# Human-object interaction detection via interactive visual-semantic graph learning

Tongtong WU[1], Fuqing DUAN[1*], Liang CHANG[1] & Ke LU[2]

[1]*College of Artificial Intelligence, Beijing Normal University, Beijing 100875, China;*
[2]*School of Computer and Control Engineering, University of Chinese Academy of Sciences, Beijing 100049, China*

Dear editor,

Human-object interaction (HOI) detection is an important human-centric visual understanding task with several applications in visual monitoring, intelligent robot, etc. It aims at localizing and inferring interaction relationships between humans and objects in images. As a result of their success in object detection, many object detectors can be used to localize human and object instances. Therefore, the key to HOI detection mainly lies in the second part, namely interaction recognition. Many factors such as occlusion or self-occlusion, variation of the human gesture, and camera viewpoint, which bring ambiguity of the interaction recognition, cause it to be challenging.

Recently, owing to the success of deep learning and the availability of large-scale datasets, many deep learning-based HOI models have been proposed [1,2]. These methods usually use context information to improve model performance. However, they mainly adopt visual context information, and semantic context, which can provide effective prior knowledge, is still overlooked.

To address the issue mentioned above, we propose modeling the visual and semantic contexts using a visual graph and a semantic graph, and learning and refining the contextual information using a set of proposed graph update-modules, including graph inner update modules and graph cross update modules. This method reduces the model's burden to extract clues from raw images and provides key HOI prediction evidence.

Given an input image, the task is to detect all humans and objects' bounding boxes and predict their interactions. Like many existing methods, we use an off-the-shelf detector for detection and focus on the stage of interaction classification. We first use a detector to detect all human bounding boxes $\boldsymbol{B}_\mathrm{h}$ and object bounding boxes $\boldsymbol{B}_\mathrm{o}$. Our model's input is the image features from a backbone CNN and these bounding boxes. For each human bounding box $\boldsymbol{b}_\mathrm{h} \in \boldsymbol{B}_\mathrm{h}$ and object bounding box $\boldsymbol{b}_\mathrm{o} \in \boldsymbol{B}_\mathrm{o}$, we predict the HOI score $S_\mathrm{h,o,a}$ where "a" denotes an interaction. The score $S_\mathrm{h,o,a}$ depends on three parts, human score $S_\mathrm{h}$, object score $S_\mathrm{o}$

and interaction score $S_\mathrm{a}$. Specifically, the HOI score $S_\mathrm{h,o,a}$ has the following form:

$$S_\mathrm{h,o,a} = S_\mathrm{h} \times S_\mathrm{o} \times S_\mathrm{a}. \tag{1}$$

In this study, we only focus on the action score $S_\mathrm{a}$.

*Overview.* Figure 1 shows a high-level overview of the proposed network architecture. The network mainly consists of two branches, a visual branch which extracts visual features from the human-object pairs, and a graph branch which models contextual information using graphs and extracts contextual features from graphs. After these branches, the visual and contextual features are fused for classification.

*Visual branch.* This branch extracts visual characteristics from human-object pairs. Given a human/object bounding box, we crop the region from the image feature, and then send it through an RoI pooling layer, a residual block, and a global average pooling layer. After that, each human/object is defined as a visual embedding vector $\boldsymbol{f}$. Letting $\boldsymbol{b}_\mathrm{h} = \{x_\mathrm{h}^1, y_\mathrm{h}^1, x_\mathrm{h}^2, y_\mathrm{h}^2\}$, $\boldsymbol{b}_\mathrm{o} = \{x_\mathrm{o}^1, y_\mathrm{o}^1, x_\mathrm{o}^2, y_\mathrm{o}^2\}$, we embed the spatial relationship of the human-object pair as a vector by

$$
\boldsymbol{f}_\mathrm{sp} = \left\{ \frac{x_\mathrm{h}^1}{W_\mathrm{u}}, \frac{y_\mathrm{h}^1}{H_\mathrm{u}}, \frac{x_\mathrm{h}^2}{W_\mathrm{u}}, \frac{y_\mathrm{h}^2}{H_\mathrm{u}}, \frac{x_\mathrm{o}^1}{W_\mathrm{u}}, \frac{y_\mathrm{o}^1}{H_\mathrm{u}}, \frac{x_\mathrm{o}^2}{W_\mathrm{u}}, \right.
$$
$$
\left. \frac{y_\mathrm{o}^2}{H_\mathrm{u}}, \frac{A_\mathrm{h}}{A_\mathrm{u}}, \frac{A_\mathrm{o}}{A_\mathrm{u}}, \frac{A_\mathrm{h}}{A_\mathrm{i}}, \frac{A_\mathrm{o}}{A_\mathrm{i}}, \frac{x_\mathrm{h}^1 - x_\mathrm{o}^1}{x_\mathrm{o}^2 - x_\mathrm{o}^1}, \frac{y_\mathrm{h}^1 - y_\mathrm{o}^1}{y_\mathrm{o}^2 - y_\mathrm{o}^1} \right\}, \tag{2}
$$

where $W_\mathrm{u}$ and $H_\mathrm{u}$ represent the width and height of the union bounding box, $A_\mathrm{h}$ and $A_\mathrm{o}$ represent the area of the two bounding boxes and $A_\mathrm{i}$ and $A_\mathrm{u}$ are the area of the intersection and union bounding boxes, respectively. The spatial relationship vector is applied to a fully connected layer, and outputs the spatial embedding vector $\boldsymbol{f}_\mathrm{sp}^{'}$. Finally, human embedding vector $\boldsymbol{f}_\mathrm{h}$, object embedding vector $\boldsymbol{f}_\mathrm{o}$, and spatial embedding vector $\boldsymbol{f}_\mathrm{sp}^{'}$ are concatenated and the concatenated vector is sent into a fully connected layer.

$$\boldsymbol{f}_\mathrm{vis} = \mathrm{relu}(\boldsymbol{W}_1 \times (\boldsymbol{f}_\mathrm{h} \oplus \boldsymbol{f}_\mathrm{sp}^{'} \oplus \boldsymbol{f}_\mathrm{o}) + \boldsymbol{b}_1), \tag{3}$$

where $\oplus$ denotes the concatenate operation.

* Corresponding author (email: fqduan@bnu.edu.cn)

**Figure 1** (Color online) Model architecture.

*Graph branch.* This branch uses graph convolution networks to capture useful contextual information. Graph convolution networks model the inputs through a graph composed of nodes and edges and extract features by traversing and updating the nodes. We propose two kinds of graphs, visual graph and semantic graph. In visual graph, human embedding vectors $\boldsymbol{f}_\text{h}$ and object embedding vectors $\boldsymbol{f}_\text{o}$ are defined as nodes. In semantic graph, we use Word2Vec [3] to generate word embedding vectors based on human and object categories. We add a virtual object node whose bounding box is the same size as the input image to retain information integrity. Thus, all contextual information, including background, is preserved. The virtual node is constructed based on the entire image feature map in visual graph and is a unit vector in semantic graph. The spatial embedding vectors $\boldsymbol{f}_\text{sp}'$ of node pairs are defined as edges. Finally, these visual nodes, semantic nodes, and edges are used to construct a fully connected visual graph and a semantic graph, respectively.

Graph update consists of two modules, a graph inner update module and a graph cross update module. Messages are transferred within each graph in graph inner update module and between the visual graph and semantic graph in graph cross update module. In graph inner update module, we use graph attention [4] method to update graph, while in graph cross update module, we use corresponding pairs of visual and semantic nodes to update each other's features through a basic attention method. The details of graph update modules are shown in Appendix A.

*Classification.* Given a human-object pair, we extract a visual feature vector from the output of the visual branch, and from the final updated graphs, we obtain corresponding human node feature vector $\boldsymbol{h}_\text{h}$, object node feature vector $\boldsymbol{h}_\text{o}$, and human-object edge feature vector $\boldsymbol{h}_\text{e}$. The feature vectors $\boldsymbol{h}_\text{h}$, $\boldsymbol{h}_\text{o}$, and $\boldsymbol{h}_\text{e}$ are concatenated, and the concatenated feature vector is sent into a fully connected layer to generate contextual features $\boldsymbol{h}_\text{g}$. For visual graph and semantic graph, we calculate the contextual features $\boldsymbol{h}_\text{g,vis}$ and $\boldsymbol{h}_\text{g,sem}$, respectively, and concatenate $\boldsymbol{h}_\text{g,vis}$ and $\boldsymbol{h}_\text{g,sem}$ as $\boldsymbol{f}_\text{g}$.

$$\boldsymbol{h}_\text{g} = \text{relu}(\boldsymbol{W}_2 \times (\boldsymbol{h}_\text{h} \oplus \boldsymbol{h}_\text{e} \oplus \boldsymbol{h}_\text{o}) + \boldsymbol{b}_2), \tag{4}$$

$$\boldsymbol{f}_\text{g} = \boldsymbol{h}_\text{g,vis} \oplus \boldsymbol{h}_\text{g,sem}. \tag{5}$$

The concatenated vector $\boldsymbol{f}_\text{g}$ represents the final contextual features, and is fused with $\boldsymbol{f}_\text{vis}$ by a basic attention method. The attention output $\boldsymbol{f}_\text{out}$ is used directly through a fully connected layer for interaction classification by a multi-hot sigmoid cross entropy function.

*Experiments.* Experiments on prominent datasets HICO-DET and V-COCO show that our model is effective. The details of the experiment results and an ablation study are shown in Appendixes B and C, respectively.

*Conclusion.* In this study, we propose modeling the context visually and semantically by combining a visual graph and a semantic graph and learning a vital context in the HOI problem using a group of graph update-modules, including graph inner update modules and graph cross update modules. We fuse the contextual features from the visual graph and semantic graph with the visual characteristics of the human-object pairs in a network to detect HOIs. We evaluate our proposed model on two challenging datasets, HICO-DET and V-COCO, and demonstrate excellent performance. Our work can provide a reference for modeling contextual information in the HOI problem.

**References**

1 Gao C, Zou Y L, Huang J-B. iCAN: instance-centric attention network for human-object interaction detection. 2018. ArXiv:1808.10437

2 Wang T, Anwer R M, Khan M H, et al. Deep contextual attention for human-object interaction detection. In: Proceedings of IEEE International Conference on Computer Vision, 2019

3 Mikolov T, Sutskever I, Chen K, et al. Distributed representations of words and phrases and their compositionality. In: Proceedings of Advances in Neural Information Processing Systems, 2013. 3111–3119

4 Veličković P, Cucurull G, Casanova A, et al. Graph attention networks. 2017. ArXiv:1710.10903