objects. Compared with the previous methods, this proposed approach achieved better performance of support relation prediction.

Aditya et al.[46] present an intermediate knowledge structure called Scene Description Graph (SDG) for predicting the object interactions in a scene. The objects, scenes and other visual constituents are first detected by a deep learning perception system from input images. Furthermore, Based on the image annotations, a common-sense knowledge base is built by a Bayesian Network, and the scene constituents are finally inferred to predict the object interactions.

Zhang et al.[47] made a research on relationship recognition at an unprecedented scale, where the total number of visual entities is more than 80,000. An image is input to the visual module, and three visual embeddings x_s , x_p , and x_o for subject, relation, and object can be obtained. To this end, a continuous output space is used for objects and relations instead of discrete labels, and a new relationship detection model is developed to embed objects and relations into two vector spaces, and learns a visual and a semantic module to map the features from the two modalities into a shared space.

Shi et al. [106] proposed explainable and explicit Neural Modules (XNMs) for image reasoning based on scene graphs. For a given image and a question, the image is first parsed into a scene graph and the question into a module program, which would be then executed over the scene graph. Furthermore, generic base modules are proposed to conduct reasoning over scene graphs. Besides, XNMs are designed with attention mechanism so as to make all the reasoning steps transparent.

Generating semantic layout from scene graph is a crucial intermediate task of connecting textual descriptions to the relevant images. There would have widely applications based on learning the relation from semantic descriptions to visual appearances, such as text-image matching and text-based image retrieval [108], [109].

The task of referring Expression Grounding (REF) is to localize an specific image region, which is correspondingly described by natural language. To achieve this task, Liu et al.[110] present a REF framework (Marginalized Scene Graph Likelihood (MSGL)) to jointly model all the objects mentioned by language expressions. Compared with traditional models which neglect the rich linguistic structure, MSGL fully exploits the linguistic structures.

In [45], adversarial adaptation of the scene graph models is proposed for reasoning the civic issues, which are mainly reflected by the relationships of the objects. In this model, Faster R-CNN provides the object labels and their locations first, and the contextualized representation of each object is generated based on the object contexts. Edge contexts and the representation of the object pairs can be used for the contextualized representation generation. During adversarial training, information of edge contexts are passed on to the discriminator to distinguish between the seen and unseen object pairs, and the training objective would result in gradients flowing into the discriminator.

Another visual relationship recognition model (Zoom-Net) is proposed in [44] by interpreting the rich interactions between pairs of detected objects from the images and mining deep feature interactions. The method of Spatiality-Context-Appearance Module (SCA-M) is proposed as the core of Zoom-Net, and attempts to extract the contextual features by directly fusing

pairwise object visual features. Furthermore, SCA-M integrates the global and local contextual features, and three classifiers with intra-hierarchy structures are applied to the features for the final visual relationship recognition.

Furthermore, to solve the problem of the diverse interactions among the objects, Plesse et al. [111] proposed guided proposal framework, which makes use of Semantic knowledge and Internal knowledge distillation. In this framework, object detection is the first step towards image understanding, as images can not be fully understood only by recognizing the objects, but the relationships between these objects are more important for the tasks of image understanding and reasoning. To this end, the paper proposes a probabilistic model that translates predicate similarities into probability densities to learn different visual relationships. Moreover, a relevance prediction scheme is present to evaluate the importance of a given object pair to an annotation.

4.6 3D scene graph

3D Scene Graph [49] can provide numerically accurate quantification to the object relationships in 3D scenes. Similar to 2D scene graph generated from 2D images, 3D scene graph describes the environments compactly by abstracting the objects and their relationships in 3D space as graphs, where nodes depict the objects and edges are the relationships between objects. At present, there are few 3D scene graph construction frameworks being proposed [2], [49], while the state-of-the-art methods based on 3D scene graph have got excellent performances for several 3D visual tasks compared with traditional methods.

3-D scene graph is defined in [49] by Kim et al. to represent the physical environments in a sparse and semantic way, and a 3D scene graph construction framework is also proposed, and the 3-D scene graph can illustrate the environments in a sparse manner and cover up an extensive range of physical spaces, which guarantees the scalability. Furthermore, the applicability of the 3-D scene graph is verified by demonstrating two major applications of visual question and answering (VQA) and task planning, and achieved better performance than the traditional methods.

Another 3D scene graph construction method is proposed in [2]. The input to this method is the output of 3D scanners, which consists of 3D mesh models, registered RGB panoramas and the corresponding camera parameters, and each panorama is densely sampled for rectilinear images. Then, Mask R-CNN detection on these images are aggregated back on the panoramas with a weighted majority voting scheme. The final output is the 3D space Graph with a four layered graph. Each layer has a set of nodes with their attributes, and edges representing the relationships between nodes.

According to the current research works on 3D scene graph construction, 3D scene graph is mainly applied to construct 3D interior environments by detecting, locating and recognizing the indoor objects and predicting the positional relationships between these indoor objects.

5 CONCLUSION

It is always the goal of computer vision to have a deep understanding of a scene, and then be able to reason about relevant events, even some unseen events. Since scene graph, a new content for scene description, is proposed in 2015, subsequently, a wave of research works on scene graph generation and application has been set off. Scene graph is a type of data structure that describes

the objects, attributes and the relationship between objects in a scene, and has powerful expression for the scene. While the first scene graph is established manually. Subsequently, many scene graph generation methods are proposed to build a more complete scene graph by a variety of network models, feature extraction methods, and even by introducing the prior knowledge. Meanwhile, some relevant models and methods are designed to reduce the computational complexity of scene graph generation. Furthermore, there are also many research works on applying scene graph to different types of visual tasks, such as image retrieval, image generation, image/video caption and so on. Due to the scene graph's powerful ability of scene representation and the introduction of relevant knowledge information, the performances of these visual tasks are greatly improved. Therefore, this paper gives a systematic overview of the current researches on scene graph generation and application. For scene graph generation, the model types of object relation recognition are classified; while we categorize scene graph applications according to the visual tasks. The review of scene graph generation and application is to summarize the latest scene graph research, point out the problems that still need to be solved in future scene graph research. We expect this review can provide an overall technical reference for scene graph research.

ACKNOWLEDGMENTS

This work is supported in part by NSFC grant 61702415, Australian Research Council (ARC) Discovery Early Career Researcher Award (DECRA) under grant no. DE190100626, Air Force Research Laboratory and DARPA under agreement number FA8750-19-2-0501.

REFERENCES

- [1] J. Johnson, R. Krishna, M. Stark, L.-J. Li, D. Shamma, M. Bernstein, and L. Fei-Fei, "Image retrieval using scene graphs," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3668–3678.
- [2] I. Armeni, Z.-Y. He, J. Gwak, A. R. Zamir, M. Fischer, J. Malik, and S. Savarese, "3d scene graph: A structure for unified semantics, 3d space, and camera," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 5664–5673.
- [3] H. Qi, Y. Xu, T. Yuan, T. Wu, and S.-C. Zhu, "Scene-centric joint parsing of cross-view videos," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [4] R. Wang, Z. Wei, P. Li, Q. Zhang, and X. Huang, "Storytelling from an image stream using scene graphs."
- [5] A. Zareian, S. Karaman, and S.-F. Chang, "Bridging knowledge graphs to generate scene graphs," *arXiv preprint arXiv:2001.02314*, 2020.
- [6] E. E. Aksoy, A. Abramov, F. Wörgötter, and B. Dellen, "Categorizing object-action relations from semantic scene graphs," in *2010 IEEE International Conference on Robotics and Automation*. IEEE, 2010, pp. 398–405.
- [7] S. Aditya, Y. Yang, C. Baral, C. Fermüller, and Y. Aloimonos, "From images to sentences through scene description graphs using common-sense reasoning and knowledge," *arXiv preprint arXiv:1511.03292*, 2015.
- [8] S. Schuster, R. Krishna, A. Chang, L. Fei-Fei, and C. D. Manning, "Generating semantically precise scene graphs from textual descriptions for improved image retrieval," in *Proceedings of the fourth workshop on vision and language*, 2015, pp. 70–80.
- [9] C. Lu, R. Krishna, M. Bernstein, and L. Fei-Fei, "Visual relationship detection with language priors," in *European conference on computer vision*. Springer, 2016, pp. 852–869.
- [10] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma et al., "Visual genome: Connecting language and vision using crowdsourced dense image annotations," *International Journal of Computer Vision*, vol. 123, no. 1, pp. 32–73, 2017.

- [11] J. Peyre, J. Sivic, I. Laptev, and C. Schmid, "Weakly-supervised learning of visual relations," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 5179–5188.
- [12] B. Zhuang, Q. Wu, C. Shen, I. Reid, and A. van den Hengel, "Hcvrd: a benchmark for large-scale human-centered visual relationship detection," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [13] A. Newell and J. Deng, "Pixels to graphs by associative embedding," in *Advances in neural information processing systems*, 2017, pp. 2171–2180.
- [14] M. Klawonn and E. Heim, "Generating triples with adversarial networks for scene graph construction," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [15] L. D. Dai Bo, Zhang Yuqi, "Detecting visual relationships with deep relational networks," in *Proceedings of the IEEE conference on computer vision and Pattern recognition*, 2017, pp. 3076–3086.
- [16] W. Cong, W. Wang, and W.-C. Lee, "Scene graph generation via conditional random fields," *arXiv preprint arXiv:1811.08075*, 2018.
- [17] H. Zhang, Z. Kyaw, S.-F. Chang, and T.-S. Chua, "Visual translation embedding network for visual relation detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 5532–5540.
- [18] Z.-S. Hung, A. Mallya, and S. Lazebnik, "Union visual translation embedding for visual relationship detection and scene graph generation," *arXiv preprint arXiv:1905.11624*, 2019.
- [19] N. Gkanatsios, V. Pitsikalis, P. Koutras, A. Zlatintsi, and P. Maragos, "Deeply supervised multimodal attentional translation embeddings for visual relationship detection," in *2019 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2019, pp. 1840–1844.
- [20] J. Zhang, M. Elhoseiny, S. Cohen, W. Chang, and A. Elgammal, "Relationship proposal networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 5678–5686.
- [21] Y. Liang, Y. Bai, W. Zhang, X. Qian, L. Zhu, and T. Mei, "Vrr-vg: Refocusing visually-relevant relationships," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 10 403–10 412.
- [22] Y. Li, W. Ouyang, X. Wang, and X. Tang, "Vip-cnn: Visual phrase guided convolutional neural network," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1347–1356.
- [23] Y. Chen, Y. Wang, Y. Zhang, and Y. Guo, "Panet: A context based predicate association network for scene graph generation," in *2019 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2019, pp. 508–513.
- [24] R. Zellers, M. Yatskar, S. Thomson, and Y. Choi, "Neural motifs: Scene graph parsing with global context," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5831–5840.
- [25] K. Tang, H. Zhang, B. Wu, W. Luo, and W. Liu, "Learning to compose dynamic tree structures for visual contexts," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6619–6628.
- [26] Y. Li, W. Ouyang, B. Zhou, K. Wang, and X. Wang, "Scene graph generation from objects, phrases and region captions," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 1261–1270.
- [27] Y. Li, W. Ouyang, B. Zhou, J. Shi, C. Zhang, and X. Wang, "Factorizable net: an efficient subgraph-based framework for scene graph generation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 335–351.
- [28] M. Qi, W. Li, Z. Yang, Y. Wang, and J. Luo, "Attentive relational networks for mapping images to scene graphs," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3957–3966.
- [29] Y.-S. Wang, C. Liu, X. Zeng, and A. Yuille, "Scene graph parsing as dependency parsing," *arXiv preprint arXiv:1803.09189*, 2018.
- [30] X. Liang, L. Lee, and E. P. Xing, "Deep variation-structured reinforcement learning for visual relationship and attribute detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 848–857.
- [31] Z. Cui, C. Xu, W. Zheng, and J. Yang, "Context-dependent diffusion network for visual relationship detection," in *Proceedings of the 26th ACM international conference on Multimedia*, 2018, pp. 1475–1482.
- [32] T. Chen, W. Yu, R. Chen, and L. Lin, "Knowledge-embedded routing network for scene graph generation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6163–6171.
- [33] J. Gu, H. Zhao, Z. Lin, S. Li, J. Cai, and M. Ling, "Scene graph generation with external knowledge and image reconstruction," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1969–1978.
- [34] B. A. Plummer, A. Mallya, C. M. Cervantes, J. Hockenmaier, and S. Lazebnik, "Phrase localization and visual relationship detection with comprehensive image-language cues," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 1928–1937.
- [35] S. Wang, R. Wang, Z. Yao, S. Shan, and X. Chen, "Cross-modal scene graph matching for relationship-aware image-text retrieval," *arXiv preprint arXiv:1910.05134*, 2019.
- [36] J. Johnson, A. Gupta, and L. Fei-Fei, "Image generation from scene graphs," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 1219–1228.
- [37] R. Herzig, A. Bar, H. Xu, G. Chechik, T. Darrell, and A. Globerson, "Learning canonical representations for scene graph to image generation," *arXiv preprint arXiv:1912.07414*, 2019.
- [38] X. Yang, K. Tang, H. Zhang, and J. Cai, "Auto-encoding scene graphs for image captioning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 10 685–10 694.
- [39] X. Li and S. Jiang, "Know more say less: Image captioning based on scene graphs," *IEEE Transactions on Multimedia*, vol. 21, no. 8, pp. 2117–2130, 2019.
- [40] J. Gu, S. Joty, J. Cai, H. Zhao, X. Yang, and G. Wang, "Unpaired image captioning via scene graph alignments," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 10 323–10 332.
- [41] H. Ben-Younes, R. Cadene, N. Thome, and M. Cord, "Block: Bilinear superdiagonal fusion for visual question answering and visual relationship detection," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 8102–8109.
- [42] S. Ghosh, G. Burachas, A. Ray, and A. Ziskind, "Generating natural language explanations for visual question answering using scene graphs and visual attention," *arXiv preprint arXiv:1902.05715*, 2019.
- [43] J. Li, Y. Wong, Q. Zhao, and M. S. Kankanhalli, "Visual social relationship recognition," *arXiv preprint arXiv:1812.05917*, 2018.
- [44] G. Yin, L. Sheng, B. Liu, N. Yu, X. Wang, J. Shao, and C. Change Loy, "Zoom-net: Mining deep feature interactions for visual relationship recognition," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 322–338.
- [45] S. Kumar, S. Atreja, A. Singh, and M. Jain, "Adversarial adaptation of scene graph models for understanding civic issues," in *The World Wide Web Conference*, 2019, pp. 2943–2949.
- [46] S. Aditya, Y. Yang, C. Baral, Y. Aloimonos, and C. Fermüller, "Image understanding using vision and reasoning through scene description graph," *Computer Vision and Image Understanding*, vol. 173, pp. 33–45, 2018.
- [47] J. Zhang, Y. Kalantidis, M. Rohrbach, M. Paluri, A. Elgammal, and M. Elhoseiny, "Large-scale visual relationship understanding," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 9185–9194.
- [48] G. Wang, P. Luo, L. Lin, and X. Wang, "Learning object interactions and descriptions for semantic image segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 5859–5867.
- [49] U.-H. Kim, J.-M. Park, T.-J. Song, and J.-H. Kim, "3-d scene graph: A sparse and semantic representation of physical environments for intelligent agents," *IEEE transactions on cybernetics*, 2019.
- [50] B. Thomee, D. A. Shamma, G. Friedland, B. Elizalde, K. Ni, D. Poland, D. Borth, and L.-J. Li, "The new data and new challenges in multimedia research," *arXiv preprint arXiv:1503.01817*, vol. 1, no. 8, 2015.
- [51] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European conference on computer vision*. Springer, 2014, pp. 740–755.
- [52] F. A. Sadeghi Mohammad Amin, "Recognition using visual phrases," in *Computer Vision and Pattern Recognition (CVPR)*, 2011.
- [53] D. Xu, Y. Zhu, C. B. Choy, and L. Fei-Fei, "Scene graph generation by iterative message passing," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 5410–5419.
- [54] B. R. M. Y. Y. Wentong Liao, Lin Shuai, "Natural language guided visual relationship detection," in *arXiv preprint arXiv:1711.06032*, 2017, pp. 1–12.
- [55] R. Yu, A. Li, V. I. Morariu, and L. S. Davis, "Visual relationship detection with internal and external linguistic knowledge distillation," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 1974–1982.

- [56] S. Woo, D. Kim, D. Cho, and I. S. Kweon, "Linknet: Relational embedding for scene graph," in *Advances in Neural Information Processing Systems*, 2018, pp. 560–570.
- [57] G. R. S. J. Ren Shaoqing, He Kaiming, "Faster r-cnn: towards real-time object detection with region proposal networks," *IEEE Transactions on Pattern and Analysis Machine Intelligence*, vol. 39, no. 6, pp. 1137–1149, 2015.
- [58] A. G.-D. Antoine Bordes, Nicolas Usunier, "Translating embeddings for modeling multi-relational data," in *NIPS*, 2013.
- [59] D. E.-S. R. C.-Y. F. A. C. B. Wei Liu, Dragomir Anguelov, "Ssd: Single shot multibox detector," in *ECCV*, 2016, pp. 21–37.
- [60] D. S. G. R. . F. A. Redmon, J., "You only look once: Unified, real-time object detection," in *CVPR*, 2015.
- [61] H. Wan, Y. Luo, B. Peng, and W.-S. Zheng, "Representation learning for scene graph completion via jointly structural and visual embedding," in *IJCAI*, 2018, pp. 949–956.
- [62] Y. Liu, R. Wang, S. Shan, and X. Chen, "Structure inference net: Object detection using scene-level context and instance-level relationships," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 6985–6994.
- [63] J. Zhang, K. Shih, A. Tao, B. Catanzaro, and A. Elgammal, "An interpretable model for scene graph generation," *arXiv preprint arXiv:1811.09543*, 2018.
- [64] Z. A. Simonyan Karen, "Very deep convolutional networks for large-scale image recognition," 2014.
- [65] A. Kolesnikov, A. Kuznetsova, C. Lampert, and V. Ferrari, "Detecting visual relationships using box attention," in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2019, pp. 0–0.
- [66] R. Hu, M. Rohrbach, J. Andreas, T. Darrell, and K. Saenko, "Modeling relationships in referential expressions with compositional modular networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1115–1124.
- [67] W. Gao, Y. Zhu, W. Zhang, K. Zhang, and H. Gao, "A hierarchical recurrent approach to predict scene graphs from a visual-attention-oriented perspective," *Computational Intelligence*, vol. 35, no. 3, pp. 496–516, 2019.
- [68] S. Jae Hwang, S. N. Ravi, Z. Tao, H. J. Kim, M. D. Collins, and V. Singh, "Tensorize, factorize and regularize: Robust visual relationship learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1014–1023.
- [69] J. Yang, J. Lu, S. Lee, D. Batra, and D. Parikh, "Graph r-cnn for scene graph generation," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 670–685.
- [70] R. Herzig, M. Raboh, G. Chechik, J. Berant, and A. Globerson, "Mapping images to scene graphs with permutation-invariant structured prediction," in *Advances in Neural Information Processing Systems*, 2018, pp. 7211–7221.
- [71] A. Dornadula, A. Narcomey, R. Krishna, M. Bernstein, and F.-F. Li, "Visual relationships as functions: Enabling few-shot scene graph prediction," in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2019, pp. 0–0.
- [72] J. Zhang, K. J. Shih, A. Elgammal, A. Tao, and B. Catanzaro, "Graphical contrastive losses for scene graph generation," *arXiv preprint arXiv:1903.02728*, 2019.
- [73] L. Chen, H. Zhang, J. Xiao, X. He, S. Pu, and S.-F. Chang, "Counterfactual critic multi-agent training for scene graph generation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 4613–4623.
- [74] M. Raboh, R. Herzig, J. Berant, G. Chechik, and A. Globerson, "Differentiable scene graphs," in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2019, pp. 0–0.
- [75] B. Schroeder, S. Tripathi, and H. Tang, "Triplet-aware scene graph embeddings," in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2019, pp. 0–0.
- [76] D. S. A. A. R.-J. V. M. G. B. A. G. M. R. A. K. F. G. O. V. Mnih, K. Kavukcuoglu, "Human-level control through deep reinforcement learning," vol. 518, no. 7540, pp. 529–533, 2015.
- [77] M. J. Peter Anderson, Basura Fernando and S. Gould, "Spice: Semantic propositional image caption evaluation," *ECCV*, 2016.
- [78] J. Peyre, I. Laptev, C. Schmid, and J. Sivic, "Detecting unseen visual relations using analogies," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 1981–1990.
- [79] B. M. Z. R. Li Yujia, Tarlow Daniel, "Gated graph sequence neural networks," *Computer Science*, 2015.
- [80] E. Belilovsky, M. Blaschko, J. Kiros, R. Urtasun, and R. Zemel, "Joint embeddings of scene graphs and images," 2017.
- [81] I. M. L. S. Gong Yunchao, Ke Qifa, "A multi-view embedding space for modeling internet images, tags, and their semantics," *International Journal of Computer Vision*, vol. 106, no. 2, 2012.
- [82] L. Z. H. E. X. E. Hu Zhiting, Ma Xuezhe, "Harnessing deep neural networks with logic rules," pp. 2410–2420, 2016.
- [83] W. Wang, R. Wang, S. Shan, and X. Chen, "Exploring context and visual pattern of relationship for scene graph generation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8188–8197.
- [84] J. Ji, R. Krishna, L. Fei-Fei, and J. C. Niebles, "Action genome: Actions as composition of spatio-temporal scene graphs," *arXiv preprint arXiv:1912.06992*, 2019.
- [85] M. Qi, Y. Wang, and A. Li, "Online cross-modal scene retrieval by binary representation and semantic graph," in *Proceedings of the 25th ACM international conference on Multimedia*, 2017, pp. 744–752.
- [86] S. Ramnath, A. Saha, S. Chakrabarti, and M. M. Khapra, "Scene graph based image retrieval—a case study on the clevr dataset," *arXiv preprint arXiv:1911.00850*, 2019.
- [87] S. Wang, R. Wang, Z. Yao, S. Shan, and X. Chen, "Cross-modal scene graph matching for relationship-aware image-text retrieval," in *The IEEE Winter Conference on Applications of Computer Vision*, 2020, pp. 1508–1517.
- [88] B. Zhao, L. Meng, W. Yin, and L. Sigal, "Image generation from layout," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8584–8593.
- [89] A. A. M. S. M. T. Mittal Gaurav, Agrawal Shubham, "Interactive image generation using scene graphs," 2019.
- [90] A. B. H. T. Subarna Tripathi, Anahita Bhiwandiwala, "Using scene graph context to improve image generation," 2019.
- [91] L. Yikang, T. Ma, Y. Bai, N. Duan, S. Wei, and X. Wang, "Pastegan: A semi-parametric method to generate image from scene graph," in *Advances in Neural Information Processing Systems*, 2019, pp. 3950–3960.
- [92] S. Tripathi, S. Nittur Sridhar, S. Sundaresan, and H. Tang, "Compact scene graphs for layout composition and patch retrieval," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 0–0.
- [93] H. F. S. H. C. X. Fei Fang, Miao Yi, "Narrative collage of image collections by scene graph recombination," *IEEE Trans. Vis. Comput. Graph.*, vol. 24, no. 9, pp. 2559–2572, 2018.
- [94] L. Gao, B. Wang, and W. Wang, "Image captioning with scene-graph based semantic concepts," in *Proceedings of the 2018 10th International Conference on Machine Learning and Computing*, 2018, pp. 225–229.
- [95] K.-H. Lee, H. Palangi, X. Chen, H. Hu, and J. Gao, "Learning visual relation priors for image-text matching and image captioning with neural scene graph generators," *arXiv preprint arXiv:1909.09953*, 2019.
- [96] C. Sur, "Tpsgr: Neural-symbolic tensor product scene-graph-triplet representation for image captioning," *arXiv preprint arXiv:1911.10115*, 2019.
- [97] N. Xu, A.-A. Liu, J. Liu, W. Nie, and Y. Su, "Scene graph captioner: Image captioning based on structural visual representation," *Journal of Visual Communication and Image Representation*, vol. 58, pp. 477–485, 2019.
- [98] C. Zhang, W.-L. Chao, and D. Xuan, "An empirical study on leveraging scene graphs for visual question answering," *arXiv preprint arXiv:1907.12133*, 2019.
- [99] Z. Yang, Z. Qin, J. Yu, and Y. Hu, "Scene graph reasoning with prior visual relationship for visual question answering," *arXiv preprint arXiv:1812.09681*, 2018.
- [100] M. Zhang, X. Liu, W. Liu, A. Zhou, H. Ma, and T. Mei, "Multi-granularity reasoning for social relation recognition from images," in *2019 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2019, pp. 1618–1623.
- [101] B. Xu, Y. Wong, J. Li, Q. Zhao, and M. S. Kankanhalli, "Learning to detect human-object interactions with knowledge," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [102] L. Shen, S. Yeung, J. Hoffman, G. Mori, and L. Fei-Fei, "Scaling human-object interaction recognition through zero-shot learning," in *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2018, pp. 1568–1576.
- [103] Y.-L. Li, S. Zhou, X. Huang, L. Xu, Z. Ma, H.-S. Fang, Y.-F. Wang, and C. Lu, "Transferable interactivity prior for human-object interaction detection," *arXiv preprint arXiv:1811.08264*, 2018.
- [104] G. Gkioxari, R. Girshick, P. Dollár, and K. He, "Detecting and recognizing human-object interactions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8359–8367.

- [105] S. Qi, W. Wang, B. Jia, J. Shen, and S.-C. Zhu, "Learning human-object interactions by graph parsing neural networks," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 401–417.
- [106] J. Shi, H. Zhang, and J. Li, "Explainable and explicit visual reasoning over scene graphs," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8376–8384.
- [107] M. Y. Yang, W. Liao, H. Ackermann, and B. Rosenhahn, "On support relations and semantic scene graphs," *ISPRS journal of photogrammetry and remote sensing*, vol. 131, pp. 15–25, 2017.
- [108] B. Li, B. Zhuang, M. Li, and J. Gu, "Seq-sg2sl: Inferring semantic layout from scene graph through sequence to sequence learning," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 7435–7443.
- [109] A. Talavera, D. S. Tan, A. Azcarraga, and K.-L. Hua, "Layout and context understanding for image synthesis with scene graphs," in *2019 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2019, pp. 1905–1909.
- [110] D. Liu, H. Zhang, Z.-J. Zha, and F. Wang, "Referring expression grounding by marginalizing scene graph likelihood," *arXiv preprint arXiv:1906.03561*, 2019.
- [111] F. Plesse, A. Ginsca, B. Delezoide, and F. Prêteux, "Visual relationship detection based on guided proposals and semantic knowledge distillation," in *2018 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2018, pp. 1–6.