# Improving Visual Relation Detection using Depth Maps

Sahand Sharifzadeh
*Ludwig Maximilian University of Munich*
sharifzadeh@dbs.ifi.lmu.de

Sina Moayed Baharlou
*Sapienza University of Rome*
baharlou@dis.uniroma1.it

Max Berrendorf
*Ludwig Maximilian University of Munich*
berrendorf@dbs.ifi.lmu.de

Rajat Koner
*Ludwig Maximilian University of Munich*
koner@dbs.ifi.lmu.de

Volker Tresp
*Ludwig Maximilian University of Munich*
*& Siemens AG*
volker.tresp@siemens.com

*Abstract*—Visual relation detection methods rely on object information extracted from RGB images such as 2D bounding boxes, feature maps, and predicted class probabilities. We argue that depth maps can additionally provide valuable information on object relations, e.g. helping to detect not only spatial relations, such as `standing behind`, but also non-spatial relations, such as `holding`. In this work, we study the effect of using different object features with a focus on depth maps. To enable this study, we release a new synthetic dataset of depth maps, *VG-Depth*, as an extension to Visual Genome (VG). We also note that given the highly imbalanced distribution of relations in VG, typical evaluation metrics for visual relation detection cannot reveal improvements of under-represented relations. To address this problem, we propose using an additional metric, calling it *Macro Recall@K*, and demonstrate its remarkable performance on VG. Finally, our experiments confirm that by effective utilization of depth maps within a simple, yet competitive framework, the performance of visual relation detection can be improved by a margin of up to $8\%$.

*Index Terms*—scene graph, visual relation detection, depth maps

## I. INTRODUCTION

Scene Graph Generation, i.e. detecting objects and their relations in images in form of (`subject`, `predicate`, `object`), is a fundamental task in scene understanding and can play an important role in recommender systems, visual question answering, decision making, etc. For example, detecting whether a man is `on` a bike or `next to` a bike is a crucial challenge in autonomous driving. Most works in this area rely on image-based object information such as class labels, bounding boxes and RGB features. We argue that depth maps can additionally provide valuable information about an object's relations as they provide the objects' distance from the camera. This information can help to distinguish between many relations such as `behind`, `in front of` and even improve detection in situations where the objects are nearby such as `covered in`. Figure 1 shows a successfully detected example of the relation (`fence`, `behind`, `dog`) after employing its depth map, and using our model. The goal of this work is to study the effect of using different object features on visual relation detection, with a focus on depth maps.
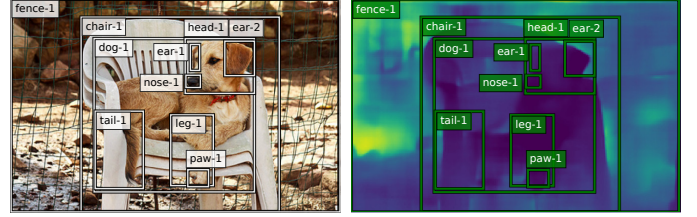


Fig. 1. An image from the VG dataset (left), and the corresponding synthetically generated depth map from VG-Depth dataset (right), annotated by the scene graph. Bright colors in the depth map indicate a larger distance to the camera. Utilizing depth maps allows us to successfully predict the relation (`fence-1, behind, dog-1`).

Unfortunately, most available image datasets, specifically the ones with relational annotations such as Visual Relation Detection (VRD) [1] and Visual Genome (VG) [2], do not provide depth maps, because the acquisition of depth maps is a cumbersome task requiring specialized hardware. We tackle this issue by synthetically generating the corresponding *pseudo* depth maps from 2D images of Visual Genome. This is possible thanks to the large corpora of available RGB-D pairs, i.e. NYU-Depth-v2 [3] dataset. Using RGB-D pairs in NYU-Depth-v2 and a fully convolutional neural network, allow us to learn the mapping function of RGB images to their corresponding depth maps. We can then apply this network to the images from VG, generating their corresponding depth maps. We release the depth maps that are generated from VG, as an extention to it, calling it *VG-Depth*[1]. The object information extracted from depth maps and RGB images, i.e. class labels, location vectors, RGB and depth features, are the basis for relation detection in our simple yet effected framework.

Additionally, we note that the typically employed Recall@K metric (Micro Recall@K), cannot properly reveal the improvements of under-represented relations in highly imbalanced datasets such as VG. This might be an issue in applications such as autonomous driving where it is important to ensure

[1]The dataset and our framework are publicly available at https://github.com/Sina-Baharlou/Depth-VRD.

that the model is capable of predicting also important but less represented predicates such as *walking on* (648 in VG test set) and not just *wearing* (20,148 in VG test set). We address this issue by proposing to employ **Macro Recall@K**, where we compute the mean over Micro Recall@K per predicate, thereby eliminating the effect that over-represented classes have in Micro Recall@K setting.

In summary, our contributions are as follows:

1) We perform an extensive study on the effect of using different sources of object information in visual relation detection. We show in our empirical evaluations using the VG dataset, that our model can outperform competing methods by a margin of up to $8\%$ points, even those using external language sources or contextualization.
2) We release a new synthetic dataset *VG-Depth*, to compensate for the lack of depth maps in Visual Genome.
3) We propose *Macro Recall@K* as a competitive metric for evaluating the visual relation detection performance in highly imbalanced datasets such as Visual Genome.

## II. Related Works

*a) Knowledge Graph (KG) Modeling:* In Knowledge Graph modeling, the aim is typically to find embeddings or latent representations for entities and predicates, which then can serve to predict the probability of unseen triples. These methods mostly differ in how they model relations. In RESCAL [4] each relation is defined as a transformation in the embedding space of entities, producing a triple probability. TransE [5] employs a similar idea but limits each relation to a translation. In comparison to RESCAL, it has fewer parameters; as a disadvantage, it cannot model symmetric relations. DistMult [6] considers each relation as a vector, similar to TransE, but minimizes the trilinear dot product of subject, predicate and object vector. DistMult can be understood as a form of RESCAL, where the transformation matrix is diagonal. ComplEx [7] extends DistMult to complex-valued vectors of embeddings. A multilayer perceptron (MLP) architecture [8] extends these methods to non-linear transformations and has shown to be competitive to the other discussed approaches on most benchmarks [9], [10]. For an extensive review and study on different KG models refer to [9], [11], [12].

*b) Scene Graph (SG) Generation:* SG Generation started with the release of Visual Relation Detection (VRD) [1] and the VG [2]. In VRD, Word2Vec representations of the subject, object, and the predicate were used to train a model jointly with the corresponding image region that describes the predicate. In particular, they consider the joint bounding box of subject and object as the image representation for the predicate. Follow-up work achieved improved performance by incorporating a knowledge graph, constructed from the image annotations [13]. Later, VTransE employed TransE [14] to model visual relations. More recently, Yu et al. [15] proposed a teacher-student model to distill external language knowledge to improve visual relation detection. Iterative Message Passing [16], Neural Motifs [17] (NM) and Graph R-CNN [18]

incorporate context within each prediction using RNNs and graph convolutions respectively. For an extensive discussion on the connection between scene graphs and knowledge graphs refer to [19], [20].

*c) Depth Maps:* While several works have leveraged depth maps to improve *object* detection [21]–[23], the idea of using depth maps in the *relation* detection task has only been explored recently: Yang et al. [24] employ a basic framework for visual relation detection, with handcrafted depth map features, i.e. the mean and mode over pixel values of each depth map. They have a limited experimental setting, where they consider only human-centered relations. In this work, we explore the usability of depth maps in a larger domain and using a convolutional neural network for feature extraction. Furthermore, we provide a more extensive study, release a relevant dataset, and propose a more suitable metric.

## III. Framework

In this section, we introduce the framework that we employed for this study. Let $\mathcal{E} = \{e_1, e_2, ..., e_n\}$ be the set of all entities, including subjects ($s$) and objects ($o$), and $\mathcal{P} = \{p_1, p_2, ..., p_m\}$ the set of all predicates. Each entity $e_i$ can appear in images within a bounding box $\mathbf{bb_i} = [x_i, y_i, w_i, h_i]$, from an image $\mathbf{I}$, where $[x_i, y_i]$ are the coordinates of the bounding box and $[w_i, h_i]$ are its width and height. In this work we apply Faster R-CNN [25], on each image $\mathbf{I}$ to extract a feature map $\mathbf{fmap_I}$, together with object proposals as a set of bounding boxes $\mathbf{bb}$ and class probability distributions $\mathbf{c}$. For each RGB image, we generate a depth map $\mathbf{D}$ where the same bounding box areas encompass the entities' distance from the camera. In the next section, we first describe the synthetic generation of $\mathbf{D}$s and then the feature extraction from generated depth maps. In the end, we describe the relation detection module, where the pairwise features are fused and then employed for relation detection.

### A. Depth Maps for Relation Detection

*1) Generation:* We incorporate an RGB-to-Depth model within our visual relation detection framework. As shown in Figure 2, this is a fully convolutional neural network (CNN) that takes an RGB image as input and generates its predicted depth map. This model can be pre-trained on any datasets containing pairs of RGB and depth maps regardless of having the class annotations for objects or predicates. This enables us to work with the already available visual relation detection datasets without requiring to collect additional data, and also mitigates the need for specialized hardware in real-world applications. The architectural details are explained in Section IV and the generated depth maps from VG are separately released as a dataset called *VG-Depth*.

*2) Feature Extraction:* Depth maps have been employed in tasks such as *object detection* and *segmentation* [22], [26]. In these works, it is common to simply render a depth map as an RGB image, and extract depth features using a CNN, that has been pre-trained with RGB images (for object detection). They argue that the edges in depth maps might
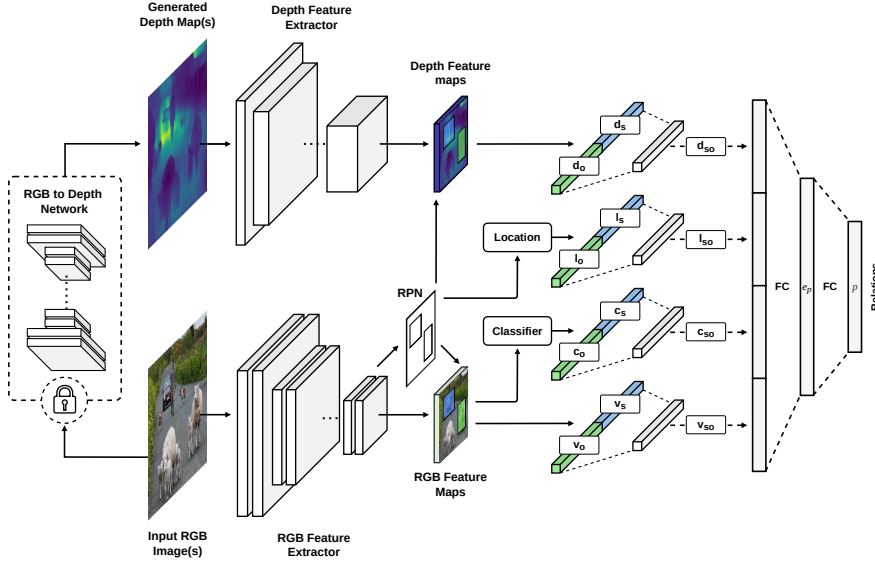
Fig. 2. We study the effect of object information, i.e. class labels, location vectors, RGB and depth features in visual relation detection by employing the simple yet effective framework presented in this figure. We generate depth maps synthetically using an RGB-to-Depth model, eliminating the need for specialized hardware. On the left side, we see the RGB image and its generated depth map, fed into CNNs to extract feature maps from both modalities. We create pairwise feature vectors $\mathbf{d_{so}}$ (pooled from depth feature maps), $\mathbf{l_{so}}$ (from bounding boxes), $\mathbf{c_{so}}$ (from class labels) and $\mathbf{v_{so}}$ (pooled from RGB features) and feed them into a relation detection layer to infer the predicate.

yield better object contours than the edges in cluttered RGB images and that one may combine edges from both RGB and depth to obtain more information [26]. Therefore, they aim to get similar, complementary features from both modalities. However, the practice of employing a model pre-trained on a particular source modality, e.g. RGB, and applying it on a different target modality, e.g. depth map, is sub-optimal in many applications (one should also keep in mind that even fine-tuning some layers of a network does not change the very early convolutional filters). Hence, unlike other works, we train a feature extractor CNN directly on depth maps and specifically for the task of relation detection. Given a depth map $\mathbf{D}$, this network generates a feature map $\mathbf{fmap_D}$. The architectural details of this network is presented in Section IV.

### B. Relation Model

In the previous section, we described methods for the extraction of $\mathbf{fmap_I}$, $\mathbf{fmap_D}$, $\mathbf{c}$ and $\mathbf{bb}$. Here, we outline the model that infers relations using pairwise combinations of these features. For each pair of detected objects within an image, we create a scale-invariant location feature $\mathbf{l_s} = [t_x, t_y, t_w, t_h]$ with: $t_x = (x_s - x_o)/w_o, t_y = (y_s - y_o)/h_o, t_w = \log(w_s/w_o), t_h = \log(h_s/h_o)$ and similarly $\mathbf{l_o}$. We then pool the corresponding features $\mathbf{v_s}$ and $\mathbf{v_o}$ from $\mathbf{fmap_I}$ and create a visual feature vector $[\mathbf{v_s}; \mathbf{v_o}]$. Similarly, we create a depth feature vector $[\mathbf{d_s}; \mathbf{d_o}]$, by pooling features from $\mathbf{fmap_D}$, within $\mathbf{bb}_s$ and $\mathbf{bb}_o$. Additionally, we create $[\mathbf{c_s}; \mathbf{c_o}]$ and $[\mathbf{l_s}; \mathbf{l_o}]$. Each of these vectors are fed into separate fully connected layers, followed by ReLUs, yielding $\mathbf{v_{so}}$, $\mathbf{l_{so}}$,

$\mathbf{c_{so}}$ and $\mathbf{d_{so}}$ before being fed to the relation head which projects them to the relation space such that:

$$\mathbf{e_p} = f(\mathbf{W}[\mathbf{v_{so}}; \mathbf{l_{so}}; \mathbf{c_{so}}; \mathbf{d_{so}}]) \tag{1}$$

Here, $\mathbf{W}$ describes a linear transformation and $f(.)$ is a non-linear function. We realize them as a fully connected layer in a neural network with ReLU activations and dropout. $\mathbf{e_p}$ is an embedding vector of pairwise features. This simple relation prediction model is inspired by the work of [8] to predict links in knowledge graphs. Therefore, we call it **ERMLP-E**, short for ERMLP-Extended. The input of their proposed model is a triple and the output is a single Bernoulli variable, whereas in our work the inputs are *subject* and *object* and we have a Bernoulli variable for each predicate class in the output. This gives us fewer parameters compared to that model, and simplifies training by imposing an implicit negative sampling through the cross-entropy loss.

As shown in earlier works, using more sophisticated models for context propagation between objects with RNNs or graph convolutions, can further improve the prediction accuracy. However, the aim here is to study the effect of including depth maps as additional object features in visual relation detection and as will be shown later, even with this simple model, utilizing depth maps can be more effective than e.g. propagating context. Clearly, those other models can also further enrich their understanding of object relations by employing depth maps.

To learn the parameters, we consider each relation (subject, predicate, object) with an associated Bernoulli variable that takes 1 if the triple is observed and 0 otherwise, following a locally closed world assumption [9].

TABLE I
PREDICATE PREDICTION RECALL VALUES ON VG TEST SET. WHEN THE DEPTH MAPS ARE UTILIZED TOGETHER WITH ALL OTHER FEATURES (*Ours*-l, **c**, **v**, **d**), WE GAIN A LARGE IMPROVEMENT COMPARED TO THE STATE-OF-THE-ART. ONE CAN ALSO SEE THAT EVEN REPLACING DEPTH MAPS WITH VISUAL FEATURES (*Ours*-l, **c**, **d** COMPARED TO *Ours*-l, **c**, **v**) CAN YIELD BETTER RESULTS. ADDITIONALLY, COMPARING *Ours*-l, **c**, **v** TO *VTransE* AND *Neural Motifs* REVEALS THE ADVANTAGE OF OUR SIMPLE MODEL REGARDLESS OF DEPTH MAPS.

| | Strategy | **Macro** | | | **Micro** | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | Task | Predicate Pred. | | | Predicate Pred. | | |
| | Metric | R@100 | R@50 | R@20 | R@100 | R@50 | R@20 |
| models | VTransE [27] | - | - | - | 62.87 | 62.63 | - |
| | Yu's-S [15] | - | - | - | 49.88 | - | - |
| | Yu's-S+T [15] | - | - | - | 55.89 | - | - |
| | IMP [16] | - | - | - | 53.00 | 44.80 | - |
| | Graph R-CNN [18] | - | - | - | 59.10 | 54.20 | - |
| | NM [17] | 14.39 | 13.20 | 10.25 | 67.10 | 65.20 | 58.50 |
| ablations | Ours - $d$ | 9.51 | 8.46 | 6.35 | 54.72 | 51.90 | 43.86 |
| | Ours - $c$ | 15.65 | 13.09 | 8.56 | 64.82 | 60.54 | 49.89 |
| | Ours - $v$ | 13.88 | 12.24 | 8.99 | 61.72 | 58.50 | 50.41 |
| | Ours - $l$ | 5.19 | 4.66 | 3.57 | 49.07 | 46.13 | 37.48 |
| | Ours - $v, d$ | 15.47 | 14.04 | 10.83 | 62.88 | 60.52 | 53.07 |
| | Ours - $l, v, d$ | 15.76 | 14.40 | 11.07 | 63.06 | 60.83 | 53.55 |
| | Ours - $l, c, d$ | 21.67 | 19.56 | 15.12 | 67.97 | 66.09 | 59.13 |
| | Ours - $l, c, v$ | 19.16 | 17.72 | 13.93 | 67.94 | 66.06 | 59.14 |
| | Ours - $l, c, v, d$ | **22.72** | **20.74** | **16.40** | **68.00** | **66.18** | **59.44** |

Given the set of observed triples $\mathcal{T}$, the loss function is the categorical cross entropy between the one-hot targets and the distribution obtained by softmax over the network's output defined as:

$$\mathcal{L} = \sum_{(s,p,o) \in \mathcal{T}} - \log \frac{\exp\left(\mathbf{w'}_p^{\mathrm{T}} \mathbf{e_p}\right)}{\sum_{p' \in \mathcal{P}} \exp\left(\mathbf{w'}_{p'}^{\mathrm{T}} \mathbf{e_p}\right)} \quad (2)$$

where $\mathbf{w'}_p$ is the weight vector corresponding to $p$ in the last layer (linear classification layer).

## IV. EVALUATION

In our study, we are interested to answer the following questions:

1) If we are given *only* depth maps of some objects in a scene (and not even object labels), how accurately can we infer the distribution of possible pairwise relations? How do other sources of object information compare to it?
2) Current visual relation detection frameworks commonly rely on extensive object information such as class labels, bounding boxes, RGB features, contextual information, etc. Do depth representations bring any additional information or would they only contribute redundant scene knowledge?

Additionally, we study whether Recall@K can sufficiently reflect the improvements of under-represented relations within a highly imbalanced dataset such as VG.

In what follows, we introduce the dataset, metrics, architectural details and experiments to answer these questions.

### A. Dataset

We test our approach on the *Visual Genome* [2] dataset. We use the more commonly used subset of VG dataset proposed by [16] which contains 150 object classes and 50 relations.

### B. Metrics

*a) Micro Recall@K:* This metric is defined as the mean prediction accuracy in each image given the top $K$ predictions and is typically called *Recall@K*. We assigned the *Micro* prefix to its name to distinguish this metric with *Macro Recall@K*. Recall@K is a popular choice in most of the visual relation detection studies. The main reason is the incompleteness of visual relation detection datasets, i.e. some relations might not be annotated in the test set, while due to the model's generalization, they might get higher prediction values than the annotated ones. This sensitivity is handled by the $K$ parameter in Recall@K.

*b) Macro Recall@K:* We define this metric as:

$$\text{MACRO RECALL@K} = \sum_{(s,p,o) \in \mathcal{T}_p} \frac{\text{MICRO R@K}(p)}{|\mathcal{T}_p|} \quad (3)$$

where $\mathcal{T}_p \subset \mathcal{T}$ is set of all relations with predicate $p$, and MICRO R@K$(p)$ is computed on $\mathcal{T}_p$. The motivation behind this metric is the highly imbalanced distribution of classes in some datasets such as VG. In these datasets Micro Recall@K score gets dominated by frequently labeled relations and might not reflect the improvements in some important but under-represented classes. However, in Macro R@K, the prediction accuracy of under-represented classes can have a stronger effect on the output. This metric is inspired from the Macro F1 measure [28].
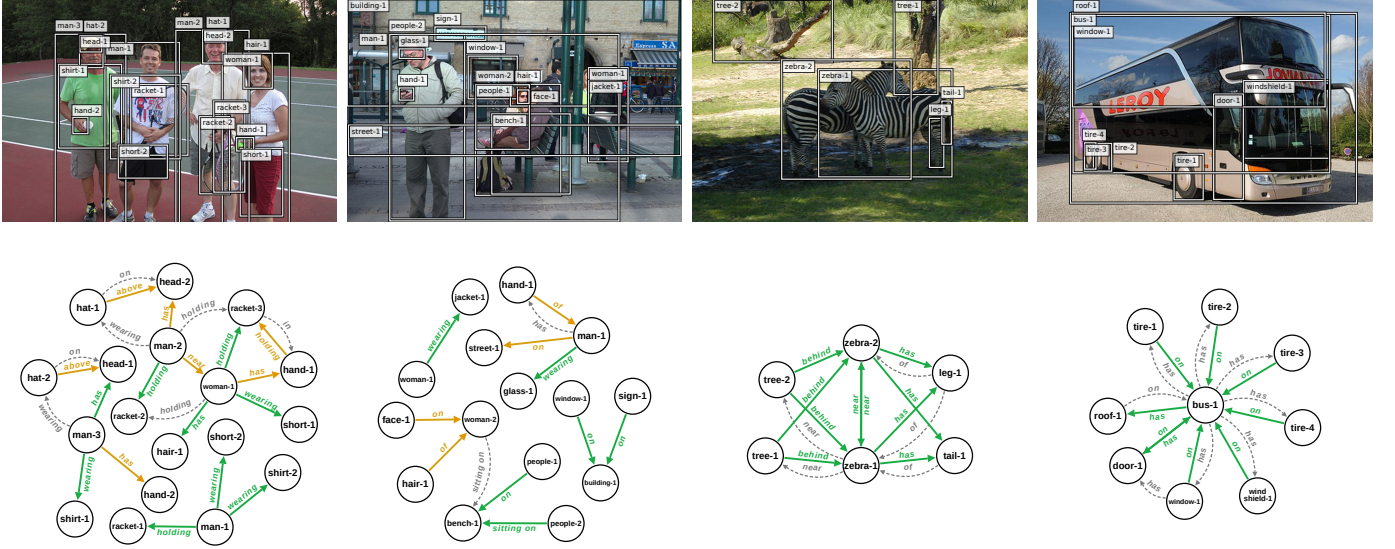
Fig. 3. Some of the qualitative results from our model's predictions. Green arrows indicate the successfully detected predicates (true positives), orange arrows indicate the false negatives and gray arrows indicate predicted links which are not annotated in the ground truth.

## C. Architectures

*a) RGB-to-Depth Network:* We employ the RGB-to-Depth architecture that has been introduced in [29]. The model is a fully convolutional neural network built on ResNet-50 [30], and trained in an end-to-end fashion on data from NYU Depth Dataset v2 [3]. In our experiments, we also trained the model from the outdoor images of Make3D dataset [31]. However, the model that was trained on this dataset, did not show promising results for relation detection. This observation is not surprising because unlike Visual Genome, Make3D images contain mostly outdoor scenes with very few objects.

*b) RGB Feature Extraction:* To extract embeddings and class probabilities of RGB images, we use the VGG-16 architecture [32] pre-trained on ImageNet [33] and fine-tuned on VG by Zellers et al. [17].

*c) Depth Map Feature Extraction:* For depth map extraction we use ResNet-18 proposed in [30]. We trained this model from scratch following the earlier discussions in Subsection III-A2. This network was trained separate from other inputs and on a pure depth-based, relation detection task using Adam [34], with a learning rate of $10^{-4}$ and batch size of 32 for 30 epochs.

*d) Relation Detection Network:* In relation detection head, each extracted feature pair goes to a separate, fully connected hidden layer of 64 neurons (∼12K learnable weights) for class probabilities, 512 for RGB feature maps (∼4M learnable weights), 4096 for depth feature maps (∼4M learnable weights) and 20 for location features (160 learnable weights). Each of them with a dropout rate of 0.1, 0.8, 0.6 and 0.1. The concatenated outputs are then connected to a fully connected hidden layer of 4096 neurons with 0.1 dropout and then to the classification layer. We trained this network by Adam [34], with a learning rate of $10^{-5}$. We used a batch size of 16 and

30 epochs of training. All of the layers were initialized with Xavier weights [35].

## D. Comparing Methods

We compare our results with *VTransE* [27] that takes visual embeddings and projects them to relation space using TransE. We also compare to the student network of [15] (*Yu's-S*), and their full model (*Yu's-S+T*) that employs external language data from Wikipedia. From the context propagating methods, we report Neural Motifs [17], Graph R-CNN [18] and IMP [16]. In an ablation study, we report our relation prediction results under several settings in which different combinations of object information are employed for prediction.

## E. Experiments

As our main goal is to investigate the role of depth maps and other features in relation detection, we report *predicate prediction* results. In this setting, the relation detection performance is analyzed by isolating it from the object detector's error. Therefore, the goal is to evaluate the relation detection accuracy given the objects in an image. We carried on our experiments by training each model 8 times with different random seeds. The maximum variance of the results was no more than 0.01. The results are shown in Table I. In what follows, we provide a discussion over the quantitative and qualitative results.

The upper part of the table demonstrates the results directly reported from other works while the lower part presents the results from the ablation study on our model. For NM, we have computed the Macro R@K results using their publicly available code. We can see that our full model with depth maps, achieves the highest accuracy in comparison to the others in all settings. It is also interesting to note that when
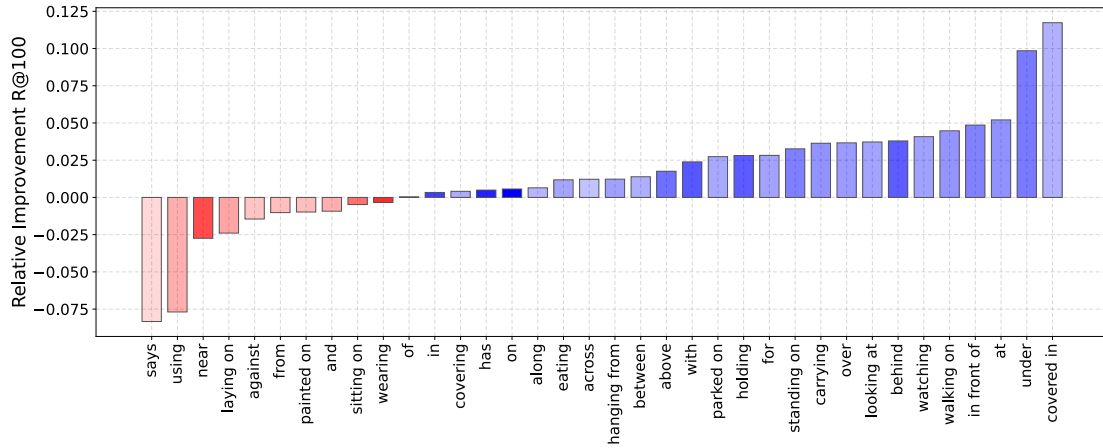
Fig. 4. This plot shows the prediction changes per predicate, going from *Ours-v* to *Ours-v, d*. The classes with zero changes are omitted from the plot. The darker shades indicate larger number of that class within the test set whereas the lighter shades are under-represented classes. An improvement in predicates with more frequency has a larger effect on the Micro R@K whereas this effect is eliminated within Macro R@K. We can see that indeed the improvements by using depth maps are mostly happening within the less-represented classes.

using *only* depth maps we can already achieve a significant accuracy in predicate prediction, emphasizing the value of relational information that are stored within the depth maps alone. By comparing *Ours-v* to *Ours-v, d*, we can observe the improvements that depth maps bring. Also comparing *Ours-l, c, d* to *Ours-l, c, v* is specially informative from two aspects: (1) It shows that while some results are almost equal in Micro settings, one can observe a significant difference in the Macro setting, demonstrating the effectiveness of this metric in presenting the improvements of under-represented classes. (2) We observe that $v$ alone has a higher R@K than $d$ alone. However, when we add them separately to $c, l$ we can see that $d$ has more to offer. In other words, $v$ brings more redundant information to $c, l$ compared to $d$. To get a better intuition of the improvements that we gain after including depth maps (*Ours-v, d* compared to *Ours-v*), we plotted the changes in prediction accuracy for each predicate in Figure 4. We used darker shades for over-represented classes and lighter shades for under-represented ones. This helps to also gain a better intuition of improvement versus frequency of data. For example we can see that in general the accuracy of relations including the predicates such as under, in front of and behind has been improved. These predicates appear much less often in the dataset than on or has, having less effect in the computed *Micro* accuracy. Figure 5 presents some samples of synthetically generated depth maps in VG-Depth dataset including both high quality and faulty ones. Additionally, we present some of the predicted relations by our model in Figure 3.

## V. CONCLUSION

We employed an RGB-to-Depth network, trained on a large corpus of data, to generate depth maps for Visual Genome dataset, releasing a new extension called *VG-Depth*. We provided a metric, *Macro R@K* for better evaluation of relation detection in Visual Genome and other highly imbalanced datasets. In extensive empirical evaluations, we demonstrated the effect of different object features in visual relation detection and showed that by using depth information, we achieve significantly better performance compared to other state-of-the-art methods.

## REFERENCES

[1] 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.

[2] 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.

[3] P. K. Nathan Silberman, Derek Hoiem and R. Fergus, "Indoor segmentation and support inference from rgbd images," in *ECCV*, 2012.

[4] M. Nickel, V. Tresp, and H.-P. Kriegel, "A three-way model for collective learning on multi-relational data." in *ICML*, vol. 11, 2011, pp. 809–816.

[5] A. Bordes, N. Usunier, A. Garcia-Duran, J. Weston, and O. Yakhnenko, "Translating embeddings for modeling multi-relational data," in *Advances in Neural Information Processing Systems 26*, C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, Eds. Curran Associates, Inc., 2013, pp. 2787–2795. [Online]. Available: http://papers.nips.cc/paper/5071-translating-embeddings-for-modeling-multi-relational-data.pdf

[6] B. Yang, W.-t. Yih, X. He, J. Gao, and L. Deng, "Embedding entities and relations for learning and inference in knowledge bases," *arXiv preprint arXiv:1412.6575*, 2014.

[7] T. Trouillon, J. Welbl, S. Riedel, É. Gaussier, and G. Bouchard, "Complex embeddings for simple link prediction," in *International Conference on Machine Learning*, 2016, pp. 2071–2080.

[8] X. Dong, E. Gabrilovich, G. Heitz, W. Horn, N. Lao, K. Murphy, T. Strohmann, S. Sun, and W. Zhang, "Knowledge vault: A web-scale approach to probabilistic knowledge fusion," in *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2014, pp. 601–610.
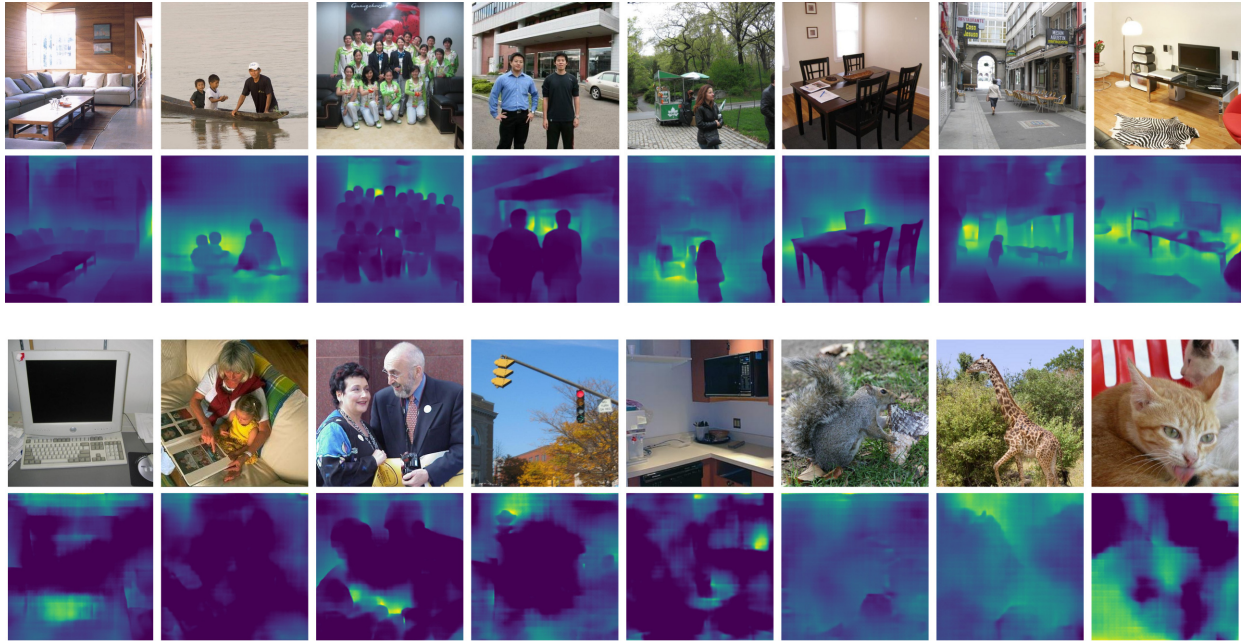
Fig. 5. The first two rows are the examples of visual genome images and their synthetically generated high quality depth maps. The second two rows are the examples of visual genome images and their synthetically generated noisy depth maps.

[9] M. Nickel, K. Murphy, V. Tresp, and E. Gabrilovich, "A review of relational machine learning for knowledge graphs," *Proceedings of the IEEE*, vol. 104, no. 1, pp. 11–33, 2016.

[10] R. Socher, D. Chen, C. D. Manning, and A. Ng, "Reasoning with neural tensor networks for knowledge base completion," in *Advances in neural information processing systems*, 2013, pp. 926–934.

[11] M. Ali, M. Berrendorf, C. T. Hoyt, L. Vermue, M. Galkin, S. Sharifzadeh, A. Fischer, V. Tresp, and J. Lehmann, "Bringing light into the dark: A large-scale evaluation of knowledge graph embedding models under a unified framework," *arXiv preprint arXiv:2006.13365*, 2020.

[12] M. Ali, M. Berrendorf, C. T. Hoyt, L. Vermue, S. Sharifzadeh, V. Tresp, and J. Lehmann, "Pykeen 1.0: A python library for training and evaluating knowledge graph emebddings," *arXiv preprint arXiv:2007.14175*, 2020.

[13] S. Baier, Y. Ma, and V. Tresp, "Improving visual relationship detection using semantic modeling of scene descriptions," in *International Semantic Web Conference*. Springer, 2017, pp. 53–68.

[14] A. Bordes, N. Usunier, A. Garcia-Duran, J. Weston, and O. Yakhnenko, "Translating embeddings for modeling multi-relational data," in *Advances in neural information processing systems*, 2013, pp. 2787–2795.

[15] R. Yu, A. Li, V. I. Morariu, and L. S. Davis, "Visual relationship detection with internal and external linguistic knowledge distillation," in *IEEE International Conference on Computer Vision (ICCV)*, 2017.

[16] 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.

[17] 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.

[18] 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.

[19] V. Tresp, S. Sharifzadeh, and D. Konopatzki, "A model for perception and memory."

[20] V. Tresp, S. Sharifzadeh, D. Konopatzki, and Y. Ma, "The tensor brain: Semantic decoding for perception and memory," *arXiv preprint arXiv:2001.11027*, 2020.

[21] L. Bo, X. Ren, and D. Fox, "Unsupervised feature learning for rgb-d based object recognition," in *Experimental Robotics*. Springer, 2013, pp. 387–402.

[22] A. Eitel, J. T. Springenberg, L. Spinello, M. Riedmiller, and W. Burgard, "Multimodal deep learning for robust rgb-d object recognition," in *Intelligent Robots and Systems (IROS), 2015 IEEE/RSJ International Conference on*. IEEE, 2015, pp. 681–687.

[23] S. Gupta, R. Girshick, P. Arbeláez, and J. Malik, "Learning rich features from rgb-d images for object detection and segmentation," in *European Conference on Computer Vision*. Springer, 2014, pp. 345–360.

[24] H.-K. Yang, A.-C. Cheng, K.-W. Ho, T.-J. Fu, and C.-Y. Lee, "Visual relationship prediction via label clustering and incorporation of depth information," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 0–0.

[25] 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.

[26] C. Hazirbas, L. Ma, C. Domokos, and D. Cremers, "Fusenet: Incorporating depth into semantic segmentation via fusion-based cnn architecture," in *Asian Conference on Computer Vision*. Springer, 2016, pp. 213–228.

[27] H. Zhang, Z. Kyaw, S. Chang, and T. Chua, "Visual translation embedding network for visual relation detection," in *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*. IEEE Computer Society, 2017, pp. 3107–3115. [Online]. Available: https://doi.org/10.1109/CVPR.2017.331

[28] H. Schütze, C. D. Manning, and P. Raghavan, "Introduction to information retrieval," in *Proceedings of the international communication of association for computing machinery conference*, 2008, p. 260.

[29] I. Laina, C. Rupprecht, V. Belagiannis, F. Tombari, and N. Navab, "Deeper depth prediction with fully convolutional residual networks," in *3D Vision (3DV), 2016 Fourth International Conference on*. IEEE, 2016, pp. 239–248.

[30] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *CoRR*, vol. abs/1512.03385, 2015. [Online]. Available: http://arxiv.org/abs/1512.03385

[31] A. Saxena, M. Sun, and A. Y. Ng, "Learning 3-d scene structure from a single still image," in *Computer Vision, 2007. ICCV 2007. IEEE 11th International Conference on*. IEEE, 2007, pp. 1–8.

[32] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[33] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, "ImageNet Large Scale Visual Recognition Challenge,"

*International Journal of Computer Vision (IJCV)*, vol. 115, no. 3, pp. 211–252, 2015.

[34] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[35] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," in *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, 2010, pp. 249–256.