# Weakly-supervised learning of visual relations

**Julia Peyre**[1,2]  **Ivan Laptev**[1,2]  **Cordelia Schmid**[2]  **Josef Sivic**[1,2,3]

[1]DI ENS, École normale supérieure, PSL Research University  [2]INRIA  [3]CIIRC, CTU in Prague

## Goal

- Detect visual relations using only **image-level annotations**
- Consider visual relations of the form (**subject**, **predicate**, **object**)

dog taller than person | car under elephant | dog ride bike

## Challenges

- Diversity of visual relations
- Prohibitive cost of exhaustive manual annotation

## Contributions

- Learning visual relations from image-level annotations
- **UnRel** : a new evaluation dataset with clean labels for evaluating the performance and generalization of relation detection

## Overview

*Image-level triplets*

person stand on surfboard
person carry person
person above surfboard
dog on surfboard
person taller than dog

**Training**

**Test**

person on surfboard | dog in front of person | dog stand on surfboard

## Visual representation of a relation

- **Appearance features** for individual objects from fc7 output of Fast-RCNN object detector
- **Quantized spatial configuration** between boxes (with GMM)

$$r(o_s, o_o) = [\frac{x_o - x_s}{\sqrt{w_s h_s}}, \frac{y_o - y_s}{\sqrt{w_s h_s}}, \sqrt{\frac{w_o h_o}{w_s h_s}}, \sqrt{\frac{o_s \cap o_o}{o_s \cup o_o}}, \frac{w_o}{h_o}, \frac{w_s}{h_s}]$$

translation between boxes | ratio of box sizes | overlap | aspect ratio

Visualization of the spatial clusters

## Learning with image-level labels

person ride bike
person hold bike

Discriminative clustering framework [1] :

$$\min_{Z \in \mathcal{Z}} \min_{W \in \mathbb{R}^{d \times R}} \frac{1}{N} \|Z - XW\|^2_F + \lambda \|W\|^2_F$$

matching latent assignments | regularization

$A_1 Z \geq 1$   $A_k Z \geq 1$   **"at least one" constraints**

$Z_{1,,} = 1$   $Z_{2,,} = 1$   $\cdots$   $Z_{N,,} = 1$   **multiclass constraints**

Joint optimization of W and Z with Block-coordinate Frank-Wolfe algorithm [2]
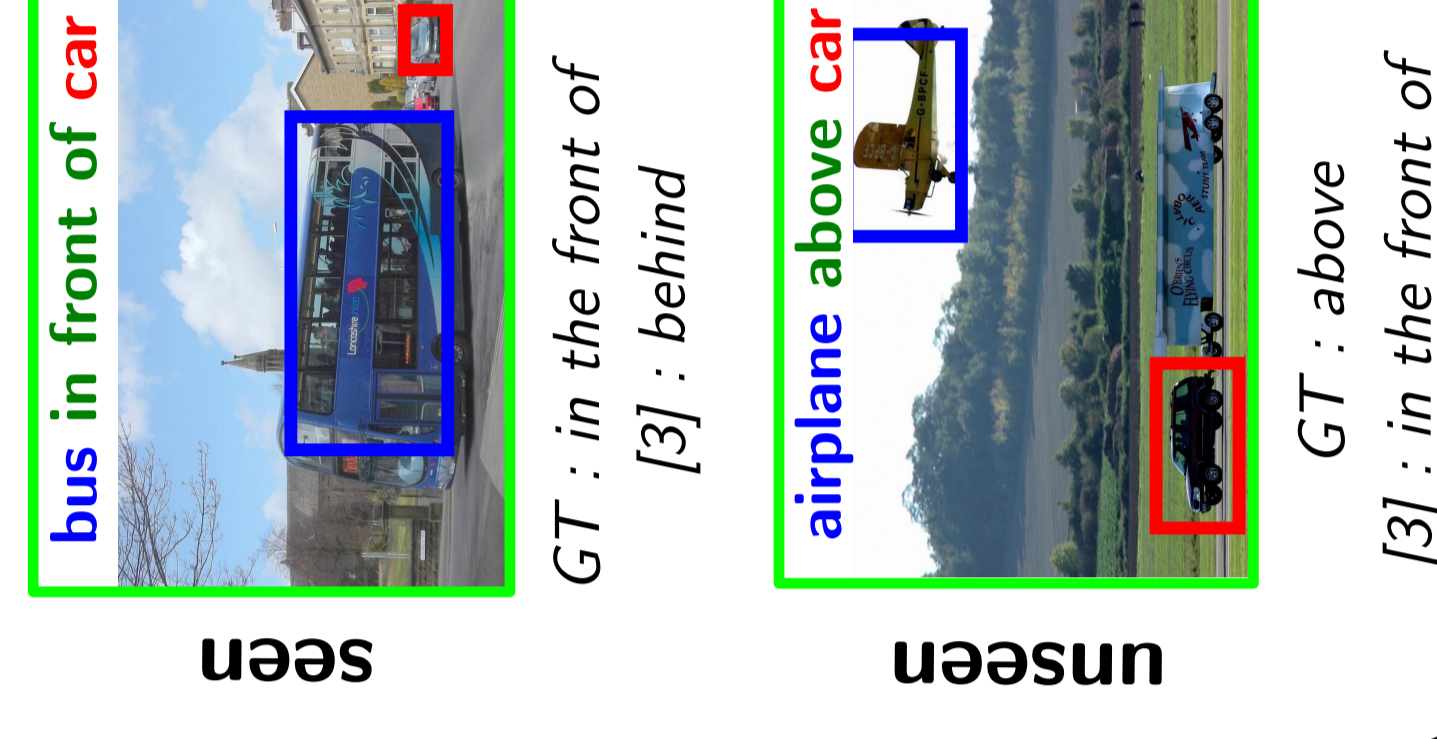
## Recall on Visual Relationship Dataset

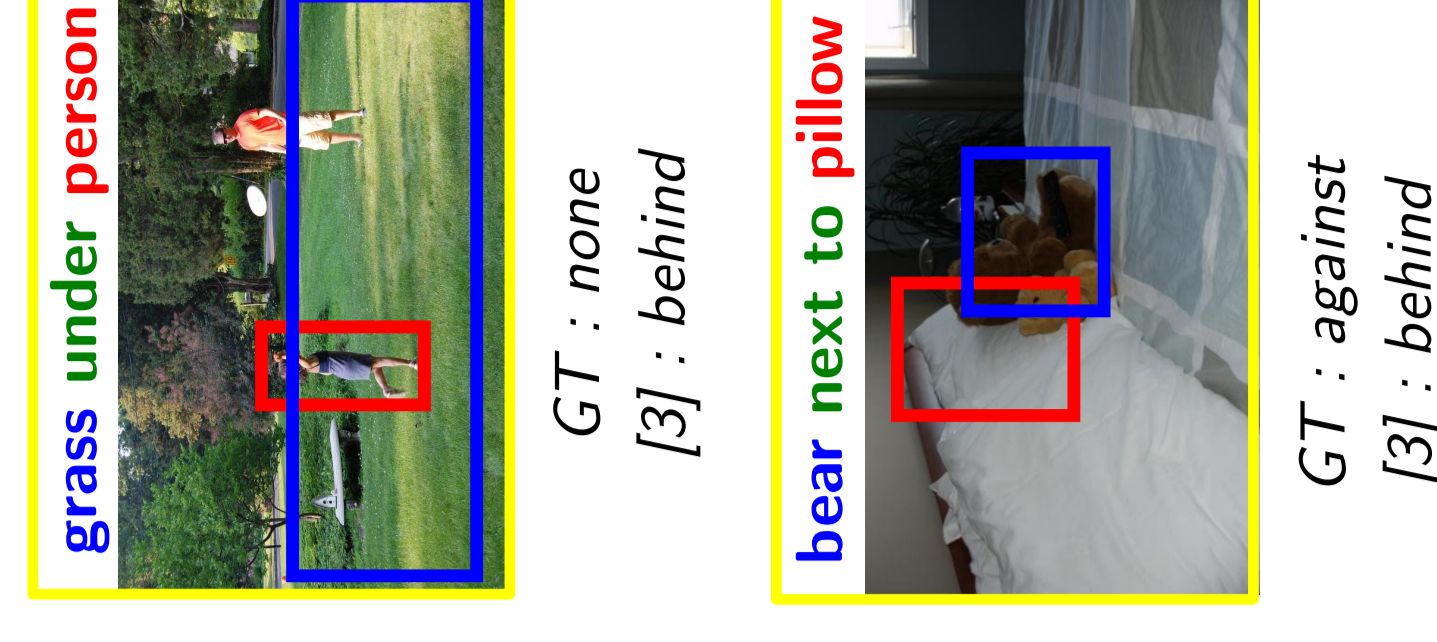**Metric** : recall@k, proportion of ground truth triplets retrieved among the top k detections

| Recall@50 | Predicate Det. | | Phrase Det. | | Relationship Det. | |
|---|---|---|---|---|---|---|
| | All | Unseen | All | Unseen | All | Unseen |
| Visual Phrases [4] | 0.9 | - | 0.04 | - | - | - |
| Language Prior [3] | 47.9 | 8.5 | 16.2 | 3.4 | 13.9 | 3.1 |
| Ours full | **50.4** | **23.6** | **16.7** | **7.4** | **14.9** | **7.1** |
| Ours weak | 46.8 | 19.0 | 16.0 | 6.9 | 14.1 | 6.7 |

*full sup.* | *weak sup.*

**correctly recognized relations**

bus in front of car | person hold umbrella | person wear hat | grass under person

GT : in the front of | GT : hold, under | GT : wear |
[3] : in the front of | [3] : hold | GT : has, wear | GT : none
| | | [3] : wear | [3] : behind

airplane above car | bear on motorcycle | horse have hat | bear next to pillow

GT : above | GT : on | GT : has, wear | GT : against
[3] : in the front of | [3] : drive | [3] : wear | [3] : behind

**seen** / **unseen**

**missing GT** | person on table | **incorrect**

van have car

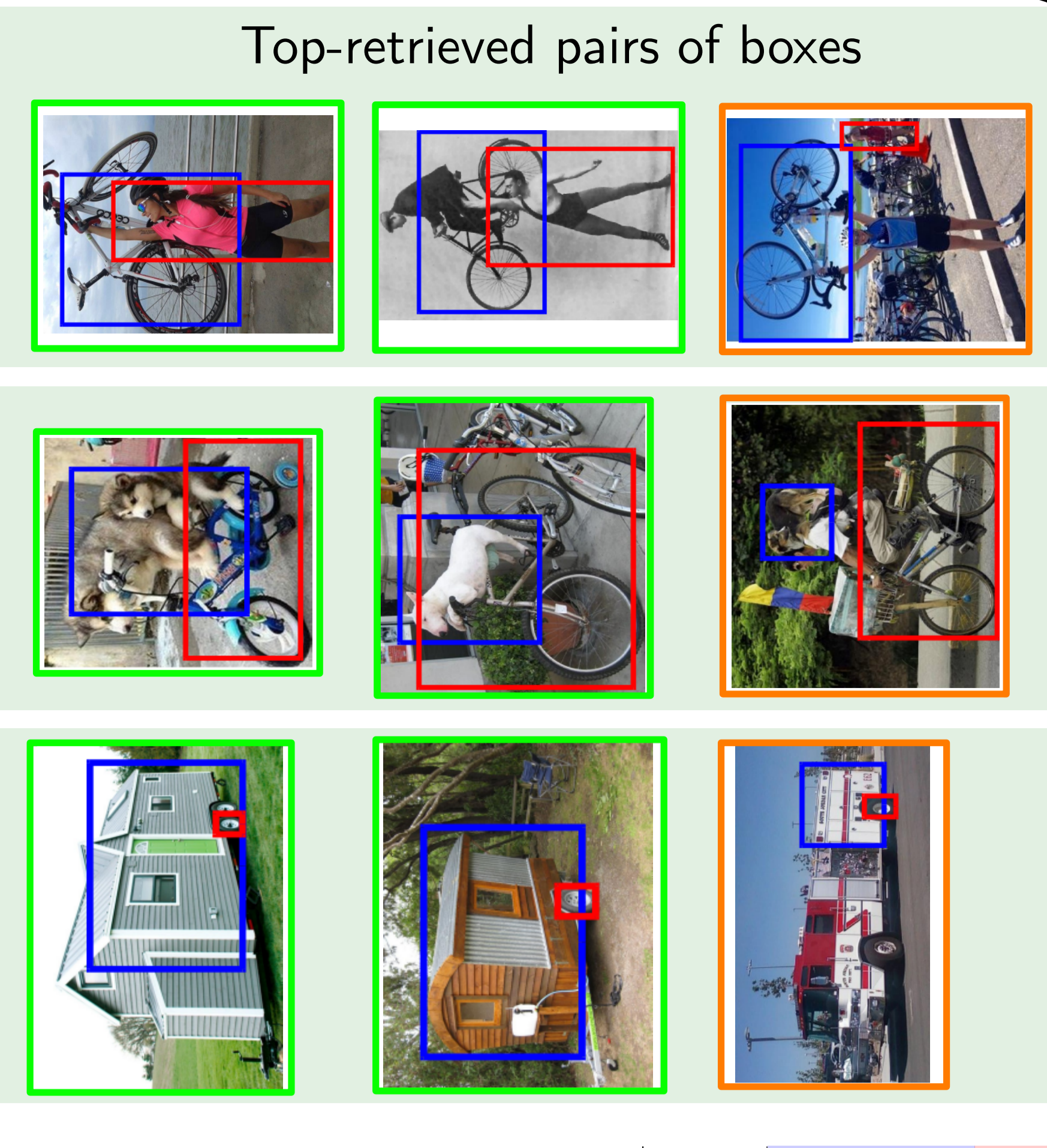GT : next to, behind | GT : behind, left of
[3] : on | [3] : behind

## Retrieval for Unusual Relations (UnRel) Dataset

- in the spirit of Out-of-Context dataset [Choi et al. Context Models and Out-of-Context Objects. In Pattern Recognition Letters, 2012]
- 76 rare queries to test **generalization**
- 1071 images with box-level annotations for relations to evaluate in a **clean setup**
- 1533 relations annotated in total

| mAP | With GT | | With candidates | | |
|---|---|---|---|---|---|
| | union | subj/obj | union | subj | subj/obj |
| DenseCap [5] | - | - | 6.2 | 6.8 | 7.2 |
| Language Prior [3] | 50.6 | | 12.0 | 10.0 | 9.9 |
| Ours full | **62.6** | | **14.1** | **12.1** | |
| Ours weak | 58.5 | | 13.4 | 11.0 | 8.7 |

**Top-retrieved pairs of boxes**

building have wheel | dog ride bike | bike above person

## References

[1] F. R. Bach and Z. Harchaoui. Diffrac: a discriminative and flexible framework for clustering. In NIPS, 2007
[2] A. Osokin, J.-B. Alayrac, I. Lukasewitz, P. K. Dokania, and S. Lacoste-Julien. Minding the gaps for block Frank-Wolfe optimization of structured SVMs.
[3] C. Lu, R. Krishna, M. Bernstein, and L. Fei-Fei. Visual relationship detection with language priors. In ECCV, 2016
[4] M. A. Sadeghi and A. Farhadi. Recognition using visual phrases. In CVPR, 2011
[5] J. Johnson, A. Karpathy, and L. Fei-Fei. Densecap: Fully convolutional localization networks for dense captioning. In CVPR, 2016