

# Visual Relationship Detection with Language prior and Softmax

Jaewon Jung

Visual Intelligence Research Group  
University of Science and Technology, ETRI SCHOOL  
Daejeon, Republic of Korea  
woodcook486@naver.com

Jongyoul Park

Visual Intelligence Research Group  
ETRI  
Daejeon, Republic of Korea  
jongyoul@etri.re.kr

**Abstract**—Visual relationship detection is an intermediate image understanding task that detects two objects and classifies a predicate that explains the relationship between two objects in an image. The three components are linguistically and visually correlated (e.g. “wear” is related to “person” and “shirt”, while “laptop” is related to “table” and “on”) thus, the solution space is huge because there are many possible cases between them. Language and visual modules are exploited and a sophisticated spatial vector is proposed. The models in this work outperformed the state of arts without costly linguistic knowledge distillation from a large text corpus and building complex loss functions. All experiments were only evaluated on Visual Relationship Detection and Visual Genome dataset.

**Index Terms**—Visual relationship, Image understanding, Deep learning.

## I. INTRODUCTION

Understanding images is important in computer vision. In deep learning for computer vision, *object classification* [1], [2], *detection* [3]–[5], *attribute detection* [6], *segmentation* [7], [8] and other tasks have improved performance for image understanding. Although these works are still insufficient for understanding images, there is room for improving their performance. Researchers changed their focus to *Scene graph* [9], *image captioning* [10], *image retrieve* [11] and other related works.

One area is *visual relationship detection* [12]–[17]. Visual relationships are a type of relationship between objects in an image and consist of subject, predicate, and object; e.g. {person, ride, motorcycle}, {person, eat, hamburger}, and {cup, on, table}. These can be considered sentences without adjectives, adverbs, or the in/definite article. The subject and object in a visual relationship are exactly the same as a subject and object in a sentence. A predicate in a visual relationship is different from a predicate in a sentence. In the dictionary, the meaning of “predicate” is the part of the sentence that contains the verb and gives information about the subject, but the predicate of a visual relationship is similar to a verb. It can be a regular verb, a preposition, a comparative, a prepositional verb, a phrasal verb or other words that could explain a connection between objects.

One previous approach [13] considered each visual relationship as a one of a class. e.g. {person, ride, motorcycle}, {person, ride, bicycle} and {person, ride, skateboard} are of

different classes. This fashion requires numerous data because all possible combinations meaning the number of predicates times the number of objects squared are different classes and it results in a huge solution space. Other previous approaches [12], [14], [16], [17] consider detecting objects (subject and object) and a predicate separately. This fashion reduces the solution space rather than the above approach [13] because solution space of the objects (subject and object) and predicate are decoupled; an object detector and a predicate classifier are only needed in this case but this way still requires large amounts of data. This work follows the later approach to solve the problem.

There are three major difficulties for visual relationship detection: first, intra-class variance; a predicate can be involved with any subject and object. e.g. {person, eat, pizza}, {elephant, eat, grass}, {person, use, phone}, {person, use, knife} and so on. These visual relationships are totally difference visually and make the solution space huge. The second difficulty is long-tail distribution. Some of the predicates may occur many times but other certain predicates may only occur once or twice throughout the whole dataset and most of the visual relationships are insufficient for training. This phenomenon brings out a biased dataset and model training result. The third difficulty is class overlapping; some of the predicates in the dataset are almost similar meaning but each data belongs to different class even though their annotations are nearly the same: (near, adjacent, around), (below, under), (look, watch), (next to, feed) etc.

This work utilizes a pair of word vectors, a spatial vector and a union box of two objects boxes to detect visual relationships in an image using a language and visual module. The proposed models significantly outperform the state of arts. All experiments in this paper are conducted on the VRD [12]<sup>1</sup> and Visual Genome dataset (VG) [18]<sup>2</sup>.

## II. RELATED WORK

Object classification [1], [2] is the basis of image understanding, is based on a Convolution Neural Network (CNN). This network learns the features of objects in images and

<sup>1</sup>VRD dataset link : <https://cs.stanford.edu/people/ranjaykrishna/vrd/>

<sup>2</sup>VG dataset link : <https://visualgenome.org>

classifies what objects are in an image. As a result of this research field, several CNNs such as VGG16 [1], ResNet [2] called the backbone network outperform object classification. In visual relationship detection, most papers [12], [14]–[17] use these networks to classify the predicate between two objects; this paper employs VGG16 [1].

Object detection [3]–[5] is the next level of object classification for image understanding. This field also achieves massive success through deep learning. The object detection network localizes objects as bounding boxes in images. R-CNN and Fast/Faster R-CNN [3]–[5] are common object detectors that follow the two-stage approach in which object candidates are proposed while working with RPN [3] and then classify what object is in candidates. Some papers [14], [15] about visual relationships utilize RPN; they show how the employment of RPN and object detection results are improved. This work employs faster R-CNN [3] based on VGG16 [1].

Humanobject interaction recognition [19], [20] is a subset of the visual relationship. In contrast to visual relationships, a subject is fixed as a person; this field focuses on the interaction between a person and an object or another person. Average Precision (AP) is an evaluation metric of this research field. Specifically, they evaluate the AP of the triplet {person, verb, object}, which is called the role AP. Moreover, Ref. [21] is focusing on the action or pose without interaction in an image.

Image captioning is an interesting field in visual tasks in which an image is given as an input and the output is a description that explains that image; this field involves natural language. Recurrent Neural Network (RNN) and Long Short-term Memory (LSTM) [22] are used with CNN. Vinyals et al. [23] proposed an architecture in which CNN encodes visual features in an image and RNN decodes it to natural language.

The scene graph [9], [24], [25] is a higher level image understanding. It is a kind of graph in computer science grounded by the visual. Nodes are objects, relationships are edges, and attributes are a sub-node coupled with objects in an image. This field is related to natural language, so some papers [24], [26] have attempted to generate a scene graph-based image description.

Visual relationship detection is a superset of humanobject interaction. Differently, relationships between any two objects are focused in an image. Some papers [12]–[17], [27], [28] do work to detect visual relationships; Lu et al. [12] established a visual relationship detection task and introduced a VRD dataset [12] which contains four categories such as verb, preposition, spatial and comparative, and only one predicate exists in a visual relationship regardless of the category. Yu et al. [17] improved the detection performance by expensive linguistic knowledge distillation from an internal and external text corpus. Li et al. [15] proposed a top-down pipeline. Unlike other approaches [12], [14], [16], [17], the visual relationships including subject, predicate, and object are detected simultaneously with RPN [3] and the phrase-guided message passing structure (PMPS). Zhang et al. [14] built an equation to embed the visual relationship into space with a class indicator, a location vector, and a visual feature. Ref.

[14], [15] employed RPN in their architecture and said that cooperating with it improved the object detection result. Liang et al. [27] proposed a novel framework called deep Variation-structured Reinforcement Learning (VRL) to detect both visual relationships and attributes to understand the global context in an image and use prior language to build a directed semantic graph. Plummer et al. [28] conjugated visual and language cues for the localization and grounding of phrases in images and gave a special attention to relationships between people and body parts or clothing. Bo et al. [16] represented the predicate by using a union box that included the subject, object, and a spatial module consisting of several convolution layers, and detected visual relationships using a deep relational network.

### III. DIFFICULTIES OF VISUAL RELATIONSHIP DETECTION

#### A. Intra-class Variance

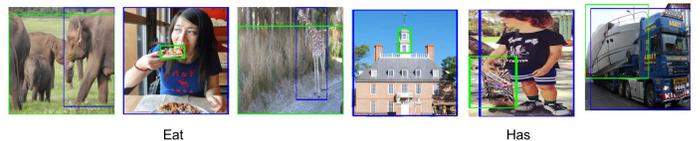


Fig. 1: The examples of Intra-class Variance

The approach that detects objects and the predicate independently requires a predicate classifier. A predicate is involved with many subjects and objects. Therefore visual appearance can have a big gap between visual relationships on the same predicate. {person, eat, pizza} and {elephant, eat, grass} are examples of this.

#### B. Long Tail Distribution

This common problem is mentioned in most papers [12]–[17] and has two aspects: the first is the number of predicates in a dataset and the second is the number of visual relationships. In the VRD [12] and VG [18] dataset, the number of predicate “on” is a huge part of the dataset, but the number of “feed” and “talk” make up a small part of the dataset. Most of the data are small to train because gathering data and annotation are difficult and expensive; subject, predicate, and object can be obtained easy individually but rarely appear together in an image: “airplane,” “next to,” and “bag” are easily obtained individually, but {airplane, next to, bag} is rare.

#### C. Class Overlapping

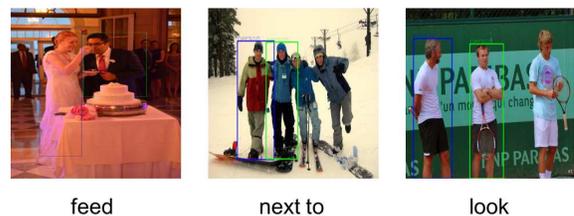


Fig. 2: The examples of Class Overlapping

In a predicate list of both datasets, there are some predicates that mean nearly same thing like “near” and “next to”. According to the dictionary, the difference in the literal meaning between them is vague. For two other cases, some predicates are a superset of others or a subset of others. In visual relationship detection, only one predicate exists between two objects regardless of category. This means that an unrelated predicate that has a totally different meaning can be chosen when the feature vector is close. These phenomena cause wrong classification results.

#### IV. APPROACH

##### A. Visual Module

A visual module is VGG16 [1] and trained similar to object fine-tuning with a softmax loss for classifying the predicate using a union box that includes two objects in an image as an input.

This work relieves ambiguous inferences using only the union box containing two objects, a pair of word vectors, and/or a spatial vector are/is used together. Therefore, variant visual modules are newly created based on the visual module such as a spatial and visual module (SV), a visual and word vector module (VW), and a spatial, visual, and word vector module (SVW).

##### B. Language Module

$$W \times [\text{wordvector}(\text{subject}), \text{wordvector}(\text{object})] + b \quad (1)$$

A language module is trained with the softmax loss instead of the K, L loss in [12], and takes a pair of subject and object word vectors which is 600 dimensions as an input. These are fed to a fully connected layer and produce 70 dimensions vector as an output and 70 is the number of predicates in the dataset.  $W$  is  $70 \times 600$  and  $b$  is 70 dimensions in (1). This simple training approach fulfils the K and, L loss. In [12], the L loss gives a higher likelihood to high-frequency data and a lower likelihood to low-frequency data in a training dataset. The K loss enforces similar visual relationships getting close, and far away from dissimilar visual relationships; e.g. {person, eat, pizza}, {person, eat, hamburger} are similar and {car, has, wheel} is dissimilar from them. It means the language module in [12] produces similar likelihood when visual relationship are close. Without the L loss, the language module in this work naturally assigns the appropriate likelihood to predicates depending on the frequency because several predicates can exist when subject and object are given, only one predicate is annotated for a visual relationship and loss is the softmax ; for an example, “wear” has a higher likelihood than “hold” when the subject is “person” and the object is “shirt”. Without the K loss, the property of word vector leads similar visual relationships to get close and further away from dissimilar visual relationships. {person, ride, bicycle} and {person, ride, motorcycle} naturally get close because “bicycle” and “motorcycle” word vectors are close in a word vector space. This

means that the language module in this work produces nearly same likelihood for “ride” when the subject is “person” and the object is “bicycle” or “motorcycle”.

As with the visual module, a spatial vector is concatenated on the pair of word vectors before the fully connected layer as a new module called the language and spatial module (LS) to relieve the ambiguous inference based on only the pair of word vectors.

##### C. Spatial Vector

$$[IOU, x, y, S_{\text{subject}}/S_{\text{image}}, S_{\text{object}}/S_{\text{image}}, cflag_{\text{subject}}, cflag_{\text{object}}] \quad (2)$$

This work proposes a sophisticated spatial vector different from [17]. The spatial vector in [17] only reflects only each objects bounding box normalized location and size in an image. This encoding is insufficient to classify predicates because predicates do not depend on location in an image. The proposed vector encodes the intersection over union (IOU) and normalized relative location (x, y) based on the subject box center, normalized subject and object size, and contain flag (cflag) for subject and object; cflag for a subject is 1 if the subject box contains the object box, and 0 otherwise, and vice versa.  $S_{\text{subject,object}}$  is the size of the bounding box for each and  $S_{\text{image}}$  is the size of an image in (2).

##### D. Model Variants

$$\text{softmax}(\text{visaulmodule} \times \text{language module}) \quad (3)$$

Several components are available in the modules including a spatial vector, word vectors and a union of the bounding box. The model consists of two modules: the language and visual module. The base component of a language module is a pair of word vectors and spatial vector can be added to the language module. The base component of a visual module is a union box, and word vector and/or spatial vector can be added to the visual module. Furthermore, the language and visual modules can be trained together or separately. Therefore, possible models are L, LS, V, VW, SVW, SV, L + V, L + VW, L + SV, L + SVW, LS + V, LS + VW, LS + SV, and LS + SVW for an experiment (“+” means that two modules are combined). When the language and visual modules are jointly trained, the loss function is (3). “ $\times$ ” means element-wise multiplication. Each module produces 70 dimensions vector.

#### V. EXPERIMENTS

In this section, predicate prediction, phrase, and relationship detection are conducted on the VRD and VG datasets [12], [18]. In predicate prediction, the models take an image and set of localized subjects and objects as an input and predict the set of possible predicates between pairs of objects. In phrase detection, the models take an image as an input, detect the phrase for a set of visual relationships and localize the entire relationship as one bounding box that has at least a 0.5 overlap

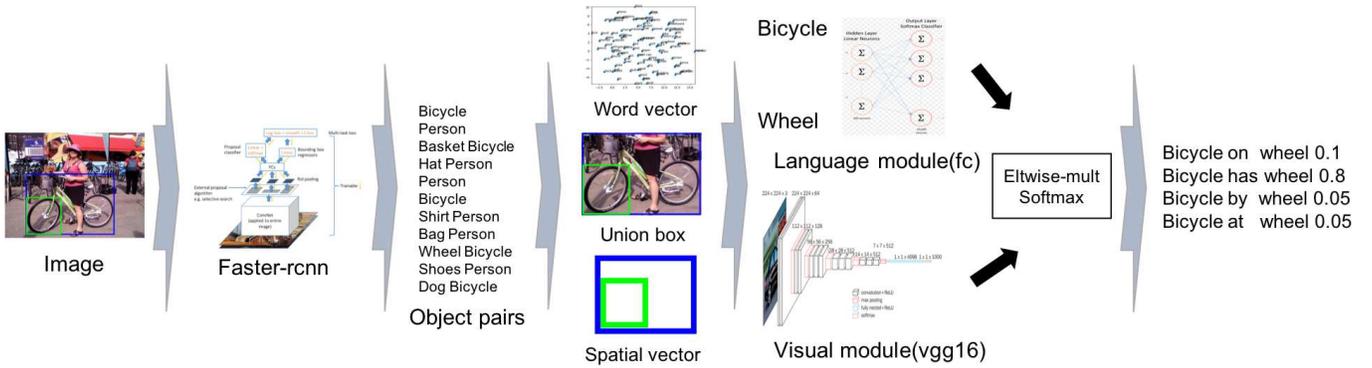


Fig. 3: An overview of a full visual relationship detection model. An image is given as an input. Faster R-CNN [3] based on VGG16 [1] detects object(s) and pairs of objects called candidates are generated from detected objects. Each pair of objects is fed to the predicate classifier with a pair of word vectors, a spatial vector, and a union box that includes subject and object. The language and visual module produce results, multiply it element-wise and softmax is applied to the result.

with the ground truth box. In relationship detection, the models take an image as an input, detect the relationship as a set of visual relationships and localize the subject and object in the image that have at least a 0.5 overlap with their ground truth boxes at the same time.

#### A. Dataset

Particularly for the VG dataset [18], previous works cleaned up the VG dataset on their way. For a fair comparison, the visual relationships that are consistent with the VRD dataset [12] are extracted from the VG dataset to compare with Yus results. For a zero-shot test, a dataset must consist of unseen visual relationships that never occurred in the training dataset for the VRD and VG datasets. The task that extracts unseen visual relationships from the original is applied on both original datasets.

#### B. Evaluation Metric

Then, Recall@n ( $R@n$ ) is chosen as a metric because it is used in [12], [17] and the evaluation algorithm is modified based on [12]. Additionally, the evaluation fashion from [17] is applied. Consistent with [17], the number of chosen predictions ( $k$ ) per object pair is hyper-parameter and shows  $R@n$  for different  $k$  for fair and equal comparison.

#### C. Predicate Prediction

Table I shows the results of predicate prediction on the VRD dataset [12]. The results of L model that only considers a pair of word vector is 44.09  $R@50,100$  when  $k = 1$ . The most referenced predicates are “on” and “wear”, as these are common predicates in the dataset. Better performance is produced from the LS model which takes a pair of word vectors and spatial vector. “on” and “wear” are the most commonly referenced but this model can distinguish spatial predicates. The L + V model outperforms the model in [12] for same condition that a pair of word vectors and a union box are used. In particular, the zero-shot result is nearly 6% higher than [12]. The reason is that the language prior is well obtained in the model in this work rather than the model

TABLE I: Predicate prediction on the VRD dataset. In [17], “U” is the union box that includes two objects, “SF” is the spatial vector in their work, “W” is the word-embedding-based semantic representations, “L” is the linguistic knowledge distillation, “S” is the student network, “T” is the teacher network and “S+T” is the combination of two networks. In this work, “L” is a language module that uses word vectors, “S” is the proposed spatial vector, “V” is the same as “U”, “W” is the word vectors in the visual module, “+” means that the two modules placed before and after the “+” are used together. Before the double vertical line is the general performance and after is the zero-shot performance

	$R@50$ k=1	$R@100$ k=1	$R@50$ k=70	$R@100$ k=70	$R@50$ k=1	$R@100$ k=1	$R@50$ k=70	$R@100$ k=70
VRD [12]	47.87	47.87	-	-	8.45	8.45	-	-
U+W+SF [17]	41.33	41.33	72.29	84.89	14.13	14.13	48.13	69.41
U+W+L:S [17]	42.98	42.98	71.83	84.94	13.89	13.89	51.37	72.53
U+W+L:T [17]	52.96	52.96	83.26	88.98	7.81	7.81	32.62	40.15
U+SF+L:S [17]	41.06	41.06	71.27	84.81	14.33	14.33	48.32	69.01
U+SF+L:T [17]	51.67	51.67	83.84	87.71	8.05	8.05	32.77	41.51
U+W+SF+L:S [17]	47.50	47.50	74.98	86.97	16.98	16.98	54.20	74.65
U+W+SF+L:T [17]	54.13	54.13	82.54	89.41	8.80	8.80	32.81	41.53
U+W+SF+L:T + S [17]	55.16	55.16	85.64	94.65	-	-	-	-
L	44.09	44.09	75.48	86.69	10.86	10.86	50.55	69.71
LS	48.19	48.19	78.31	88.40	15.82	15.82	55.09	74.85
SVW	48.57	48.57	78.04	88.30	16.85	16.85	55.77	74.85
L+V	49.77	49.77	79.99	88.81	14.88	14.88	54.40	72.51
LS+VW	53.05	53.05	85.12	93.17	20.78	20.78	64.67	81.35
LS+SV	53.37	53.37	85.61	93.74	21.21	21.21	65.78	82.37
LS+SVW	55.16	55.16	88.88	95.18	21.38	21.38	64.49	83.49

in [12]. The result from SVW model is a huge improvement over U + W + SF, U + W + L:S and U + SF + L:S. This means that coupling a spatial vector and a pair of word vectors on the visual module works better than [17] and shows the possibility that a model can perform better without linguistic knowledge distillation. Next, the language and visual modules are jointly trained, and LS + SV, LS + VW and LS + SVW



Word vector embedding space provides clusters that have semantically similar word vectors. These clusters help detect unseen visual relationships that never occurred in training dataset. For example, “jacket” and “shirt” resemble one another with regard to wearing but they are different. One is outer clothing and the other is regular clothing. If {person, wear, jacket} only occurs in test dataset, Proposed model easily detects this relationship because “jacket” and “shirt” word vectors are really close in Fig. 4. Particularly when the spatial vector is given, {person, ride, motorcycle}, which never occurs in training dataset can be detected easily rather than detection using only word vector. {person, ride, bicycle} is semantically and spatially related to {person, ride, motorcycle}. From the view of riding a vehicle, the pose is similar between them and the spatial vectors of those are naturally almost the same. The predicate “ride” can be detected with high confidence between “person” and “motorcycle”.

## VI. CONCLUSIONS

The main contribution is outperformed result which can be obtained by simple modification on [12]. The proposed spatial vector is better than the spatial vector in [17] on visual relationship detection. Especially zero-shot performance is significantly improved using proposed spatial vector and word vectors. This paper mentions class overlapping for the first time, which is difficult in visual relationship detection. This work will be shared in public.

## ACKNOWLEDGMENT

This work was supported by Electronics and Telecommunications Research Institute (ETRI) grant funded by the Korean government. [18ZS1100, Core Technology Research for Self-Improving Artificial Intelligence System]

## REFERENCES

- [1] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *CoRR*, vol. abs/1409.1556, 2014.
- [2] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [3] S. Ren, K. He, R. Girshick, and J. Sun, “Faster r-cnn: Towards real-time object detection with region proposal networks,” in *Advances in neural information processing systems*, 2015, pp. 91–99.
- [4] R. Girshick, J. Donahue, T. Darrell, and J. Malik, “Region-based convolutional networks for accurate object detection and segmentation,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 38, no. 1, pp. 142–158, 2016.
- [5] R. Girshick, “Fast r-cnn,” in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1440–1448.
- [6] N. Zhang, M. Paluri, M. Ranzato, T. Darrell, and L. Bourdev, “Panda: Pose aligned networks for deep attribute modeling,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 1637–1644.
- [7] K. He, G. Gkioxari, P. Dollár, and R. Girshick, “Mask r-cnn,” in *Computer Vision (ICCV), 2017 IEEE International Conference on*. IEEE, 2017, pp. 2980–2988.
- [8] J. Long, E. Shelhamer, and T. Darrell, “Fully convolutional networks for semantic segmentation,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3431–3440.
- [9] 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.

- [10] Z. Ren, X. Wang, N. Zhang, X. Lv, and L.-J. Li, “Deep reinforcement learning-based image captioning with embedding reward,” in *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*. IEEE, 2017, pp. 1151–1159.
- [11] J. Wan, D. Wang, S. C. H. Hoi, P. Wu, J. Zhu, Y. Zhang, and J. Li, “Deep learning for content-based image retrieval: A comprehensive study,” in *Proceedings of the 22nd ACM international conference on Multimedia*. ACM, 2014, pp. 157–166.
- [12] 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.
- [13] M. A. Sadeghi and A. Farhadi, “Recognition using visual phrases,” in *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*. IEEE, 2011, pp. 1745–1752.
- [14] 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*, vol. 2, no. 3, 2017, p. 4.
- [15] Y. Li, W. Ouyang, X. Wang, and X. Tang, “Vip-cnn: Visual phrase guided convolutional neural network,” in *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*. IEEE, 2017, pp. 7244–7253.
- [16] B. Dai, Y. Zhang, and D. Lin, “Detecting visual relationships with deep relational networks,” in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2017, pp. 3298–3308.
- [17] R. Yu, A. Li, V. I. Morariu, and L. S. Davis, “Visual relationship detection with internal and external linguistic knowledge distillation,” in *2017 IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2017, pp. 1068–1076.
- [18] 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.
- [19] G. Gkioxari, R. Girshick, P. Dollár, and K. He, “Detecting and recognizing human-object interactions,” *CVPR*, 2018.
- [20] A. Gupta, A. Kembhavi, and L. S. Davis, “Observing human-object interactions: Using spatial and functional compatibility for recognition,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 10, pp. 1775–1789, 2009.
- [21] G. Gkioxari, R. Girshick, and J. Malik, “Contextual action recognition with r\* cnn,” in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1080–1088.
- [22] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [23] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, “Show and tell: A neural image caption generator,” in *Computer Vision and Pattern Recognition (CVPR), 2015 IEEE Conference on*. IEEE, 2015, pp. 3156–3164.
- [24] 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 Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1261–1270.
- [25] 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*, vol. 2, 2017.
- [26] 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.
- [27] X. Liang, L. Lee, and E. P. Xing, “Deep variation-structured reinforcement learning for visual relationship and attribute detection,” in *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*. IEEE, 2017, pp. 4408–4417.
- [28] 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 Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1928–1937.

This figure "class\_overlapping.png" is available in "png" format from:

<http://arxiv.org/ps/1904.07798v1>

This figure "overview.png" is available in "png" format from:

<http://arxiv.org/ps/1904.07798v1>