

Explanation-based Weakly-supervised Learning of Visual Relations with Graph Networks

Federico Baldassarre, Kevin Smith, Josephine Sullivan, and Hossein Azizpour

KTH - Royal Institute of Technology, Stockholm, Sweden
{fedbal,ksmith,sullivan,azizpour}@kth.se

Abstract. Visual relationship detection is fundamental for holistic image understanding. However, the localization and classification of (subject, predicate, object) triplets remain challenging tasks, due to the combinatorial explosion of possible relationships, their long-tailed distribution in natural images, and an expensive annotation process.

This paper introduces a novel weakly-supervised method for visual relationship detection that relies on minimal image-level predicate labels. A graph neural network is trained to classify predicates in images from a graph representation of detected objects, implicitly encoding an inductive bias for pairwise relations. We then frame relationship detection as the *explanation* of such a predicate classifier, i.e. we obtain a complete relation by recovering the subject and object of a predicted predicate.

We present results comparable to recent fully- and weakly-supervised methods on three diverse and challenging datasets: HICO-DET for human-object interaction, Visual Relationship Detection for generic object-to-object relations, and UnRel for unusual triplets; demonstrating robustness to non-comprehensive annotations and good few-shot generalization.

1 Introduction

Visual perception systems, built to understand the world through images, are not only required to identify objects, but also their interactions. Visual relationship detection aims at forming a holistic representation by identifying triplets in the form (subject, predicate, object). Subject and object are localized and classified instances such as a cat or a boat, and predicates include actions such as *pushing*, spatial relations such as *above*, and comparatives such as *taller than*.

In recent years, we have witnessed unprecedented development in various forms of object recognition; from classification to detection, segmentation, and pose estimation. Yet, the higher-level visual task of inter-object interaction recognition remains unsolved, mainly due to the combinatorial number of possible interactions w.r.t. the number of objects. This issue not only complicates the inference procedure, but also complicates data collection – the cost of gathering and annotating data that spans a sufficient set of relationships is enormous. In this work, we propose a novel inference procedure that requires minimal labeling thereby making it easier and cheaper to collect data for training.¹

¹ PyTorch implementation, data and experiments: github.com/baldassarreFe/ws-vrd



Fig. 1: **Weakly-supervised relationship detection:** detecting all $\langle \text{subj}, \text{pred}, \text{obj} \rangle$ triplets by training only on weak image-level predicate annotations $\{\textit{push}, \textit{wear}, \textit{drive}\}$

Consider the problem of adding a predicate category to a small vocabulary of 20 objects. A single predicate could introduce up to 20^2 new relationship categories, for which samples must be collected and models should be trained. Moreover, we know that the distribution of naturally-occurring triplets is long-tailed, with combinations such as *person ride dog* rarely appearing [29]. This exposes standard training methods to issues arising from extreme class imbalance. These challenges have prompted modern techniques to take a compositional approach [24,34,15,29] and to incorporate visual and language knowledge [24,31,29], improving both training and generalization.

Although some progress has been made towards recognition of rare triplets, successful methods require training data with exhaustive annotation and localization of $\langle \text{subj}, \text{pred}, \text{obj} \rangle$ triplets. This makes weakly-supervised learning a promising research direction to mitigate the costs and errors associated with data collection. Nonetheless, we identified only two weakly-supervised works tackling general visual relation detection [30,48], both requiring image-level triplet annotation. In this work, we use an even weaker setup for visual relationship detection that relies only on *image-level predicate* annotations (figure 1).

To achieve that, we decompose a probabilistic description of visual relationship detection into the subtasks of object detection, predicate classification and retrieval of localized relationship triplets. Due to considerable progress in object detection, we focus on the last two and use existing pre-trained models for object detection. For predicate classification, we use graph neural networks operating on a graph of object instances, encoding a strong inductive bias for object-object relations. Finally, we use backward explanation techniques to attribute the graph network’s predicate predictions to pairs of objects in the input.

Contributions. The main contributions of this work are threefold:

- I) We tackle visual relation detection using a weaker form of label, i.e. *image-level predicate annotations* only, which reduces data collection cost, is more robust to non-exhaustive annotations, and helps generalization w.r.t. rare/unseen triplets.
- II) We propose a novel explanation-based weakly-supervised approach for relationship detection. We believe this is the first work to (a) use weakly-supervised learning beyond object/scene recognition, and (b) employ explanation techniques on graph networks as the *key component* of a relationship detection pipeline.
- III) Despite using weaker supervision, we show comparable results to state-of-the-art methods with stronger labels on several visual relation benchmarks.

2 Related Works

We are interested in weakly-supervised learning of visual relations. We achieve this by employing graph network explanation techniques. In this section, we cover the related papers corresponding to the different aspects of our work.

Visual Relationship Detection. Visual relation detection involves identifying groups of objects that exhibit semantic relations, in particular (subject, predicate, object) triplets. Relations are usually either comparative attributes/relative spatial configurations [12] which are useful for referral expression [26] and visual question answering [17], or, inter-object interactions [39] which is crucial for scene understanding. Due to the importance of human-centered image recognition for various applications, many of such works focus on human-object interactions [46,7,6,34,15,51].

Visual relation detection has been initially tackled by considering the whole relationship triplet as a single-phrase entity [39]. However, this approach comes with high computational costs and data inefficiency due to the combinatorial space of possible phrases. It is therefore important to devise methods that improve data efficiency and better generalize to rare or unseen relations.

Most modern works take a compositional approach [24,31,30,34,15,29], where objects and predicates are modelled in their own right, which enables better and more efficient generalization. Leveraging language through construction of priors, text embeddings, or joint textual-visual embeddings has also been shown to improve generalization [24,31,29]. The recent work of Peyre *et al.* [29] deals with the combinatorial growth of relation triplets using visual-language analogies. While this approach generalizes well to unseen combinations of seen entities, it adopts a fully-supervised training procedure that demands a considerable amount of annotated triplets for training.

In contrast, our approach improves data efficiency by only requiring image-level predicate labels, and instead learning relation triplets through weakly-supervised learning. Our non-reliance on the subject/object entities, in turn, improves generalization to unseen relations as, importantly, we do not require subject/object entities to appear in the training set.

Weakly-Supervised Learning. Weakly-supervised learning is generally desirable since it reduces the need for costly annotations. It has already proven effective for various visual recognition tasks including object detection [28,5], semantic segmentation [10,20], and instance segmentation [52,14]. Relationship detection can benefit from weakly-supervised learning even more than object/scene recognition, since the number of possible relation triplets grows quadratically with the number object categories. Despite this, weakly-supervised learning of visual relations has received surprisingly less attention than object-centric tasks.

Weakly-Supervised Learning of Visual Relations. The early work of Prest *et al.* [33], similar to our work, only requires image-level action labels. But Prest *et al.* focused on human-object interactions using part detectors, as opposed to general visual relationship detection. More recent works [30,48] learn visual relations in a weakly-supervised setup where triplets are annotated at the image level and not localized through bounding boxes. Peyre *et al.* [30] repre-

sents object pairs by their individual appearance as well as their relative spatial configuration. Then, they use discriminative clustering with validity constraints to assign object pairs to image-level labels. In [48], three separate pipelines are used, one for object detection, one for object-object relation classification and the third for object-object pair selection for each relation. The softmax output of the latter is then used as an attention mechanism over object pairs to account for the weak labels.

Both [48,30] work with non-localized triplets annotated at the image-level². Our weaker supervision setup, by not requiring subject and object annotations, allows for potentially simpler, more general, and less costly construction of large training datasets using search engines or image captions. Furthermore, our method is based on object-centric explanations of graph networks, which sets it apart from previous works on weakly-supervised learning of visual relations.

Explanation Techniques. In mission-critical applications such as medical prognosis, a real-world deployment of trained AI systems require explanations of the predictions. Thus, many explanation techniques have been developed based on local approximation [37], game theory [25], or gradient propagation [2,50,41]. Recently, following the success of graph networks, explanation methods have been extended to those models as well [32,4,47]. We use graph networks to obtain image-level predicate predictions and then apply graph explanation techniques to obtain the corresponding subject and object in an unsupervised manner.

Explanation-based weakly-supervised learning. The idea of using explanations to account for weak labels has been previously used for object recognition. Class Activation Mapping (CAM) uses a specific architecture with fully-convolutional layers and global average pooling to obtain object localization at the average pooling layer [50]. [52] extends this approach by backpropagating the maximum response of the CAM back to the image space for weakly-supervised instance segmentation. Grad-CAM [41] generalizes CAM and extends its applicability to a wider range of architectures by pushing the half-rectified gradient backward and using channel-wise average pooling to obtain location-wise importance. Similar to CAM, Grad-CAM is applied to ILSVRC [38] for weakly-supervised object localization. Finally, [14] develops a cascaded label propagation setup with conditional random fields and object proposals to obtain object instance segmentation from image-level predictions, using excitation back-propagation [49] for the backward pass. Our work is an extension to this line of research: we consider a more complicated application, namely visual relationship detection, and use explanation techniques on graph networks.

3 Method

Detecting visual relationships in an image consists in identifying triplets $\tau = \langle \text{subj}, \text{pred}, \text{obj} \rangle$ of subject, predicate and object. For example, *person drive car* or *tree next to building*. To formalize this, we denote the set of objects in an

² It should be noted that [30] can be extended to work with only predicate annotations, using a new set of more relaxed constraints.

image by \mathcal{O} , where each object instance, i , has a corresponding bounding box b_i and is categorized as c_i according to a vocabulary of object classes $\{1 \dots C\}$. Predicates belong to a vocabulary of predicate classes $\{1 \dots K\}$ that include actions such as *eating*, spatial relations such as *next to* and comparative terms such as *taller than*.

With this notation, detecting visual relations from an image \mathcal{I} corresponds to determining high-density regions of the following joint probability distribution:

$$P(\tau|\mathcal{I}) \triangleq P(c_{\text{subj}} = c_i, k_{\text{pred}} = k, c_{\text{obj}} = c_j, b_{\text{subj}} = b_i, b_{\text{obj}} = b_j | \mathcal{I}), \quad (1)$$

where c_{subj} and c_{obj} indicate resp. the class of the subject and the object, k_{pred} indicates the class of the predicate, b_{subj} and b_{obj} indicate resp. the location of the subject and the object, and $i, j = 1 \dots |\mathcal{B}|$ index the bounding boxes.

To accommodate weakly-supervised learning, we propose the following approximate factorization based on object detection and predicate classification:

$$\begin{aligned}
 P(\tau|\mathcal{I}) &= \\
 P(c_{\text{subj}} = c_i | \mathcal{I}, b_{\text{subj}} = b_i) P(c_{\text{obj}} = c_j | \mathcal{I}, b_{\text{obj}} = b_j) & \quad \text{object detection} \quad (2) \\
 P(k_{\text{pred}} = k | \mathcal{I}) & \quad \text{predicate classification} \quad (3) \\
 P(b_{\text{subj}} = b_i, b_{\text{obj}} = b_j | \mathcal{I}, k_{\text{pred}} = k) & \quad \text{likelihood of a pair} \quad (4) \\
 P(c_{\text{subj}} = c_i, c_{\text{obj}} = c_j | k_{\text{pred}} = k). & \quad \text{prior over relations} \quad (5)
 \end{aligned}$$

For equation 2, we use an object detection pipeline to localize and classify objects in an image. The two terms, then, refer to the confidence scores assigned by the object detector to the subject and object of the relationship (section 3.1).

Equation 3 corresponds to a predicate classifier that predicts the presence of predicate k in the image. This component only relies on image-level predicate annotations during training, and does not explicitly attribute its predictions to pairs of input objects. However, by carefully designing the architecture of the predicate classifier, we introduce a strong inductive bias towards objects and relations, which we can later exploit to recover $\langle \text{subj}, \text{pred}, \text{obj} \rangle$ triplets (sec. 3.2).

Given a certain predicate k , equation 4 recovers the likelihood of object pairs to be the semantic subject and object of that predicate. In other words, we wish to identify *all* possible $(\text{subj}, \text{obj})$ pairs by their likelihood equation 4 w.r.t. a given predicate. Therefore, we use an explanation technique to compute unnormalized scores that associate predicates to pairs of objects (section 3.3).

Term 5, which we refer to as *prior over relationships*, represents the co-occurrence of certain classes as subjects or objects of a predicate, and the directionality of such relationship. For instance, it can indicate that *(person, truck)*, with *person* as the subject, is a more likely pair for *drive* than *(fork, sandwich)*. As such, this term is optional, and excluding it would be the same as assuming a uniform prior. However, this term assumes great importance in a weakly-supervised setup, since isolated predicate labels provide no clue on the directionality of the relation between subject and object (section 3.4).

3.1 Object detection

We use an object detection module to extract a set of objects \mathcal{O} from a given image \mathcal{I} . We describe each object bounding box by the visual appearance features and the classification scores obtained from the detector. These objects will then be used to classify the predicates present in \mathcal{I} and, later on, serve as targets for explanations that identify relevant relationship triplets. Similar to the weakly-supervised setup of Peyre *et al.* [30] we assume the availability of pre-trained object detectors [36] as there is substantial progress in that field.

3.2 Predicate classification

Predicate classification as described in equation 3 is a mapping from image to predicate(s) and as such does not necessarily require an understanding of objects. Thus, a simple choice for the classifier would be a convolutional neural network (CNN) trained on image-level predicate labels, e.g. ResNeXt [44]. However, the raw representation of images as pixels does not explicitly capture the compositional nature of the task. Instead, we introduce a strong inductive bias towards objects and relationships in both the data representation and the architecture. Specifically, the module is implemented as a graph neural network (GNN) with architecture similar to [40], that takes as input a graph representation of the image $\mathcal{G} = (\mathcal{O}, \mathcal{E})$, aggregates information by passing messages over the graph, and produces image-level predicate predictions. This design choice allows us to later explain the predictions in terms of objects, rather than raw pixels.

Each node in the image graph represents an object $i \in \mathcal{O}$ with its spatial and visual features extracted by the object detector, which together we denote as the tuple $\mathbf{n}_i = (\mathbf{n}_i^s, \mathbf{n}_i^v)$. The image graph is built as fully-connected and therefore impartial to relations between objects. Directed edges $i \rightarrow j$ are placed between every pair of nodes, excluding self loops, resulting in $|\mathcal{O}|^2 - |\mathcal{O}|$ edges.

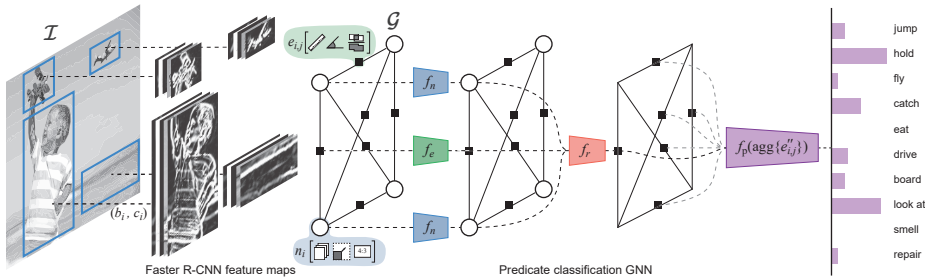


Fig. 2: **A graph neural network (GNN) trained to classify the predicates depicted in a scene.** Object detections extracted through Faster R-CNN are represented as a fully-connected graph. The GNN classifier aggregates local information across nodes and produces an image-level predicate prediction. The input representation and architecture implicitly encode an inductive bias for pairwise relationships

Node \mathbf{n}_i and edge $\mathbf{e}_{i,j}$ representations are first transformed through two small networks f_n and f_e :

$$\mathbf{n}'_i = f_n(\mathbf{n}_i) \quad (6)$$

$$\mathbf{e}'_{i,j} = f_e(\mathbf{e}_{i,j}). \quad (7)$$

Then, a relational function f_r aggregates local information by considering pairs of nodes and the edge connecting them:

$$\mathbf{e}''_{i,j} = f_r(\mathbf{n}'_i, \mathbf{e}'_{i,j}, \mathbf{n}'_j). \quad (8)$$

This pairwise function induces an architectural bias towards object-object relationships, which hints at the ultimate goal of relationship detection.

In a fully-supervised scenario, a classification head could be applied to each of the $\mathbf{e}''_{i,j}$ edges and separate predicate classification losses could be computed using ground-truth pairwise labels $\mathbf{p}_{i,j}$, e.g. [34]. Instead, we consider image-level labels $\mathbf{p} \in \{0, 1\}^K$, where p_k indicates the presence of predicate k in the image, e.g. \mathbf{p} would contain 1s at the locations of *push*, *wear*, *drive* for figure 1. Therefore, we aggregate all edge vectors and apply a final prediction function that outputs a binary probability distribution over predicates as in equation 3:

$$\mathbf{y} = f_p(\text{agg}\{\mathbf{e}''_{i,j}\}) \in [0, 1]^K, \quad (9)$$

where *agg* is a permutation-invariant pooling function such as *max*, *sum* or *mean*.

Designed as such, the graph-based predicate classifier can be trained by minimizing the binary cross entropy between predictions and ground-truth labels:

$$\mathcal{L} = - \sum_{k=1}^K \{p_k \log(y_k) + (1 - p_k) \log(1 - y_k)\}. \quad (10)$$

3.3 Explanation-based relationship detection

Once the predicate classifier is trained, we wish to use it to detect complete relationship triplets $\langle \text{subj}, \text{pred}, \text{obj} \rangle$. This is where the relational inductive bias introduced for the predicate classifier plays a key role. In fact, had the predicate classifier been a simple CNN, we would only be able to attribute its predictions to the input pixels, e.g. through sensitivity analysis [3] or Grad-CAM [41]. Figure 3 shows an example of Grad-CAM explanations obtained for a ResNeXt architecture [44] trained for predicate classification on the Visual Relationship Detection dataset (see appendix B.3). While it is possible to guess which areas of the image are relevant for the predicted predicate, it is undoubtedly hard to identify a distinct (subj, obj) pair from the pixel-wise heatmaps.

Thanks to the GNN architecture of the previous module, we can instead attribute predicate predictions to the nodes of the input graph, evaluating the importance of *objects* rather than *pixels*. We can then consider all pairs of nodes representing the candidate subject and object of a predicate of interest, score them with a backward explanation procedure and select the top-ranking triplets.



Fig. 3: **Grad-CAM heatmap visualization of a ResNet predicate classifier.** Ground-truth annotations contain *person wear jacket* and *person above snowboard*, but it would be hard to identify subjects and objects from the pixel-level explanation.

Specifically, we apply *sensitivity analysis* [3] to compute the relevance of a node (r_i^k) and of an edge ($r_{i,j}^k$) with respect to a predicate k :

$$r_i^k = \left\| \frac{\partial y_k}{\partial \mathbf{n}_i} \right\|_1 \quad \text{single-object relevance} \quad (11)$$

$$r_{i,j}^k = \left\| \frac{\partial y_k}{\partial e_{i,j}} \right\|_1 \quad \text{object-pair relevance} \quad (12)$$

We experimented with different ways to compute these relevances, including $\text{gradient} \times \text{input}$, $\max(\text{gradient} \times \text{input}, 0)$, and the L1, L2 norms, but no significant differences were noticed on the validation set.

The product of these relevances is then used as a proxy for the unnormalized likelihood of a subject-object pair given a predicate (equation 4):

$$P(b_{\text{subj}} = b_i, b_{\text{obj}} = b_j | k_{\text{pred}} = k) \propto r_i^k \cdot r_{i,j}^k \cdot r_j^k. \quad (13)$$

Rather than computing this quantity for every predicate and for every pair of objects, we limit the search to the N top-scoring predicates, reducing the number of candidates from $K(|\mathcal{O}|^2 - |\mathcal{O}|)$ to $N(|\mathcal{O}|^2 - |\mathcal{O}|)$ relationships. A Big O complexity that scales as $|\mathcal{O}|^2$ might seem unappealing, yet with $|\mathcal{O}| < 30$ we could process batches of 128 image graphs in a single pass.

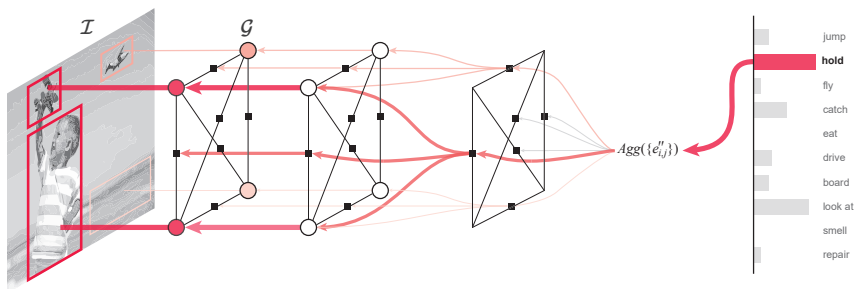


Fig. 4: **Relationship detection through explanation.** A predicate prediction is *explained* by attributing it to the pair of objects in the input that are most relevant for it, effectively recovering a full relationship triplet in the form $\langle \text{subj}, \text{pred}, \text{obj} \rangle$

3.4 Prior over relationships

Learning to detect $\langle \text{subj}, \text{pred}, \text{obj} \rangle$ relations using image-level predicate labels is inherently ill-defined. Consider the task of learning a new predicate, e.g. *squanch*. By observing a sufficient number of labeled images, we could learn that two specific objects are often in a *squanch* relationship. However, we would not be able to determine which should be the subject and which the object, i.e. the direction of such relation, without semantic knowledge about the new word (can things be *squanchier than* others? can objects *squanch* each other?).

Equation 5 represents the belief over which categories can act as subject and objects of a certain predicate. In fully- or weakly-supervised scenarios, where $\langle \text{subj}, \text{pred}, \text{obj} \rangle$ triplets are available during training, a relationship detector would learn such biases directly from data. Our graph-based predicate classifier, trained only image-level predicate annotations, can indeed learn to recognize object-object relations and to assign high probability to meaningful pairs (equation 13), but neither the training signal nor the inductive biases contain hints about directionality. In fact, the relevance $r_{i,j}^k$ is in no way constrained to represent the relationship that has i as subject and j as object, even though equation 8 considers the edge $i \rightarrow j$. Thus the explanation (equation 13) for *hold* might score both *person hold pencil* and its semantic opposite *pencil hold person* equally.

Previous work [24] use `word2vec` [27] embeddings of $\langle \text{subj}, \text{pred}, \text{obj} \rangle$ triplets from the training set to form a semantically-grounded prior. Instead, we compute a simple frequency-based prior $\text{freq}(c_i, c_j | k)$ over a small validation set, to avoid including exclusive relationship information from the training set (app. C.3).

4 Experiments

In this section, we test our weakly-supervised method for visual relationship detection on three different datasets, each one presenting specific challenges and different evaluation metrics. Before discussing the individual experiments, we provide further implementation details about the object detection, predicate classification and visual relationship explainer modules. Additional experiments and ablation studies can be found in appendix C.

4.1 Setup

Object detection. Our object detection module is based on the `detectron2` [43] implementation of Faster R-CNN [36]. Given an object i and its bounding box b_i , either from the ground-truth annotations or detected by Faster R-CNN, we use `ROIALIGN` [16] to pool a $256 \times 7 \times 7$ feature volume \mathbf{n}_i^v from the pyramid of features [22] built on top of a ResNeXt-101 backbone [44]. Furthermore, we compute a feature vector \mathbf{n}_i^s that represents the spatial configuration of b_i . Specifically, the tuple of spatial and visual features $\mathbf{n}_i = (\mathbf{n}_i^s, \mathbf{n}_i^v)$ is defined as:

$$\mathbf{n}_i^s = \left[\frac{w_i}{h_i}, \frac{h_i}{w_i}, \frac{w_i h_i}{WH} \right] \quad \text{spatial features} \quad (14)$$

$$\mathbf{n}_i^v = \text{ROIALIGN}(\text{FPN}(\mathcal{I}), b_i), \quad \text{visual features} \quad (15)$$

where (w_i, h_i) and (W, H) represent width and height of the box b_i and of the image \mathcal{I} respectively, FPN is the feature pyramid network used to extract visual features from the whole image, and ROIALIGN is the pooling operation applied to the feature pyramid to extract features relative to the box b_i .

Edge attributes $e_{i,j}$ are chosen to represent the spatial configuration of the pair of objects they connect:

$$e_{i,j} = \left[\frac{\|\mathbf{x}_j - \mathbf{x}_i\|}{\sqrt{WH}}, \sin(\angle_{ij}), \cos(\angle_{ij}), \text{IoU}(b_j, b_i), \frac{w_j h_j}{w_i h_i} \right], \quad (16)$$

where $\mathbf{x}_i \in \mathbb{R}_+^2$ is the center of b_i , \angle_{ij} is the angle between $\mathbf{x}_j - \mathbf{x}_i$ and the positive horizontal axis, and IOU is the intersection over union of the two boxes.

Predicate classifier. At training time, the input of the predicate classifier described in sec. 3.2 is a fully-connected graph of ground-truth objects. During inference, we apply the object detector and build a graph with all objects having confidence score of 30% or more. For each dataset, the hyperparameters of the GNN-based predicate classifier are selected on a validation split of 15% training images. The following values apply to the HICO-DET dataset, more details about the hyperparameter space are available in appendix B.2.

The input node function f_n is implemented as i) a $2 \times (\text{CONV} + \text{RELU})$ network applied to \mathbf{n}^v , where the convolutional layers employ 256 kernels of size 3×3 each, and ii) a $\text{LINEAR} + \text{RELU}$ operation that transforms \mathbf{n}^s into a 1024-vector. The input edge functions f_e consist of a $\text{LINEAR} + \text{RELU}$ operation that outputs a 1024-vector of transformed edge features. The relational function f_r in equation 8 is implemented as a $\text{LINEAR} + \text{RELU}$ operation where the features of two nodes and of the directed edge between them are concatenated at the input. The output of f_r is a 1024-vector for each ordered pair of nodes. For all datasets, the aggregation function in equation 9 is element-wise \max , and f_p is a $\text{LINEAR} + \text{SIGMOID}$ operation that outputs a K -vector of binary probabilities.

We train the weights of the predicate classifier by minimizing the loss in equation 10 with the Adam optimizer [19] with 10^{-3} initial learning rate and 10^{-5} weight decay. During training, we track `recall@5`, i.e. the fraction of ground-truth predicates retrieved among the top-5 confident predictions for an image. We let the optimization run on batches of 128 graphs for 18 epochs, at which point the classifier achieves 94% recall on a validation split.

Relationship detector. The explanation-based relationship detection algorithm described in section 3.3 does not have many hyperparameters. We tried i) whether to multiply the gradient with the input when computing relevances, ii) which norm to use between L1, L2 and $\max(\text{L1}, 0)$, and iii) the number N of top-scoring predicates whose gradient is traced back to the input to identify relevant triplets. As observed in [8], optimizing these parameters on the whole training set would violate the premise of weakly-supervised learning by accessing fully-labeled data. Therefore, we employ once again the 15% validation split used to optimize the classifier, assuming that in a real-world scenario it should always be possible to manually annotate a small subset of images for validation purposes. The best choice of N for all datasets was found to be 10, while the other two parameters seem to have little effect on performance.

4.2 HICO-DET

The Humans Interacting with Common Objects (HICO-DET) dataset contains $\sim 50K$ exhaustively annotated images of *human-object interactions* (HOI), split into $\sim 40K$ train and $\sim 10K$ test images [7,6]. The subject of a relationship is always *person*, the 117 predicates cover a variety of human-centric actions (e.g. *cook*, *wash*, *paint*), and the 80 objects categories are those defined as `thing_classes` in MS-COCO [23]. We can therefore use the pre-trained object detector from [43], of which we report performances in appendix A.1.

The nature of this dataset allows us to embed the relationship prior in the graph itself. A fully-connected graph encodes a uniform prior, i.e. no preference about subject-object pairs, while a sparse graph containing only edges from humans to objects encodes a bias towards *human-object interactions*.

The metric for this dataset is the 11-point interpolated mean Average Precision (mAP) [11] computed over the 600 human-object interaction classes of the dataset [6]. The following criteria should be met for a detected triplet to match with a ground-truth triplet: a) subject, predicate and object categories match, and b) subject boxes overlap with $\text{IoU} > .5$, and c) object boxes overlap with $\text{IoU} > .5$, and d) the ground-truth triplet was not matched with a previously-considered detected triplet. Table 1 reports mAP for the standard splits of HICO-DET[6]: all 600 human-object interactions, 138 rare triplets, and 462 non-rare triplets (10 or more training samples). We compare with various fully-supervised baselines including the original HO-RCNN from [6] and the method from [29] that uses semantic and visual analogies to improve detection of rare and unseen triplets. Despite the weaker supervision signal, the strong inductive bias towards pairwise relationships allows our explanation method to achieve higher mAP for both the uniform and human-object priors.

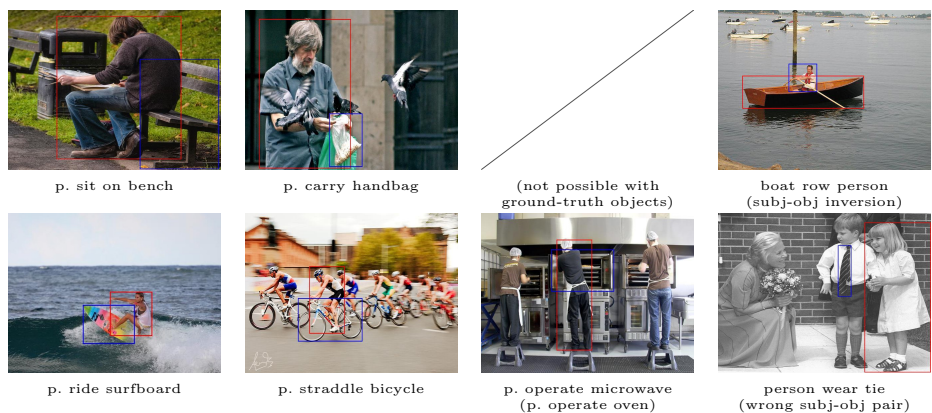


Fig. 5: **Relationship detection on HICO-DET.** Top row uses GT objects, bottom row uses Faster R-CNN objects. Left to right: correct relationship detection, correct but missing ground-truth, incorrect due to object misdetection, incorrect detection (selected predictions of our model using a uniform relationship prior)

Table 1: **Mean Average Precision on the HICO-DET dataset.** The choice of relationship prior embedded in the graph is indicated in parentheses

	Full (600)	Rare (138)	Non-rare (462)
Fully supervised			
Chao [6]	7.81	5.37	8.54
InteractNet [15]	9.94	7.16	10.77
GPNN [34]	13.11	9.34	14.23
iCAN [13]	14.84	10.45	16.15
Analogies [29]	19.40	14.60	20.90
Weakly supervised			
Ours (uniform)	24.25	20.23	25.45
Ours (human-object)	28.04	24.63	29.06

4.3 Visual Relationship Detection dataset

The Visual Relationship Detection dataset (VRD) contains ~ 5000 annotated images, split into ~ 4000 train and ~ 1000 test images [18,24]. The 70 predicates in this dataset include both verbs and spatial relationships, e.g. *carry*, *next to*. The 100 object categories cover both well spatially-defined objects such as *bottle* and concepts like *sky* and *road*, that are harder to localize. For this set of objects there is no ready-to-use object detector, therefore we finetune a **detectron2** model using annotations from the training set (details in appendix A.2).

The standard metric for VRD [24] is **recall@x** i.e. fraction of ground-truth triplets retrieved among the x top-ranked detections [1]. Here, recall is preferred over mAP since it does not penalize the retrieval of triplets that exist in the image, but are missing in the ground-truth. Criteria for true positive in VRD follow those of HICO-DET, and are used in the following settings [24]:

Predicate detection: objects for the image graph come from ground-truth annotations, allowing to test the explanation-based retrieval of relationships under perfect object detection conditions (classification and localization).

Phrase detection: objects come from Faster R-CNN proposals, but $\text{IoU} > .5$ is evaluated on the union box of subject and object, effectively localizing the entire relationship as a single image region, or *visual phrase* [39].

Relationship detection: objects come from Faster R-CNN proposals, subject and object boxes are required to individually overlap with their corresponding boxes in the ground-truth (same as HICO-DET).

As shown in table 2, our method achieves recall scores $R@100$ close to a fully-supervised baseline [24], despite the weaker training signal. By analyzing the top 100 predictions of a model with uniform prior, we often observed the co-appearance of a relationship and its semantic opposite, e.g. *person drive car* and *car drive person*, which possibly “wastes” half of the top-x detection due to incorrect directionality (corroborated by the gap between $R@50$ and $R@100$ of ours-uniform). Importantly, moving from a uniform to a frequency-based prior almost doubles $R@50$, which highlights the importance of the relationship prior in connection with our method. We expect that including a stronger prior, e.g.

Table 2: **Recall at 50 and 100 on the VRD dataset.** Comparison of fully- and weakly-supervised methods. The choice of relationship prior is indicated in parentheses

	GT objects		R-CNN objects			
	Predicate det. R@50	Phrase det. R@100	Phrase det. R@50	Relation. det. R@100	Relation. det. R@50	Relation. det. R@100
Fully supervised						
Visual Phrases [39]	0.9	1.9	0.04	0.07	-	-
Visual [24]	3.5	3.5	0.7	0.8	1.0	1.1
Visual+Language [24]	47.9	47.9	16.2	17.0	13.9	14.7
Sup. PPR-FCN [48]	47.4	47.4	19.6	23.2	14.4	15.7
Peyre [30]	52.6	52.6	17.9	19.5	15.8	17.1
Weakly sup. (subj,pred,obj)						
PPR-FCN [48]	-	-	6.9	8.2	5.9	6.3
Peyre [30]	46.8	46.8	16.0	17.4	14.1	15.3
Weakly sup. (pred only)						
Ours (uniform)	27.3	47.1	6.8	13.0	5.3	8.4
Ours (frequentist)	43.0	57.4	14.8	20.2	10.6	13.2

based on natural-language embeddings of objects and predicates, would further improve detection of semantically-correct relationships.

The test set of VRD contains a triplets that never occur during training and can be used to evaluate zero-shot generalization. As shown in table 3, our method performs on a par with other methods that use stronger annotations and explicitly improve generalization through language embeddings [24] or visual analogy transformations [29]. Expectedly, the freq-based prior computed on the validation set does not improve recall of unseen triplets. To verify importance of this term, we show that a simple prior with access to a few zero-shot triplets readily improves recall. Clearly, peeking at the test set is not correct practice, but serves as a proxy for what could be achieved by improving this term, e.g. via incorporating language or visual analogies. The next experiment better demonstrates the generalization of our method to unseen triplets.

4.4 Unusual Relations dataset

The Unusual Relations dataset (UnRel) is an evaluation-only collection of ~ 1000 images, which shares the same vocabulary as VRD and depicts rarely-occurring relationships [30]. For relationship detection methods trained on $\langle \text{subj}, \text{pred}, \text{obj} \rangle$ triplets, this dataset represents a benchmark for zero-shot retrieval of triplets not seen during training. E.g. our predicate classifier trained on VRD has clearly encountered *hold* during training, but never in *person hold plane* (figure 2).

In table 4 we report mAP over the 76 unusual triplets of UnRel. We follow the evaluation setup of [30]: the test set of VRD is mixed in to act as distractor, up to 500 candidate triplets per image are retained, and they are matched if $\text{IoU} > .3$. Since the average number of detected objects per image is small, ~ 4 , we increase the number of top-scoring predicates considered in the explanation step to $N = 50$. Differently from [30,29], we use obj. detection scores for ranking triplets, and we do not introduce a *no-interaction* predicate. Compared to recall, mAP is less affected by unseen triplets and the prior from VRD results effective.

Table 3: **Zero-shot recall on the VRD dataset:** triplets from the test set that are never seen during training. The choice of relationship prior is indicated in parentheses

	GT objects		R-CNN objects			
	Predicate det.		Phrase det.		Relation. det.	
	R@50	R@100	R@50	R@100	R@50	R@100
Fully supervised						
Visual [24]	3.5	3.5	0.7	0.8	1.0	1.1
Visual+Language [24]	8.5	8.5	3.4	3.8	3.1	3.5
Peyre 2017 [30]	21.6	21.6	6.8	7.8	6.4	7.4
Weakly sup. (subj,pred,obj)						
Peyre 2017 [30]	19.0	19.0	6.9	7.4	6.7	7.1
Weakly sup. (pred only)						
Ours (uniform)	13.7	29.2	3.8	6.5	2.8	4.6
Ours (VRD freq.)	13.5	28.2	4.4	6.4	3.3	4.6
Ours (Zero freq.)	20.5	37.0	4.7	8.2	4.0	6.4

Table 4: **Mean Average Precision on UnRel with VRD as a distractor**

	GT objects		R-CNN objects	
	Predicate	Phrase	Subj. only	Relationship
Fully supervised				
Peyre 2017 [30]	62.6	14.1	12.1	9.9
Analogies [29]	63.9	17.5	15.9	13.4
Weakly sup. (subj,pred,obj)				
Peyre 2017 [30]	58.5	13.4	11.0	8.7
Weakly sup. (pred only)				
Ours (uniform)	70.9	19.8	18.1	14.9
Ours (frequency)	70.6	20.0	18.3	15.1

5 Conclusion

We considered the task of learning visual relationship detection with weak image-level predicate labels. While this makes learning significantly harder, it enables collecting datasets that are more representative of possible relations without suffering from combinatorial scaling of search queries and annotation cost.

Using pretrained object detectors, strong inductive bias via graph networks, backward explanations, and a direction prior, we showed that it is possible to achieve results on par with recent works that benefit from stronger supervision.

An issue with predicate-only annotation is the lack of directional information, which can only be provided using auxiliary sources such as language. While we mitigated this issue through a simple frequentist prior, an important future direction is to solve it in a principled way. For instance, one can annotate a subset of images with unlocalized image-level triplets, only to disambiguate the direction of the relations. Note that, since such a dataset does not have to be exhaustively annotated for all triplets, the collection cost would be negligible.

Finally, another interesting direction is to study the proposed explanation-based weakly-supervised method in other domains such as situation recognition [21], video recognition [45], segmentation [35], chemistry [9] and biology [42].

Acknowledgements. Funded by Swedish Research Council project 2017-04609 and by Wallenberg AI, Autonomous Systems and Software Program (WASP).

References

1. Alexe, B., Deselaers, T., Ferrari, V.: Measuring the objectness of image windows. *IEEE transactions on pattern analysis and machine intelligence* **34**(11), 2189–2202 (2012)
2. Bach, S., Binder, A., Montavon, G., Klauschen, F., Müller, K.R., Samek, W.: On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PloS one* **10**(7) (2015)
3. Baehrens, D., Schroeter, T., Harmeling, S., Kawanabe, M., Hansen, K., Müller, K.R.: How to explain individual classification decisions. *Journal of Machine Learning Research* **11**(Jun), 1803–1831 (2010)
4. Baldassarre, F., Azizpour, H.: Explainability techniques for graph convolutional networks. In: *International Conference on Machine Learning (ICML) Workshops, 2019 Workshop on Learning and Reasoning with Graph-Structured Representations* (2019)
5. Bilen, H., Vedaldi, A.: Weakly supervised deep detection networks. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (June 2016)
6. Chao, Y.W., Liu, Y., Liu, X., Zeng, H., Deng, J.: Learning to detect human-object interactions. In: *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*. pp. 381–389. IEEE (2018)
7. Chao, Y.W., Wang, Z., He, Y., Wang, J., Deng, J.: Hico: A benchmark for recognizing human-object interactions in images. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 1017–1025 (2015)
8. Choe, J., Oh, S.J., Lee, S., Chun, S., Akata, Z., Shim, H.: Evaluating weakly supervised object localization methods right. *arXiv preprint arXiv:2001.07437* (2020)
9. Do, K., Tran, T., Venkatesh, S.: Graph transformation policy network for chemical reaction prediction. In: *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. pp. 750–760 (2019)
10. Durand, T., Mordan, T., Thome, N., Cord, M.: Wildcat: Weakly supervised learning of deep convnets for image classification, pointwise localization and segmentation. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 642–651 (2017)
11. Everingham, M., Van Gool, L., Williams, C.K., Winn, J., Zisserman, A.: The pascal visual object classes (voc) challenge. *International journal of computer vision* **88**(2), 303–338 (2010)
12. Galleguillos, C., Rabinovich, A., Belongie, S.: Object categorization using co-occurrence, location and appearance. In: *2008 IEEE Conference on Computer Vision and Pattern Recognition*. pp. 1–8. IEEE (2008)
13. Gao, C., Zou, Y., Huang, J.B.: ican: Instance-centric attention network for human-object interaction detection (2018)
14. Ge, W., Guo, S., Huang, W., Scott, M.R.: Label-penet: Sequential label propagation and enhancement networks for weakly supervised instance segmentation. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 3345–3354 (2019)
15. Gkioxari, G., Girshick, R., Dollár, P., He, K.: Detecting and recognizing human-object interactions. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 8359–8367 (2018)
16. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: *Proceedings of the IEEE international conference on computer vision*. pp. 2961–2969 (2017)

17. Johnson, J., Hariharan, B., van der Maaten, L., Fei-Fei, L., Lawrence Zitnick, C., Girshick, R.: Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2901–2910 (2017)
18. Johnson, J., Krishna, R., Stark, M., Li, L.J., Shamma, D., Bernstein, M., Fei-Fei, L.: Image retrieval using scene graphs. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3668–3678 (2015)
19. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
20. Lee, J., Kim, E., Lee, S., Lee, J., Yoon, S.: Ficklenet: Weakly and semi-supervised semantic image segmentation using stochastic inference. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 5267–5276 (2019)
21. Li, R., Tapaswi, M., Liao, R., Jia, J., Urtasun, R., Fidler, S.: Situation recognition with graph neural networks. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 4173–4182 (2017)
22. Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2117–2125 (2017)
23. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: European conference on computer vision. pp. 740–755. Springer (2014)
24. Lu, C., Krishna, R., Bernstein, M., Fei-Fei, L.: Visual relationship detection with language priors. In: European conference on computer vision. pp. 852–869. Springer (2016)
25. Lundberg, S.M., Lee, S.I.: A unified approach to interpreting model predictions. In: Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R. (eds.) Advances in Neural Information Processing Systems 30, pp. 4765–4774. Curran Associates, Inc. (2017)
26. Mao, J., Huang, J., Toshev, A., Camburu, O., Yuille, A.L., Murphy, K.: Generation and comprehension of unambiguous object descriptions. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 11–20 (2016)
27. Mikolov, T., Chen, K., Corrado, G., Dean, J.: Efficient estimation of word representations in vector space. arXiv preprint arXiv:1301.3781 (2013)
28. Oquab, M., Bottou, L., Laptev, I., Sivic, J.: Is object localization for free?-weakly-supervised learning with convolutional neural networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 685–694 (2015)
29. Peyre, J., Laptev, I., Schmid, C., Sivic, J.: Detecting unseen visual relations using analogies. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1981–1990 (2019)
30. Peyre, J., Sivic, J., Laptev, I., Schmid, C.: Weakly-supervised learning of visual relations. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 5179–5188 (2017)
31. Plummer, B.A., Mallya, A., Cervantes, C.M., Hockenmaier, J., Lazebnik, S.: Phrase localization and visual relationship detection with comprehensive image-language cues. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1928–1937 (2017)
32. Pope, P.E., Kolouri, S., Rostami, M., Martin, C.E., Hoffmann, H.: Explainability methods for graph convolutional neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 10772–10781 (2019)

33. Prest, A., Schmid, C., Ferrari, V.: Weakly supervised learning of interactions between humans and objects. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **34**(3), 601–614 (2011)
34. Qi, S., Wang, W., Jia, B., Shen, J., Zhu, S.C.: Learning human-object interactions by graph parsing neural networks. In: *Proceedings of the European Conference on Computer Vision (ECCV)*. pp. 401–417 (2018)
35. Qi, X., Liao, R., Jia, J., Fidler, S., Urtasun, R.: 3d graph neural networks for rgb-d semantic segmentation. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 5199–5208 (2017)
36. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. In: *Advances in neural information processing systems*. pp. 91–99 (2015)
37. Ribeiro, M.T., Singh, S., Guestrin, C.: "why should I trust you?": Explaining the predictions of any classifier. In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016*. pp. 1135–1144 (2016)
38. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al.: Imagenet large scale visual recognition challenge. *International journal of computer vision* **115**(3), 211–252 (2015)
39. Sadeghi, M., Farhadi, A.: Recognition using visual phrases. In: *Proceedings of the 2011 IEEE Conference on Computer Vision and Pattern Recognition*. pp. 1745–1752 (2011)
40. Santoro, A., Raposo, D., Barrett, D.G., Malinowski, M., Pascanu, R., Battaglia, P., Lillicrap, T.: A simple neural network module for relational reasoning. In: *Advances in neural information processing systems*. pp. 4967–4976 (2017)
41. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Grad-cam: Visual explanations from deep networks via gradient-based localization. In: *Proceedings of the IEEE international conference on computer vision*. pp. 618–626 (2017)
42. Tsubaki, M., Tomii, K., Sese, J.: Compound–protein interaction prediction with end-to-end learning of neural networks for graphs and sequences. *Bioinformatics* **35**(2), 309–318 (2019)
43. Wu, Y., Kirillov, A., Massa, F., Lo, W.Y., Girshick, R.: Detectron2. <https://github.com/facebookresearch/detectron2> (2019)
44. Xie, S., Girshick, R., Dollár, P., Tu, Z., He, K.: Aggregated residual transformations for deep neural networks. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 1492–1500 (2017)
45. Yan, S., Xiong, Y., Lin, D.: Spatial temporal graph convolutional networks for skeleton-based action recognition. In: *Thirty-second AAAI conference on artificial intelligence* (2018)
46. Yao, B., Fei-Fei, L.: Modeling mutual context of object and human pose in human-object interaction activities. In: *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. pp. 17–24. IEEE (2010)
47. Ying, Z., Bourgeois, D., You, J., Zitnik, M., Leskovec, J.: Gnnexplainer: Generating explanations for graph neural networks. In: *Advances in Neural Information Processing Systems*. pp. 9240–9251 (2019)
48. Zhang, H., Kyaw, Z., Yu, J., Chang, S.F.: Ppr-fcn: Weakly supervised visual relation detection via parallel pairwise r-fcn. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 4233–4241 (2017)

49. Zhang, J., Bargal, S.A., Lin, Z., Brandt, J., Shen, X., Sclaroff, S.: Top-down neural attention by excitation backprop. *International Journal of Computer Vision* **126**(10), 1084–1102 (2018)
50. Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., Torralba, A.: Learning deep features for discriminative localization. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 2921–2929 (2016)
51. Zhou, P., Chi, M.: Relation parsing neural network for human-object interaction detection. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 843–851 (2019)
52. Zhou, Y., Zhu, Y., Ye, Q., Qiu, Q., Jiao, J.: Weakly supervised instance segmentation using class peak response. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 3791–3800 (2018)

Supplementary material

The following pages contain: A more details about the three datasets used in this work, B more details about the architecture of the predicate classifier and hyperparameter optimization, C additional relationship detection experiments and ablation studies, and D additional qualitative results from the three datasets.

A Datasets

Table 5: Comparison of the datasets used in this work

	Number of images		Vocabulary size			Unique triplets	
	Train	Test	Subject	Predicate	Object	Train	Test
HICO-DET [6]	38118	9658	1	117	80	600	600
VRD [24]	4006	1001	100	70	100	6672	2741
UnRel [30]	-	1071	100	70	100	-	76

A.1 HICO-DET

The Humans Interacting with Common Objects dataset [7], in its detection version [6], is available at <http://www-personal.umich.edu/~ywchao/hico>. The subject of the relationships is always a *person*. The object vocabulary is the same as MS-COCO [23]. Its predicates indicate human-object interactions, e.g. *carry*. Some images from MS-COCO are also contained in HICO-DET, but the authors made sure that the test set of HICO-DET has no overlap with MS-COCO. We warn future users to ignore the EXIF rotation tags present on some of the images, in fact all bounding boxes are annotated w.r.t. the non-rotated images. See table 5 for a comparison of dataset and vocabulary size.

We use the pre-trained object detector made available through the `detectron2` implementation [43] of Faster R-CNN [36]. Since the object detector is an important part of visual relationship detection pipelines, we report object detection metrics obtained for this dataset in table 6.



Fig. 6: Ground-truth triplet annotations from the HICO-DET dataset

A.2 Visual Relationship Detection dataset

The Visual Relationship Detection Dataset (VRD) [24] is available at <https://cs.stanford.edu/people/ranjaykrishna/vrd>. Its images and annotations correspond to those in the Scene Graph dataset [18], but the vocabularies of objects and predicates have been carefully curated, e.g. figure 7. We warn future users to ignore the EXIF rotation tags present on some of the images, in fact all bounding boxes are annotated w.r.t. the non-rotated images. Also, we note that for some images the annotation file contains 0 objects and 0 relationships. See table 5 for a comparison of dataset and vocabulary size.

Since no pre-trained model is publicly available for this dataset, we fine-tune an object detector based on `detectron2` [43]. Object detection metrics are reported in table 6 for future reference.



Fig. 7: Ground-truth triplet annotations from the VRD dataset

A.3 Unusual Relationships dataset

The Unusual Relationships dataset (UnRel) [30] is available at <https://www.di.ens.fr/willow/research/unrel>. It is meant as an evaluation-only dataset for rare and unusual relationships, e.g. figure 8. See table 5 for a comparison of dataset and vocabulary size.

Since it shares the same object and predicate vocabulary of VRD, we use the same object detector, of which we report object detection metrics in table 6



Fig. 8: Ground-truth triplet annotations the UnRel dataset

Table 6: Object detection metrics for the datasets used in this work

	Mean Average Precision					
	IoU@[0.5:0.95]	IoU@0.5	IoU@0.75	small	medium	large
HICO-DET [6]	20.2	34.1	20.8	2.3	11.5	29.7
VRD [24]	21.2	35.3	22.6	4.9	14.3	25.0
UnRel [30]	21.0	35.3	22.6	4.9	14.3	25.0

	Mean Average Recall					
	top-1	top-10	top-100	small	medium	large
HICO-DET [6]	30.3	39.3	40.2	11.6	29.2	48.6
VRD [24]	34.0	45.0	45.1	14.9	33.2	48.3
UnRel [30]	34.0	45.0	45.1	14.9	33.2	48.3

B Architecture and hyperparameters

B.1 Introduction to GNNs

In our work, an image is first represented as a fully-connected graph of objects and then processed through a graph neural network to predict predicates. Specifically, we use a message-passing implementation of graph convolution. At the input, each node i is associated to a feature vector \mathbf{v}_i . Similarly, each edge $i \rightarrow j$ is associated to a feature vector $\mathbf{e}_{i,j}$. A global bias term \mathbf{u} can be used to represent information that is not localized to any specific node/edge of the graph. With this graph representation, one layer of message passing performs the following updates.

1. For every edge $i \rightarrow j$, the edge vector is updated using a function f^e that takes as input the adjacent nodes \mathbf{v}_i and \mathbf{v}_j , the edge itself $\mathbf{e}_{i,j}$ and the global attribute \mathbf{u} :

$$\mathbf{e}'_{i,j} = f^e(\mathbf{v}_i, \mathbf{v}_j, \mathbf{e}_{i,j}, \mathbf{u})$$

2. For every node i , features from incident edges $\{\mathbf{e}'_{j,i}\}$ are aggregated using a pooling function $\text{agg}^{e \rightarrow v}$:

$$\bar{\mathbf{e}}'_i = \text{agg}^{e \rightarrow v} \{ \mathbf{e}'_{j,i} \}$$

3. For every node i , the node feature vector is updated using a function f^v that takes as input the aggregated incident edges $\bar{\mathbf{e}}'_i$, the node itself \mathbf{v}_i and the global attribute \mathbf{u} :

$$\mathbf{v}'_i = f^v(\bar{\mathbf{e}}'_i, \mathbf{v}_i, \mathbf{u})$$

4. All edges are aggregated using a pooling function $\text{agg}^{e \rightarrow u}$:

$$\bar{\mathbf{e}}' = \text{agg}^{e \rightarrow u} \{ \mathbf{e}'_{i,j} \}$$

5. All nodes are aggregated using a pooling function $\text{agg}^{v \rightarrow u}$:

$$\bar{\mathbf{v}}' = \text{agg}^{v \rightarrow u} \{ \mathbf{v}'_i \}$$

6. The global feature vector is updated using a function f^u of the aggregated edges $\bar{\mathbf{e}}'$, of the aggregated nodes $\bar{\mathbf{v}}'$ and of the global attribute \mathbf{u} :

$$\mathbf{u}' = f^u(\bar{\mathbf{e}}', \bar{\mathbf{v}}', \mathbf{u})$$

These convolutional layers can be stacked to increase the receptive fields of a node. However, in this work, we used a single layer to focus on pairwise relationships. Furthermore, we did not use a global attribute \mathbf{u} , which could encode for example context and background.

B.2 Predicate classifier

For the predicate classifier we optimize the hyperparameters reported in table 7. Rather than performing a grid-search over the whole space, we perform a "guided" search: we iteratively perform parallel runs and only keep the best-performing combinations of parameters. This process of trial and elimination allows us to quickly prune unpromising regions of the search space.

Table 7: Hyperparameter space of the predicate classifier

Parameter	Choices	Final value
Optimizer		
Learning rate	$10^{-2}, 10^{-3}, 10^{-4}$	10^{-3}
Weight decay	$10^{-3}, 10^{-5}, 0$	10^{-5}
Max epochs	35	18
Model		
Linear layers	1, 2	1
Linear features	256, 512, 1024	1024
Convolutional layers	1, 2	2
Convolutional kernels	256, 512	256
Pooling function	add, max, mean	max
Bias in f_p	yes, no	yes

The best set of hyperparameters is chosen to maximize **recall@5** over a held-out validation set (15% of training data). The train/val split is made at random for every training run. Random seeds are fixed at the beginning of each run and recorded for reproducibility. Note that **recall@5** refers to the image-level predicate predictions, and relationship detection metrics are not involved in the optimization of the predicate classifier.

On the test set of HICO-DET, relative to predicate classification only, these parameters achieve a mAP of 0.44, **recall@5** of 0.90 and **recall@10** of 0.96.

B.3 ResNeXt baseline and Grad-CAM

We finetune a ResNeXt-50 [44] for predicate classification on the Visual Relationship Detection dataset. All parameters are initialized from an ImageNet [38] pretraining, except the final classification layer that is adapted to output 70-dimensional vector of predicate predictions and is initialized from a Normal distribution. Given an input image $\mathcal{I} \in [0, 1]^{3 \times H \times W}$, the convolutional architecture can be summarized as:

$$\mathbf{h} = \text{RESNEXT}(\mathcal{I}) \in \mathbb{R}^{2048 \times \tilde{H} \times \tilde{W}} \quad \text{backbone} \quad (17)$$

$$z_c = \frac{1}{\tilde{H}\tilde{W}} \sum_{i=1}^{\tilde{H}} \sum_{j=1}^{\tilde{W}} h_{c,i,j} \quad \forall c = 1, \dots, 2048 \quad \text{global average pooling} \quad (18)$$

$$\mathbf{y} = \text{softmax}(\mathbf{W}\mathbf{z} + \mathbf{b}) \in [0, 1]^K \quad \text{classification} \quad (19)$$

where \tilde{H} and \tilde{W} represent the height and width of the feature volume extracted by the backbone before global average pooling.

We use Adam optimizer [19] to minimize the same loss of the GNN-based predicate classifier described in the main text. The learning rate is set to 10^{-3} for the classification layer and to 10^{-4} for the rest of the network.

We optimize only the number of epochs and whether the final layer should include a bias term or not. Based on performances on the validation set, the best hyperparameters are training for 6 epochs and including the bias. The final CNN-based model achieves similar **recall@5** as the GNN-based classifier on the test set for predicate classification.

Grad-CAM heatmaps as in figure 3 are produced by computing:

$$\alpha_c^k = \frac{1}{\tilde{H}\tilde{W}} \sum_{i=1}^{\tilde{H}} \sum_{j=1}^{\tilde{W}} \frac{\partial y_k}{\partial h_{c,i,j}} \quad \forall c = 1, \dots, 2048 \quad (20)$$

$$s_{i,j} = \text{RELU} \left(\sum_{c=1}^{2048} \alpha_c^k h_{c,i,j} \right) \quad \forall i = 1, \dots, \tilde{H}; j = 1, \dots, \tilde{W}. \quad (21)$$

Then the 2D vector \mathbf{s} is upsampled to the $H \times W$ size of the input image, and its values are normalized to the range $[0, 1]$.

B.4 Training and inference

The graph neural network described in section 3.2 is trained to classify the predicates present in an image from image-level annotations.

Algorithm 1: Training Algorithm

Input: Pretrained object detector (`detectron2`),
 Dataset of images with image-level predicate annotations.

repeat

- | Extract objects from image \mathcal{I}
- | Build a fully-connected image graph \mathcal{G} using features from eq. 14, 15
- | Apply the predicate classifier to \mathcal{G}
- | Compute the predicate classification loss \mathcal{L} (equation 10)
- | Minimize \mathcal{L} using Adam optimizer

until *convergence*

Output: Trained predicate classifier

Once trained, the predicate classifier can be used for relationship detection. Specifically, each pred prediction is attributed to pairs of objects in the input by means of explanation, thus retrieving the full $\langle \text{subj}, \text{pred}, \text{obj} \rangle$ triplet.

Algorithm 2: Explanation-based Relationship Detection Algorithm

Input: Pretrained object detector (`detectron2`),
 Trained predicate classifier,
 Image of interest \mathcal{I} .

if *Predicate Detection* **then**

- | Extract ground-truth objects from image \mathcal{I}

else if *Phrase Detection* \vee *Relationship Detection* **then**

- | Detect objects in \mathcal{I} using the object detector

end

Build a fully-connected scene graph \mathcal{G} using features from eq. 14, 15

Apply the predicate classifier to \mathcal{G}

Visual relations $\mathcal{R} \leftarrow \emptyset$

for $\text{pred} \in \{N\text{top-scoring predicates}\}$ **do**

- | */* Predicate predictions are explained in terms of*
- | *relevant pairs of objects in the image graph \mathcal{G} */*
- | Compute node and edge relevances using eq. 11, 12
- | Score each $\langle \text{subj}, \text{obj} \rangle$ pair using equation 13
- | Multiply the score by the object detection scores of subj and obj
- | Multiply the score by the classification score of pred
- | Multiply the score by the relationship prior (equation 5)
- | Store high-scoring triplets $\langle \text{subj}, \text{pred}, \text{obj} \rangle$ in \mathcal{R}

end

Output: K top-scoring visual relations from \mathcal{R}

C Additional experiments

C.1 Pooling function

As explained in appendix B, the pooling function for equation 9 is selected according to predicate classification performances (figure 9) on a 15% split of the training set. Figure 9 shows recall@5 for *sum*, *max*, and *mean* pooling over 10 runs on the VRD dataset. Due to higher recall on the validation set, *max* pooling is selected and used for all results reported in the main text. We notice, however, that *mean* pooling performs closely to *max*.

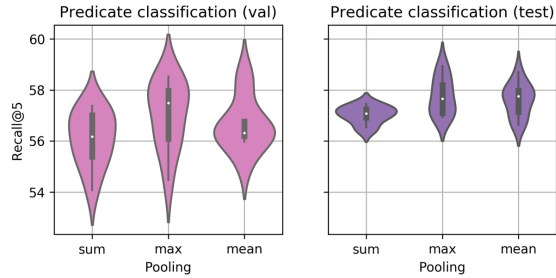


Fig. 9: Recall@5 for predicate classification on VRD using different pooling functions. Validation set (15% of training) on the left, and test set on the right

To further test the role of pooling, we evaluated relationship detection metrics for several predicate classifiers trained using *sum*, *max*, and *mean* pooling. Figure 10 shows that *mean* pooling outperforms the other two, despite performing slightly worse for predicate classification.

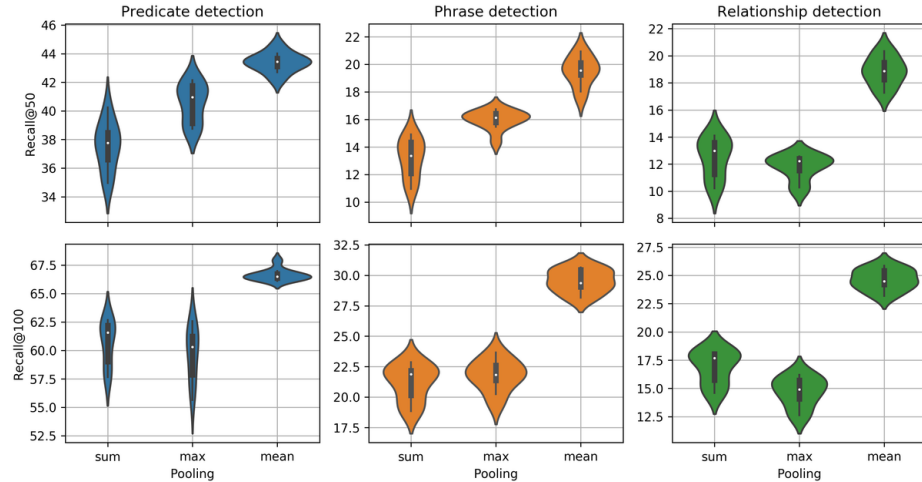


Fig. 10: Recall@50 and $@100$ for relationship detection on VRD using different pooling functions. *mean* pooling outperforms the other two, despite performing slightly worse for predicate classification

C.2 Number of explained predicates

Given an image, the GNN classifier outputs a distribution of binary probabilities over the predicates contained in the image. To recover $\langle \text{subj}, \text{pred}, \text{obj} \rangle$ triplets, we consider the top N predicates and *explain* them one at the time w.r.t. the input image graph. Therefore, the choice of N influences the diversity of predicates contained in the detected relationships, e.g. if we only explained the top scoring predicate we could still recover many triplets but they would all share the same predicate.

For the main results, we set $N = 10$, assuming that in natural images the chance of having more than ten different predicates depicted in the same picture would be rather low. To further prove this point, in figure 11 we plot `recall@50` and `recall@100` for various choices of N on the VRD dataset. Notably, considering very few predicates in the explanation phase, gives poor results on all three relationship detection scenarios. However, increasing N to consider more predicate categories yields diminishing returns after $N = 20$.

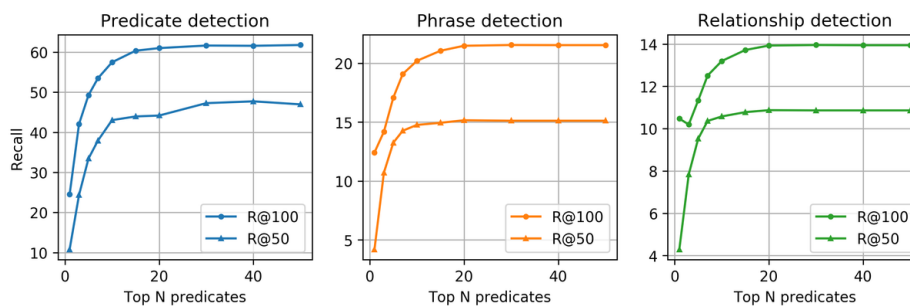


Fig. 11: Recall at 50 (R@50) and at 100 (R@100) on the VRD dataset as the number N of predicates considered for explanation increases from 1 to 50. Diminishing returns are observed, with an elbow at approximately $N = 10$

C.3 Relationship prior

As explained in section 3.4, a weakly-supervised method trained only on predicate labels is not able to learn the directionality of the relations, e.g. it could not distinguish *car on street* from *street on car*. Therefore, we introduced a simple relationship prior based on the frequency of relationships in a small subset of training data. Specifically, we compute:

$$\text{freq}(c_i, c_j | k) = \frac{|\{ \langle c_{\text{subj}}, b_{\text{subj}}, k, c_{\text{obj}}, b_{\text{subj}} \rangle \mid c_{\text{subj}} = c_i, c_{\text{obj}} = c_j, k_{\text{pred}} = k \}|}{|\{ \langle c_{\text{subj}}, b_{\text{subj}}, k, c_{\text{obj}}, b_{\text{subj}} \rangle \mid k_{\text{pred}} = k \}|}$$

In the main experiments, we use a 15% split of the training set to compute this prior, assuming that it would be enough to disambiguate most cases. In figure 12, we show how `recall@50` and `recall@100` on the VRD dataset change

according to the percentage of training triplets used to compute the relationship prior. For each percentage value, we plot the mean recall over 5 random subsets and shade the area corresponding to two standard deviations. We observe that all percentages obtain approximately the same recall, except for 0% that corresponds to a uniform prior. Notably, the randomness introduced when choosing a subset of the given percentage of training data has little effect on the result.

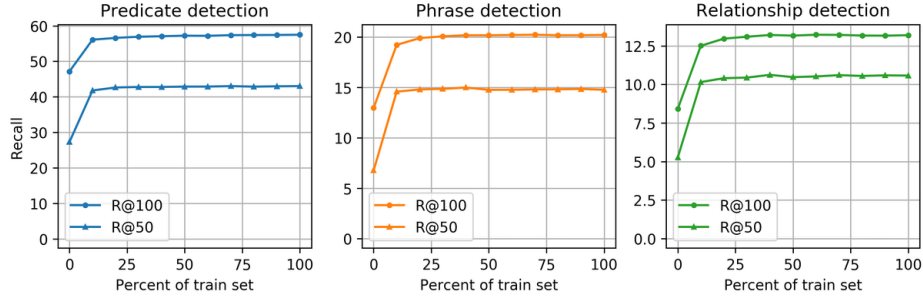


Fig. 12: Recall at 50 (R@50) and at 100 (R@100) on the VRD dataset as the percentage of training data used to compute the relationship prior increases. At each percentage, we run 5 evaluations and plot mean and two standard deviations. Each evaluation uses a different random subset to compute the prior. All percentages obtain approximately the same recall, except for 0% that corresponds to a uniform prior

D Additional results

In this section we report additional qualitative results to evaluate the relationship detection pipeline. We include examples of: correct relationship detections, correct detections missing from the ground truth, incorrect detections due to object misclassification, and incorrect detection due to subject-object inversion, wrong choice of pair, or wrong predicate. All images in figures 13, 14 and 15 are chosen at random from the test sets of each dataset. Then, representative examples are chosen from the top 10 detections of each image (top 25 for UnRel).

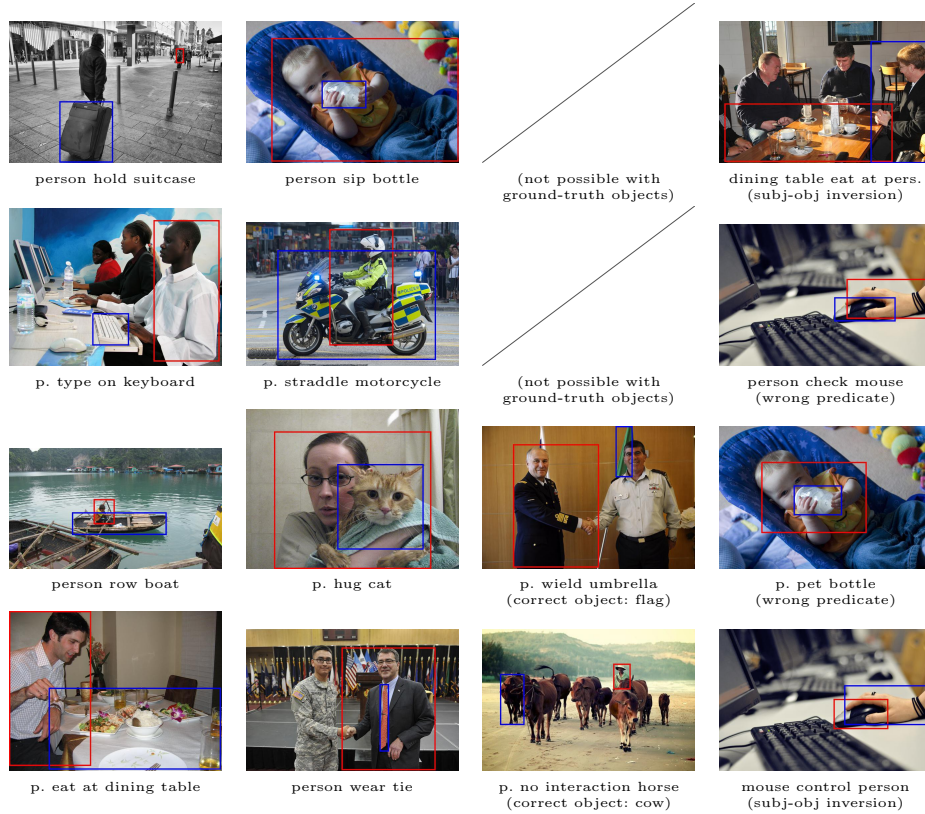


Fig. 13: **Additional detections on HICO-DET.** Top two rows use ground-truth objects, bottom two rows use Faster R-CNN objects. Subjects are framed in red, objects in blue. Left to right: correct relationship detection, correct but missing ground-truth, incorrect due to object misclassification, incorrect detection. Images are chosen at random from the test set, all depicted triplets are selected from the top 10 detections



Fig. 14: **Additional detections on VRD.** Odd rows use ground-truth objects, even rows use Faster R-CNN objects. Subjects are framed in red, objects in blue. Left to right: correct relationship detection, correct but missing ground-truth, incorrect due to object misdetection, incorrect detection. Images are chosen at random from the test set, all depicted triplets are selected from the top 10 detections of an image

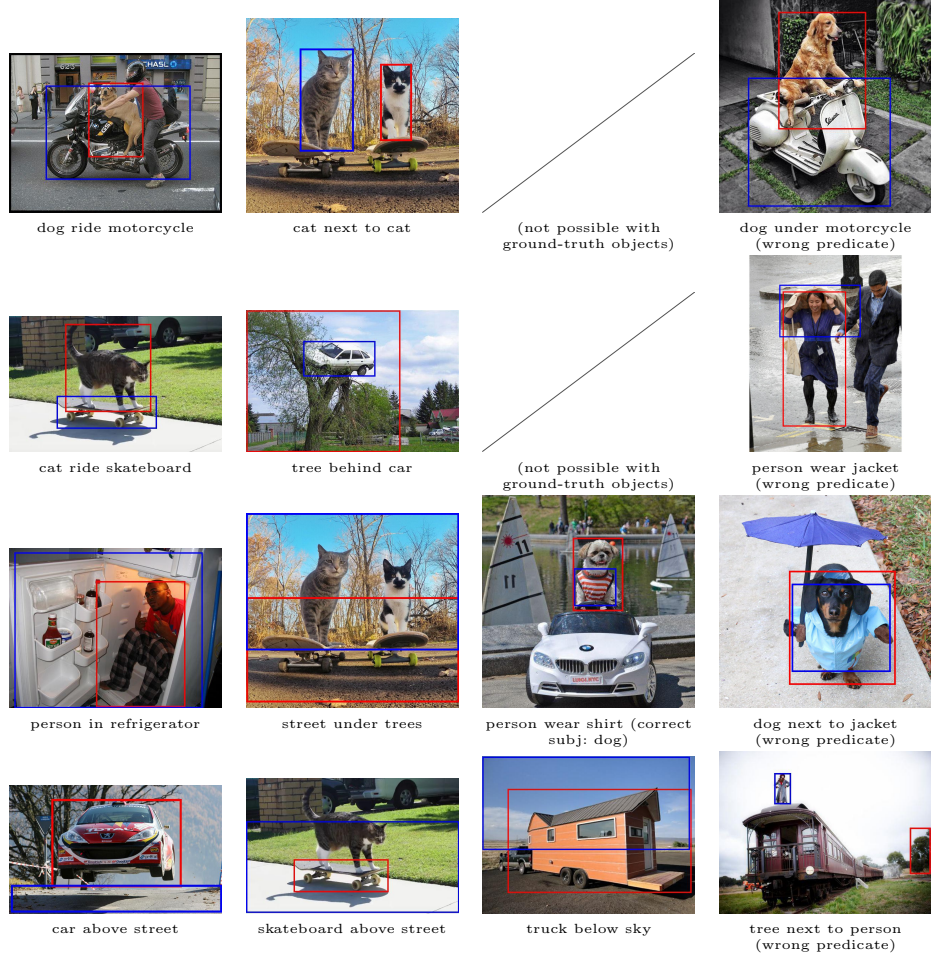


Fig. 15: **Additional detections on UnRel.** Top two rows use ground-truth objects, bottom two rows use Faster R-CNN objects. Subjects are framed in red, objects in blue. Left to right: correct relationship detection, correct but missing ground-truth, incorrect due to object misdetection, incorrect detection. Images are chosen at random from the test set, all depicted triplets are selected from the top 25 detections