# Zero-Shot Human-Object Interaction Recognition via Affordance Graphs

Alessio Sarullo
*Department of Computer Science*
*University of Manchester*
Manchester, UK
alessio.sarullo@manchester.ac.uk

Tingting Mu
*Department of Computer Science*
*University of Manchester*
Manchester, UK
tingting.mu@manchester.ac.uk

## Abstract

*We propose a new approach for Zero-Shot Human-Object Interaction Recognition in the challenging setting that involves interactions with unseen actions (as opposed to just unseen combinations of seen actions and objects). Our approach makes use of knowledge external to the image content in the form of a graph that models affordance relations between actions and objects, i.e., whether an action can be performed on the given object or not. We propose a loss function with the aim of distilling the knowledge contained in the graph into the model, while also using the graph to regularise learnt representations by imposing a local structure on the latent space. We evaluate our approach on several datasets (including the popular HICO and HICO-DET) and show that it outperforms the current state of the art.*

## 1. Introduction

Human-Object Interaction (HOI) Recognition is the task of identifying how people interact with the surrounding objects from the visual appearance of the scene and it is of paramount importance to understand the content of an image. It consists of producing a set of $\langle human, action, object \rangle$ triplets for the input image, providing a concise representation of the image semantics that can be used in higher-level tasks like Image Captioning [1] or Human-Robot Interaction [2].

One of the greatest difficulties when dealing with visual relations is that the number of possible triplets increases multiplicatively in the cardinality of the human, action and object spaces. Even if we do not distinguish between various "person" categories such as "man", "child" etc., the number of possible interactions – that is, $\langle action, object \rangle$ pairs – grows quadratically. Due to the practical challenges of building a dataset, it is common for only a subset of all possible interactions to be annotated, while a large number remains unlabelled; for instance, HICO [3] contains



Figure 1: Left to right: $\langle eating, sandwich \rangle$, $\langle eating, pizza \rangle$, $\langle cooking, pizza \rangle$. Both pairs of objects (*pizza* and *sandwich*) and actions (*eating* and *cooking*) are semantically similar, yet images that share an action look more similar than images that share an object.

only 600 interactions out of the 9360 possible pairs (among the 9360-600=8760 unlabelled interactions, some are invalid like $\langle eating, bottle \rangle$, while some are valid but missing like $\langle carrying, knife \rangle$). This is why more and more approaches are focusing on Zero-Shot Learning (ZSL) for HOI Recognition [4–7]. ZSL aims to alleviate the problems caused by the combinatorial growth of the number of possible interactions by allowing models to make predictions about previously unseen interactions.

We focus on actions, as they play a more significant role than objects in defining an interaction: several studies in Psychology [8], Neurobiology [9] and Computer Vision [7, 10] show that objects can be categorised and recognised based on their affordances, making the semantics of an object defined in term of actions, and we empirically verify this intuition via some visual examples provided in Figure 1. For this reason, we follow a challenging zero-shot setting [5] that consists of predicting interactions containing unseen action and object classes, instead of only new combinations of seen classes. We adopt a compositional strategy, as in [4, 6]: we detect objects and actions first and then combine the results to detect interactions. This is effective in the considered zero-shot setting, where many of the unseen interactions are combinations of a seen object/action with a new action/object, as the model will find it easier to predict the component containing the seen class.

Our model uses a Graph Convolutional Network (GCN)

[11] to learn unseen classes in a semi-supervised manner [12, 13]. The graph's connectivity determines how nodes are linked to each other and thus how information is aggregated in the learnt representations. We make use of a particular type of graph called an *affordance graph*, that is, a graph whose edges model *affordances* [5, 8]: action-object pairs $\langle a, o \rangle$ where $a$ can be performed on $o$ (e.g., $\langle hold, apple \rangle$, because apples can be held). Such a graph enables the model to learn what interactions are affordable regardless of whether they appear in the training set, allowing it to perform zero-shot predictions.

The focus of this paper is to propose a new training objective function that aims to improve the representations learnt by the model. More specifically, the proposed objective function enhances the loss used by state-of-the-art approaches in two ways. First, it effectively distils action affordance in the unseen class representations by making use of relations from the affordance graph to train unseen actions in a weakly-supervised way. As a result, the model learns to distinguish which unseen actions can be performed on a given object and which ones cannot. Second, it imposes a local structure on the latent space through a regulariser that clusters unseen class representations together with similar classes according to the affordance graph. Additionally, we attempt to tackle a shortcoming that affects current approaches: GCN's seen action representations are affected by unseen ones, which are not trained in a fully supervised way and thus add noise. Therefore, we learn an alternative set of representations for seen classes unaffected by unseen ones. Qualitative results demonstrate that our model (shown in Figure 2) learns representations that are effective at differentiating actions based on affordances, and our experiments show that our model outperforms the current state of the art on HICO [3], VG-HOI [5] and COCO-a [14].

## 2. Related Work

### 2.1. Knowledge Usage in HOI Recognition

Many works have been proposed to perform HOI Detection in recent years, the most similar to ours being the ones that make use of pre-existing knowledge [5–7, 15]. In [7] a language component is used to identify functionally similar objects, effectively augmenting the training data with new interaction instances. In the other works, the pre-existing knowledge is used to obtain class representations, which are used for prediction. These representations come from word embeddings that are mapped through functions implemented as a Multi-Layer Perceptron (MLP) [6] or a GCN [5, 15]. An important difference between these models lies in what representations are computed: while in [6] action, object and interactions classes are all considered and the respective scores combined in a compositional way, in the

other methods only representations for actions [15] or interactions [5] are used for prediction. Our approach is similar to [6] regarding the compositional model and to [5, 15] in the utilisation of external knowledge to build the graph used by the GCN, but differs from all of the above mainly in the way we use the graph at training time to regularise action representation and to distil affordance information into the model.

### 2.2. Zero-Shot Learning

The growing field of ZSL primarily aims to overcome the difficulties of dealing with a non-exhaustively annotated dataset. A common framework to perform ZSL [12,13,16,17] is to exploit some kind of pre-existing knowledge to transfer to unseen classes what has been learnt about seen ones in a semi-supervised way. Representations are learnt for both classes and instances and compared through a similarity function to predict output probabilities. The model is trained by feeding the output scores for seen classes into a loss function such as a ranking loss [16], least squares [17] or cross entropy [12, 13].

An interesting method to learn better representations is to add a regularisation loss [18, 19]. In particular, [18] maps label embeddings into the visual space, adding a reconstruction loss to make sure that the inverse transformation is also possible and thus the visual projection preserves semantics. A different technique is used in [19], where a cross-reconstruction loss between images and labels is added in order to "pull together" representations of the same class from the two different sources (image and labels). Inspired by these works, we formulate a different regularisation loss that uses the affordance graph and is thus better suited to our goal of modelling action affordance.

A few recent approaches tackle ZSL in HOI Recognition/Detection [4–7]. We compare our results to the works that considers unseen actions [5, 6], as they are the most closely related to ours.

## 3. Notation and Problem Statement

Let us denote by $\mathcal{O}$ and $\mathcal{A}$ the ordered set of objects and actions, respectively. For instance, we might have *apple* $\in \mathcal{O}$ and *eat* $\in \mathcal{A}$. We will denote the elements of these sets by the corresponding lowercase letter (for example $o_j$ is the $j$-th element in set $\mathcal{O}$), or sometimes by the index only (for example we will write $k \in \mathcal{A}$ instead of $a_k \in \mathcal{A}$).

Our dataset is denoted by $\mathcal{D} = \{(I_i, \mathbf{T}_i)\}_{i=1}^{M}$. Here, $I_i$ is the $i$-th image and $\mathbf{T}_i \in \{0, 1\}^{|\mathcal{O}| \times |\mathcal{A}|}$ is its label matrix, with its $jk$-th element $t_{ijk}$ being 1 if and only if example $i$ is annotated with interaction $\langle a_k, o_j \rangle$ (note that an image can have multiple labels). Under the considered Zero-Shot Learning setting, we assume that there are no available visual examples for some objects and actions. This
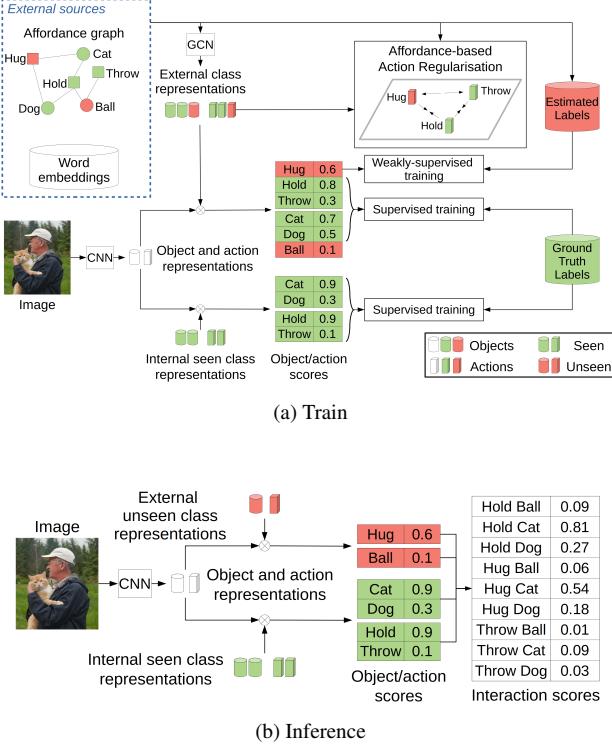
(a) Train



(b) Inference

Figure 2: Overview of the proposed model for HOI Recognition during training (a) and inference (b). $\otimes$ indicates dot product. Best viewed in colour.

is equivalent to omit the corresponding labels from all images during training, although the affected images might still be annotated with other labels that have not been omitted. The omitted class set will be denoted with $\mathcal{U}$ (they are *unseen*), while $\mathcal{S}$ is the set of *seen* (i.e., trained-on) classes. Therefore, we have $\mathcal{O} = \mathcal{S}^O \cup \mathcal{U}^O$ and $\mathcal{A} = \mathcal{S}^A \cup \mathcal{U}^A$. Note that seen and unseen classes do not intersect, i.e., $\mathcal{S}^q \cap \mathcal{U}^q = \emptyset \; \forall q \in \{O, A\}$. The task is to learn a model that is able to predict any interaction $\langle a_k, o_j \rangle$, even when $o_j \in \mathcal{U}^O$ or $a_k \in \mathcal{U}^A$ (that is, when either or both object and action are unseen).

## 4. Proposed Method

### 4.1. Affordance Graph

The main motivation of this work is to improve zero-shot interaction recognition by using structured external knowledge, which is expressed in the form of an *affordance graph.* We define it as a graph $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$ whose nodes $\mathcal{V}$ are objects and actions and edges $\mathcal{E}$ represent *affordances* [8]: object node $o_j$ is connected to action node $a_k$ only if $a_k$ can be performed on $o_j$, i.e., $\langle a_k, o_j \rangle$ constitutes a valid interaction. For example, *eat* and *apple* will be connected, but *eat* and *fork* will not because people cannot eat forks. This graph is

undirected and bipartite: all links are symmetric and there are no connections between object nodes, nor between action nodes. We construct the affordance graph by mining interactions from external sources, to simulate a real-world scenario where no interaction information regarding unseen classes is available. Details about the construction process will be provided in Section 5.2.

### 4.2. Model Architecture

#### 4.2.1 Preliminary: Graph Convolutional Networks

Let us consider a graph with $N$ nodes, adjacency matrix $\mathbf{A}$ and initial node representations $\mathbf{Z}_0 \in \mathbb{R}^{N \times d_0}$ for some dimension $d_0$. A single layer of a Graph Convolutional Network (GCN) [11] computes a new representation for each node by aggregating the ones of its neighbours according to $\mathbf{Z}_1 = \phi_1(\tilde{\mathbf{A}}\mathbf{Z}_0\boldsymbol{\Theta}_1)$, where $\phi$ is an activation function such as ReLU [20], $\boldsymbol{\Theta}_1 \in \mathbb{R}^{d_0 \times d_1}$ are the layer parameters and $\tilde{\mathbf{A}} = \mathbf{D}^{-\frac{1}{2}}(\mathbf{A} + \mathbf{I})\mathbf{D}^{-\frac{1}{2}}$ is the normalised adjacency matrix, where $\mathbf{D}$ is a diagonal matrix with $d_{ii} = 1 + \sum_{j=1}^{N} a_{ij}$ and $\mathbf{I}$ is the identity matrix. Deep GCNs can be composed by stacking $L$ of such layers, producing final representations $\mathbf{Z} = f_{GCN}(\mathbf{Z}_0) = \phi_L(\tilde{\mathbf{A}}\phi_{L-1}(\dots \phi_1(\tilde{\mathbf{A}}\mathbf{Z}_0\boldsymbol{\Theta}_1)\dots)\boldsymbol{\Theta}_L)$. We refer the reader to [11] for more details.

#### 4.2.2 Overview

Our model takes as input an image $I$, which is fed into a Convolutional Neural Network (CNN) such as ResNet [21], producing image-level visual features $\boldsymbol{v} = f_{CNN}(I)$. These features are fed into two identically structured modules indexed by variable $q$, one for objects ($q = O$) and one for actions ($q = A$). Specifically, for each module we compute a $d$-dimensional representation $\boldsymbol{x}^q = f_1^q(\boldsymbol{v})$ through a non-linear mapping $f_1^q$ (e.g., an MLP). Vector $\boldsymbol{x}^q$ is compared to a set of $d$-dimensional class representations $\mathbf{Z}^q = [\boldsymbol{z}_1^q \mid \dots \mid \boldsymbol{z}_{|\mathcal{S}^q \cup \mathcal{U}^q|}^q]$ through a similarity function $g(\boldsymbol{x}^q, \boldsymbol{z}_i^q)$, that we implement as inner product following [22]: $g(\boldsymbol{x}, \boldsymbol{z}) = \boldsymbol{x}^T\boldsymbol{z}$. Similarity scores are fed into the sigmoid function $\sigma(x) = 1/(1 + e^{-x})$ to output probabilities $\boldsymbol{y}^q = \sigma(\mathbf{Z}^q\boldsymbol{x})$. We will now describe how to compute class representations $\mathbf{Z}^q$.

#### 4.2.3 Class Representations

In our model, akin to [5, 12, 13], we use a GCN to train unseen class representations in a semi-supervised way, effectively embedding the affordance relations contained in the graph into the learnt representations. We also incorporate additional semantic information computed from word embeddings, but, differently from previous work, we do not use them to initialise GCN's input embeddings. The reason is that the affordance graph and word embeddings provide

different types of semantics: the former aims to capture affordances, while the latter co-occurrence statistics. As a result, for instance, "eat" and "drink" are distant according to affordances while close according to word embeddings, which can result in a mismatch in action similarity that brings down the performance (see supplementary). However, co-occurrence semantics carried by word embeddings can be useful for objects (e.g., "pizza" and "sandwich" have high similarity according to word embeddings, and indeed are both objects that can be eaten), so we use word embeddings to enrich the objects representations produced by the GCN. The final class representations $\mathbf{Z}_{EXT}^O \in \mathbb{R}^{|\mathcal{O}| \times d}$ and $\mathbf{Z}_{EXT}^A \in \mathbb{R}^{|\mathcal{A}| \times d}$, that we call *external representations*, are

$$\mathbf{Z}_{EXT}^O = (\mathbf{Z}_{GCN})_{\mathcal{O},:} + f_2(\mathbf{W}^O) \tag{1}$$

$$\mathbf{Z}_{EXT}^A = (\mathbf{Z}_{GCN})_{\mathcal{A},:} \tag{2}$$

$$\mathbf{Z}_{GCN} = f_3\left(f_{GCN}\left(\mathbf{Z}_0\right)\right) , \tag{3}$$

where $(\mathbf{Z}_{GCN})_{\mathcal{O},:}$ and $(\mathbf{Z}_{GCN})_{\mathcal{A},:}$ denote the rows of $\mathbf{Z}_{GCN}$ corresponding to object and action classes (respectively), $f_2$ and $f_3$ are non-linear functions (e.g., MLPs), $\mathbf{W}^O \in \mathbb{R}^{|\mathcal{O}| \times d'}$ are $d'$-dimensional word embeddings and GCN's input embeddings $\mathbf{Z}_0 \in \mathbb{R}^{(|\mathcal{O}|+|\mathcal{A}|) \times d_0}$ are randomly initialised. We use $\mathbf{Z}_{EXT}^O$ and $\mathbf{Z}_{EXT}^A$ to predict class probabilities $\boldsymbol{y}_{EXT}^q = \sigma(\mathbf{Z}_{EXT}^q \boldsymbol{x})$ for $q \in \{O, A\}$. Note that these representations (and the corresponding probabilities) are computed for both seen and unseen classes. We use $\mathbf{Z}_{EXT-\mathcal{S}}^q = (\mathbf{Z}_{EXT}^q)_{\mathcal{S}^q,:}$ and $\mathbf{Z}_{EXT-\mathcal{U}}^q = (\mathbf{Z}_{EXT}^q)_{\mathcal{U}^q,:}$ to denote the sub-matrices of $\mathbf{Z}_{EXT}^q$ that only contain rows for seen or unseen classes, respectively.

The representations computed as above contain informations aggregated from neighbours in the affordance graph. This is why they are well-suited for ZSL, but the downside is that seen class representations are affected by unseen ones. This introduces noise in the representation of seen classes, and in fact we empirically verify that it lowers performance (see supplementary). To overcome this issue, we train an alternative set of representations for seen classes $\mathbf{Z}_{INT}^q \in \mathbb{R}^{|\mathcal{S}^q| \times d}$, called *internal representation*, in the standard supervised way. This results in separate probability vectors $\boldsymbol{y}_{INT}^q = \sigma(\mathbf{Z}_{INT}^q \boldsymbol{x})$ for seen object and action classes.

#### 4.2.4 Inference

At inference time a score has to be assigned to every interaction, producing a matrix $\mathbf{Y} \in [0,1]^{|\mathcal{O}| \times |\mathcal{A}|}$ whose element $y_{jk}$ constitutes the probability for interaction $\langle a_k, o_j \rangle$. We

do so by multiplying object and action scores together:

$$\mathbf{Y} = \begin{bmatrix} \boldsymbol{y}_{INT}^O \\ \boldsymbol{y}_{EXT-\mathcal{U}}^O \end{bmatrix} \begin{bmatrix} \boldsymbol{y}_{INT}^A \\ \boldsymbol{y}_{EXT-\mathcal{U}}^A \end{bmatrix}^T$$
$$= \begin{bmatrix} \sigma\left(\mathbf{Z}_{INT}^O \boldsymbol{x}\right) \\ \sigma\left(\mathbf{Z}_{EXT-\mathcal{U}}^O \boldsymbol{x}\right) \end{bmatrix} \begin{bmatrix} \sigma\left(\mathbf{Z}_{INT}^A \boldsymbol{x}\right) \\ \sigma\left(\mathbf{Z}_{EXT-\mathcal{U}}^A \boldsymbol{x}\right) \end{bmatrix}^T . \tag{4}$$

In our model, external representations of seen classes are used during training to allow for unseen ones to be learnt in a semi-supervised fashion through the GCN, but they are not used for inference (see Figure 2b).

### 4.3. Training

Our model is trained by minimising the following composite loss function, which is designed to optimise both internal class representations $\mathbf{Z}_{INT}^q$ and all the remaining parameters $\boldsymbol{\Theta}$ (which include weights for MLPs and GCN, as well as GCN's initial representations $\mathbf{Z}_0$) through variables $\boldsymbol{y}_{INT}^q$, $\boldsymbol{y}_{EXT}^q$ and $\mathbf{Z}_{EXT}^q$:

$$\min_{\substack{\boldsymbol{\Theta} \\ \mathbf{Z}_{INT}^O \\ \mathbf{Z}_{INT}^A}} \sum_{i=1}^M \sum_{q \in \{O,A\}} \ell\left((\boldsymbol{y}_{INT}^q)_i, \boldsymbol{t}_i^q, \mathcal{S}^q\right) + \tag{5}$$
$$\sum_{i=1}^M \sum_{q \in \{O,A\}} \ell\left((\boldsymbol{y}_{EXT-\mathcal{S}}^q)_i, \boldsymbol{t}_i^q, \mathcal{S}^q\right) +$$
$$\sum_{i=1}^M \lambda\, \ell\left((\boldsymbol{y}_{EXT-\mathcal{U}}^A)_i, \hat{\boldsymbol{t}}_i^A, \mathcal{U}^A\right) +$$
$$\rho \mathcal{L}_{REG}\left(\mathbf{Z}_{EXT}^A\right) ,$$

where $\lambda$ and $\rho$ are hyperparameters that regulate the contribution of their respective terms, label vectors $\boldsymbol{t}_i^O \in \{0,1\}^{|\mathcal{O}|}$ and $\boldsymbol{t}_i^A \in \{0,1\}^{|\mathcal{A}|}$ are obtained from matrix $\mathbf{T}_i$ according to $t_{ij}^O = \max_{k \in \mathcal{A}} t_{ijk}$ and $t_{ik}^A = \max_{j \in \mathcal{O}} t_{ijk}$, and $\ell$ is the standard binary cross entropy loss:

$$\ell\left(\boldsymbol{y}, \boldsymbol{t}, \mathcal{J}\right) = \sum_{j \in \mathcal{J}} \left[t_j \log y_j + (1 - t_j) \log(1 - y_j)\right] , \tag{6}$$

where $\boldsymbol{y}$ are outputs, $\boldsymbol{t}$ target labels and $\mathcal{J}$ a set of indices. We also add $L_2$-regularisation to $\boldsymbol{\Theta}$ to prevent overfitting (not shown in Equation (5)).

The first two terms of Equation (5) implement a standard training loss, which uses ground truth labels to reward pairing instances with the corresponding seen classes and to penalise assigning the wrong class. The third term aims to train unseen actions in the same way. However, since ground truth labels are not available for unseen actions, we adopt a weakly-supervised approach and estimate labels $\hat{\boldsymbol{t}}^A$

as:

$$\hat{t}_k^A = \max_{j \in \mathcal{S}^O} m_{jk} s_{jk} \qquad \forall k \in \mathcal{U}^A \tag{7}$$

$$s_{jk} = \frac{1}{\sum_{h \in \mathcal{S}^A} t_{jh}} \sum_{h \in \mathcal{S}^A} t_{jh} \left[ \boldsymbol{w}_h^T \boldsymbol{w}_k \right]_+ , \tag{8}$$

where $[x]_+ = \max(x, 0)$, $\mathbf{M} \in \{0, 1\}^{|\mathcal{O}| \times |\mathcal{A}|}$ is the graph adjacency matrix[1] and $\boldsymbol{w}_k$ is the word embedding for the $k$-th action. Equation (8) computes a score that determines how likely unseen action $k$ describes an image containing object $j$. This score is *not* binary, but rather a real value in $[0, 1]$. This is needed because binary estimated labels would incur the risk of introducing noise, since we cannot know which of the affordable unseen actions are actually depicted in a particular image. Word embeddings are used to assign a score based on the similarity with labelled seen actions (which are compatible with object $j$, since they come from the ground truth) through the positive inner product $\left[ \boldsymbol{w}_h^T \boldsymbol{w}_k \right]_+$, so that unseen actions similar to shown seen ones will be assigned a higher score: if $o_j = $ *person* and *hug* is a labelled seen action, *kiss* and *greet* are better unseen candidates than *teach*. Action affordance is distilled into the model according to Equation (7): score $s_{jk}$ contributes to $\hat{t}_k^A$ only if $m_{jk} = 1$, that is, only if $\langle a_k, o_j \rangle$ is an affordable action. Since an image may contain multiple objects, the maximum score over objects is taken according to the Multiple Instance Learning framework [23].

Additionally, we use the affordance graph as a regulariser for action classes, with the goal of learning better representations by inducing a structure onto the latent space based on affordances. Specifically, we want to group *functionally* similar actions, that is, actions that can be performed on the same objects. To this end, we use the following ranking margin loss:

$$\mathcal{L}_{REG}\left(\mathbf{Z}_{EXT}^A\right) = \sum_{i \in \mathcal{U}^A} \sum_{j \in \mathcal{N}(i)} \sum_{k \notin \mathcal{N}(i)} \left[\gamma - c_{ij} + c_{ik}\right]_+$$

$$c_{ij} = \frac{\boldsymbol{z}_i^T \boldsymbol{z}_j}{||\boldsymbol{z}_i|| ||\boldsymbol{z}_j||} \quad \forall i, j \in \mathcal{A} ,$$
$$\tag{9}$$

where $\gamma \in \mathbb{R}$ is the margin, $c_{ij}$ is the cosine similarity between the $i$-th and $j$-th columns of $\mathbf{Z}_{EXT}^A$ ($\boldsymbol{z}_i$ and $\boldsymbol{z}_j$), and $\mathcal{N}(i)$ denotes the set of actions that are functionally similar to action node $a_i$ (i.e., actions at distance 2 from $a_i$ in the affordance graph).

We train our model using Stochastic Gradient Descent (SGD) with momentum [24] and a fixed learning rate. Further details will be provided in Section 5.3.4.

---

[1]$\mathbf{M}$ does not need to be a square matrix because the graph is bipartite.

# 5. Experiments

We compare our results to the methods reported in [5] on HICO and VG-HOI. Although our work is focused on HOI Recognition (Section 5.3 and 5.4), we also consider the Detection task (Section 5.5), in which the model is required to localise each prediction. We perform Detection experiments on both HICO-DET and COCO-a.

## 5.1. Datasets

### 5.1.1 HICO and HICO-DET

The HICO dataset [3] and its bounding-box-annotated variant HICO-DET [25] comprise 47k images, COCO's 80 object classes [26], and 117 action classes, including a null one. They are annotated with 600 interactions and each image may belong to more than one interaction class. We follow the predefined train/test split of 38,116/9,658 images. Furthermore, we randomly sample 10% of the training set for validation in every run. In our Recognition experiment we follow [5], excluding the null action during training and testing and thus restricting the dataset to 116 actions and 520 interactions.

### 5.1.2 VG-HOI

VG-HOI [5] is a dataset for Human-Object Interaction built out of Visual Genome [27]. It comprises 10,799 train images and 4251 test images, for a total of 15,050. We use 10% of the training set for validation. There are 1392 objects, 495 actions and 6643 interactions, although for testing only the 532 that have at least 10 instances are used. The much larger number of classes (compared to HICO), together with the lower number of examples, make this dataset extremely challenging.

### 5.1.3 COCO-a

COCO-a [14] contains 4413 images annotated with 145 action classes and 80 object classes (same as COCO and HICO), for a total of 1681 interactions. We use it as an evaluation dataset for our model trained on HICO-DET, following the challenging setting used in [6].

## 5.2. Affordance Graph Construction

To build the affordance graph, we mine interactions from external knowledge bases and add them to the ones that can be found in the training set. Specifically, we use four external sources: Visual Genome [27] (except for VG-HOI), ActivityNet Captions [28], imSitu [29] and HCVRD [30]. The former three contain image or video captions that we parse into action-object pairs using NLTK [31] and the dependency parser from AllenNLP [32]. On the other hand, HCVRD is annotated with triplets in the form

$\langle subject, predicate, object \rangle$. We select the ones where the $subject$ is a person and $predicate$ is an action. Note that, in all cases, we do not add extra nodes into our graph and instead discard interactions containing actions or objects not in $\mathcal{A}$ and $\mathcal{O}$ (respectively).

## 5.3. Experimental Setting

### 5.3.1 Compared Models

We use four variants of our model: our baseline ($\rho = 0$ and $\lambda = 0$) and the models obtained by only adding one of the proposed loss components ($\rho > 0$ or $\lambda > 0$) or adding both ($\rho > 0$ and $\lambda > 0$).

The most similar method to ours is [5], which performs zero-shot learning on both action and objects. We compare our models to their best results, which are denoted by "GCNCL" followed by different endings based on how the knowledge graph is built. We also report other competitive methods from [5], namely Semantic Embedding Space (SES, [33]) and Triplet Siamese. We refer the reader to the corresponding papers for more details.

### 5.3.2 Evaluation

We use the standard mean Average Precision (mAP) as evaluation metric, reporting it as a percentage. We train every model multiple times (10 for HICO and 5 for VG-HOI), reporting the average result on the test set. We run Student's t-tests against current state-of-the-art results and all reported improvements are statistically significant at the 99% confidence interval.

### 5.3.3 Zero-Shot Settings

In order to make a fair comparison, we use the same seen/unseen splits as Task 2 from [5]: the training set is made of 49 objects and 53 actions for HICO and 554 objects and 198 actions for VG-HOI. At test time all classes are included, following the *Generalised* Zero-Shot Learning setting.

### 5.3.4 Implementation Details

We use a ResNet-152 pre-trained on ImageNet [34] as image feature extractor (same as [5]). Functions $f_1$, $f_2$ and $f_3$ are implemented by two fully-connected layers with output dimensions both equal to 1024, with ReLU non-linearity. After the non-linearity we add Dropout [35] (at a 0.5 rate) for $f_1$ and $f_3$, but not for $f_2$, as suggested in [6]. We use Glorot initialisation [36] to initialise the optimisation parameters in Equation (5). Our GCN comprises two convolutional layers with output dimension 1024, the first of which is equipped with ReLU and Dropout (0.5 rate).

| Method | All | Unseen only |
|---|---|---|
| Triplet Siamese | 10.38 | 7.76 |
| SES | 11.69 | 7.19 |
| GCNCL-I | 11.93 | 7.22 |
| GCNCL+NV+A | 11.94 | 7.50 |
| Ours | 13.79 | 6.93 |
| Ours, $\lambda = 1$ | **16.02** | 10.08 |
| Ours, $\rho = 10$ | 14.02 | 7.16 |
| Ours, $\lambda = 1, \rho = 10$ | **16.02** | **10.20** |

Table 1: Results on HICO.

| Method | All | Unseen only |
|---|---|---|
| Triplet Siamese | 2.55 | 1.67 |
| SES | 2.07 | 0.96 |
| GCNCL-I+A | 4.00 | 2.63 |
| GCNCL+A | 4.07 | 2.44 |
| Ours | 4.90 | 3.51 |
| Ours, $\lambda = 0.1$ | 5.09 | 3.77 |
| Ours, $\rho = 100$ | **5.11** | **3.90** |
| Ours, $\lambda = 0.1, \rho = 100$ | 5.01 | 3.74 |

Table 2: Results on VG-HOI.

We keep the margin parameter $\gamma$ in Equation (9) fixed at 0.3, whereas we experiment with different values of $\rho$ and $\lambda$ for the two datasets. The best ones (according to validation results) are the ones shown in the respective tables.

We use GloVe [37] for our word embeddings. More specifically, we use the 300-dimensional embeddings trained on Gigaword and Wikipedia[2] and we normalise them. For compound words, we take the average of the components.

Finally, we train our model using minibatch Stochastic Gradient Descent (SGD) with momentum. We use a fixed learning rate of 0.001 and set the momentum and weight decay coefficients to 0.9 and $5 \cdot 10^{-4}$, respectively. We train our model for a maximum of 100 epochs on HICO and 150 on VG-HOI, with early stopping based on validation accuracy. We use a batch size of 64.

## 5.4. Results

### 5.4.1 Results on HICO

Our results are summarised in Table 1. We see that our baseline model already compares very favourably to all the existing approaches, and adding either or both of the proposed losses upgrades our baseline's performance considerably. The best performing model, obtained with $\lambda = 1$ and $\rho = 10$, gains more than 4% over the current state of the art (GCNCL+NV+A) for the whole test set and around

---
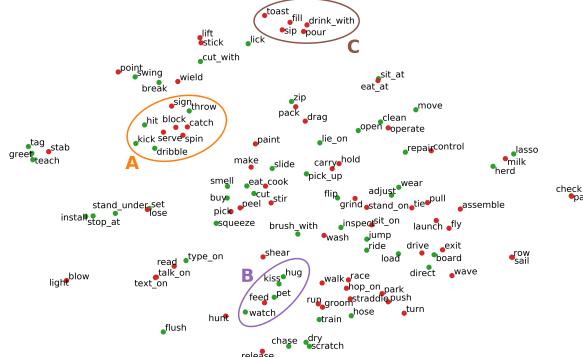
[2]Available at https://nlp.stanford.edu/projects/glove.

Figure 3: Visualisation of action class representations. Green dots represent seen actions and red dots unseen ones. Some clusters are highlighted: (*A*) sports actions (e.g., *catch*, *throw*), (*B*) actions regarding domesticated animals (e.g., *pet*, *feed*) and (*C*) action involving cups or glasses (e.g., *sip*, *pour*).

2.7% for unseen classes only. This corresponds to sizeable ∼35% relative increases. It is worth mentioning that the graph building process results in 68 missing interactions out of HICO's 520, since they cannot be mined from our external sources. Despite this, no object is completely isolated in the affordance graph, whereas only 5 actions are (*hop_on*, *hunt*, *lose*, *stab*). Most of these actions (namely *hunt*, *lose*, *stab* and *toast*) are too niche to be found in the other sources, and in fact even in HICO they only appear in one interaction each. On the other hand, *hop_on* can be found, but not with the meaning of "jumping on a ride" it has in HICO (and thus it is not paired with the same objects). Nonetheless, our model still performs very well, possibly due to the fact that additional interactions are added and they contribute to meaningful representations being learnt, even though they do not appear in HICO.

### 5.4.2 Results on VG-HOI

Results are reported in Table 2. Our baseline is better than previous models, GCNCL+A in particular: ∼.83% for all classes (∼20% relative gain) and almost 1.1% for unseen categories, corresponding to a remarkable 40% relative improvement. Adding the proposed losses improves performance, with the best one obtained by setting $\lambda = 0, \rho = 100$. While our losses improve results, they are not as effective on this dataset as on HICO. We believe this can be ascribed to the vast number of unseen categories: while in HICO there are 80 objects and 116 actions, VG-HOI contains ∼17 times as many objects and more than 4 times as many actions. The sheer number of unseen classes makes classification much more difficult; in particular, our method relies on seen object labels to estimate unseen action ones



(a)
✓ *eat_at* dining_table
✓ sit_at dining_table

(b)
✓ *hold* book
✓ open book

(c)
✓ *carry* umbrella
✓ *hold* umbrella
✗ stand_under umbrella

(d)
✓ *pull* tie
✓ *tie* tie
✓ wear tie

Figure 4: Some predictions of our best model on HICO. Marks indicate whether the prediction matches the ground truth ✓ or not ✗ . Actions in *italic* are unseen.

(see Equation (7)), therefore missing a large amount of information about objects is detrimental. The incompleteness of the affordance graph is also likely to negatively affect performance, as the graph only covers 2753 interactions (∼41%), 291 actions (∼59%) and 806 objects (∼58%). Despite these difficulties, our method still performs significantly better than previous approaches.

### 5.4.3 Qualitative Results on HICO

We show some predictions on HICO's test set examples in Figure 4, demonstrating that our model is able to correctly predict several previously unseen actions. We also show the representation space in Figure 3 using t-SNE [38] on a model trained with both proposed losses ($\lambda, \rho > 0$). Some clusters are clearly identifiable, such as cluster A, which contains actions such as *catch*, *throw* or *spin* that can be performed on small sport items like *sports_ball* or *frisbee*. This shows that the proposed approach is effective in grouping actions based on their affordance. Comparisons between representation spaces obtained in different settings can be found in the supplementary material.

### 5.5. Zero-Shot HOI Detection

We used different settings for the Detection experiments, which can be found in the supplementary material. We present the results in the following.

| Method | All |
|---|---|
| Shen *et al.* [4] | 6.46 |
| Chao *et al.* [25] | 7.81 |
| InteractNet [39] | 9.94 |
| GPNN [40] | 13.11 |
| Xu *et al.* [15] | 14.70 |
| iCAN [41] | 14.84 |
| Song *et al.* [42] | 15.27 |
| Wang *et al.* [43] | 16.24 |
| No-frills [44] | 17.18 |
| Li *et al.* [45] | 17.22 |
| RPNN [46] | 17.35 |
| PMFNet [47] | 17.46 |
| Peyre *et al.* [6] | 19.40 |
| Wang *et al.* [48] | 19.56 |
| PPDM [49] | 21.73 |
| Bansal *et al.* [7] | **21.96** |
| Ours | 18.74 |

Table 3: Results on HICO-DET in a fully supervised setting.

| Method | All | Unseen |
|---|---|---|
| Ours | 11.18 | 8.19 |
| Ours, $\lambda = 0.1, \rho = 0.1$ | **11.94** | **9.81** |

Table 4: Baseline for ZS HOI Detection on HICO-DET.

### 5.5.1 HICO-DET

While there are works on Zero-Shot Learning on HICO-DET for interactions [4, 6] and objects [7], no previous approach has dealt with zero-shot actions (to the best of our knowledge). We provide in Table 4 a baseline for future reference. We also show in Table 3 how our approach compares against other methods in a fully supervised setting as a reference, where we can see that there is a noticeable increase in mAP with respect to most methods in the literature. It is worth mentioning that some of the techniques that likely contribute to the outstanding results of the top three methods, such as fine-tuning the object detector on HICO-DET [7] or following a more intensive training regime while fine-tuning the feature extractor (50 epochs on 5 GPUs for [48], 110 epochs on 8 GPUs for [49]), are applicable to our model as well – in fact, Bansal *et al.* report that their method only achieves 16.96% mAP without such fine-tuning. We leave this for future work.

### 5.5.2 COCO-a

In Table 5 we compare our results on COCO-a (reporting our baseline plus the best hyperparameter setting for each column) against a state-of-the-art approach using the chal-

| Method | | Unseen HOIs | |
|---|---|---|
| | All | With unseen action |
| Peyre *et al.* [6] (best) | 6.9 | 7.3 |
| Ours | 9.65 | 11.00 |
| Ours, $\lambda = 0.1$ | 9.93 | **11.44** |
| Ours, $\lambda = 0.1, \rho = 10$ | **10.01** | 11.13 |

Table 5: Results on COCO-a.

lenging setting described in their paper: train on HICO-DET and evaluate on COCO-a. Under this setting, there are 1474 unseen interactions, 1048 of which involve an unseen action. Our approach performs much better than the best one from [6], gaining around 2.7 points for all unseen interactions ($\sim$40% relative gain) and 3.7 points when dealing with interactions involving unseen actions (about 50% relative gain). Performance improve even further when setting $\lambda$ and/or $\rho$ to non-zero values (the best assignment for each measure is reported).

## 6. Conclusion

We have proposed an effective approach that uses structured knowledge in the form of an affordance graph to improve Zero-Shot Human-Object Interaction Recognition. The proposed model learns regularised representations of unseen classes in a weakly supervised way using labels which are estimated through the affordance graph, while simultaneously learning representation of seen classes in a supervised fashion. Our method is able to predict unseen interactions in the very challenging case where only about half of the object and action classes are seen during training. We evaluate our results on several datasets (including standard benchmarks like HICO and HICO-DET) and show that our approach performs significantly better than the current state of the art.

## References

[1] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang, "Bottom-Up and Top-Down Attention for Image Captioning and VQA," *arXiv preprint arXiv:1707.07998*, 2017.

[2] Z. Fang, J. Yuan, and N. Magnenat-Thalmann, "Understanding Human-Object Interaction in RGB-D videos for Human Robot Interaction," in *Proceedings of Computer Graphics International 2018*, pp. 163–167, 2018.

[3] Y.-W. Chao, Z. Wang, Y. He, J. Wang, and J. Deng, "Hico: A benchmark for recognizing human-object interactions in images," in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1017–1025, 2015.

[4] L. Shen, S. Yeung, J. Hoffman, G. Mori, and L. Fei-Fei, "Scaling Human-Object Interaction Recognition through Zero-Shot Learning," in *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pp. 1568–1576, IEEE, 2018.

[5] K. Kato, Y. Li, and A. Gupta, "Compositional learning for human object interaction," in *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 234–251, 2018.

[6] J. Peyre, I. Laptev, C. Schmid, and J. Sivic, "Detecting Unseen Visual Relations Using Analogies," in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1981–1990, 2019.

[7] A. Bansal, S. S. Rambhatla, A. Shrivastava, and R. Chellappa, "Detecting Human-Object Interactions via Functional Generalization.," in *AAAI*, pp. 10460–10469, 2020.

[8] D. Norman, *The design of everyday things: Revised and expanded edition*. Basic books, 2013.

[9] L. L. Chao and A. Martin, "Representation of Manipulable Man-Made Objects in the Dorsal Stream," *NeuroImage*, vol. 12, pp. 478–484, Oct. 2000.

[10] L. Stark and K. Bowyer, "Achieving generalized object recognition through reasoning about association of function to structure," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 13, pp. 1097–1104, Oct. 1991. Conference Name: IEEE Transactions on Pattern Analysis and Machine Intelligence.

[11] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *arXiv preprint arXiv:1609.02907*, 2016.

[12] X. Wang, Y. Ye, and A. Gupta, "Zero-shot recognition via semantic embeddings and knowledge graphs," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6857–6866, 2018.

[13] J. Gao, T. Zhang, and C. Xu, "I Know the Relationships: Zero-Shot Action Recognition via Two-Stream Graph Convolutional Networks and Knowledge Graphs," 2019.

[14] M. R. Ronchi and P. Perona, "Describing Common Human Visual Actions in Images," *arXiv:1506.02203 [cs]*, June 2015. arXiv: 1506.02203.

[15] B. Xu, Y. Wong, J. Li, Q. Zhao, and M. S. Kankanhalli, "Learning to Detect Human-Object Interactions With Knowledge," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.

[16] A. Frome, G. S. Corrado, J. Shlens, S. Bengio, J. Dean, and T. Mikolov, "Devise: A deep visual-semantic embedding model," in *Advances in neural information processing systems*, pp. 2121–2129, 2013.

[17] L. Zhang, T. Xiang, and S. Gong, "Learning a deep embedding model for zero-shot learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2021–2030, 2017.

[18] A. Mishra, V. K. Verma, M. S. K. Reddy, S. Arulkumar, P. Rai, and A. Mittal, "A generative approach to zero-shot and few-shot action recognition," in *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pp. 372–380, IEEE, 2018.

[19] E. Schnfeld, S. Ebrahimi, S. Sinha, T. Darrell, and Z. Akata, "Generalized Zero-and Few-Shot Learning via Aligned Variational Autoencoders," *arXiv preprint arXiv:1812.01784*, 2018.

[20] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *Proceedings of the IEEE international conference on computer vision*, pp. 1026–1034, 2015.

[21] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.

[22] X. Li, Y. Guo, and D. Schuurmans, "Semi-supervised zero-shot classification with label representation learning," in *Proceedings of the IEEE international conference on computer vision*, pp. 4211–4219, 2015.

[23] A. Mallya and S. Lazebnik, "Learning models for actions and person-object interactions with transfer to question answering," in *European Conference on Computer Vision*, pp. 414–428, Springer, 2016.

[24] D. E. Rumelhart, G. E. Hinton, R. J. Williams, and others, "Learning representations by back-propagating errors," *Cognitive modeling*, vol. 5, no. 3, p. 1, 1988.

[25] Y.-W. Chao, Y. Liu, X. Liu, H. Zeng, and J. Deng, "Learning to Detect Human-Object Interactions," in *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pp. 381–389, IEEE, 2018.

[26] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollr, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European conference on computer vision*, pp. 740–755, Springer, 2014.

[27] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma, M. S. Bernstein, and L. Fei-Fei, "Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations," *International Journal of Computer Vision*, vol. 123, pp. 32–73, May 2017.

[28] R. Krishna, K. Hata, F. Ren, L. Fei-Fei, and J. Carlos Niebles, "Dense-captioning events in videos," in *Proceedings of the IEEE international conference on computer vision*, pp. 706–715, 2017.

[29] M. Yatskar, L. Zettlemoyer, and A. Farhadi, "Situation recognition: Visual semantic role labeling for image understanding," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5534–5542, 2016.

[30] B. Zhuang, Q. Wu, C. Shen, I. D. Reid, and A. van den Hengel, "HCVRD: A Benchmark for Large-Scale Human-Centered Visual Relationship Detection.," in *AAAI*, 2018.

[31] S. Bird, E. Klein, and E. Loper, *Natural language processing with Python: analyzing text with the natural language toolkit.* " O'Reilly Media, Inc.", 2009.

[32] M. Gardner, J. Grus, M. Neumann, O. Tafjord, P. Dasigi, N. Liu, M. Peters, M. Schmitz, and L. Zettlemoyer, "Allennlp: A deep semantic natural language processing platform," *arXiv preprint arXiv:1803.07640*, 2018.

[33] X. Xu, T. Hospedales, and S. Gong, "Semantic embedding space for zero-shot action recognition," in *2015 IEEE International Conference on Image Processing (ICIP)*, pp. 63–67, IEEE, 2015.

[34] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pp. 248–255, IEEE, 2009.

[35] G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R. R. Salakhutdinov, "Improving neural networks by preventing co-adaptation of feature detectors," *arXiv preprint arXiv:1207.0580*, 2012.

[36] 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*, pp. 249–256, 2010.

[37] J. Pennington, R. Socher, and C. Manning, "Glove: Global vectors for word representation," in *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pp. 1532–1543, 2014.

[38] L. v. d. Maaten and G. Hinton, "Visualizing data using t-SNE," *Journal of machine learning research*, vol. 9, no. Nov, pp. 2579–2605, 2008.

[39] G. Gkioxari, R. Girshick, P. Dollr, and K. He, "Detecting and recognizing human-object interactions," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8359–8367, IEEE, 2018.

[40] S. Qi, W. Wang, B. Jia, J. Shen, and S.-C. Zhu, "Learning human-object interactions by graph parsing neural networks," in *European Conference on Computer Vision*, pp. 407–423, Springer, 2018.

[41] C. Gao, Y. Zou, and J.-B. Huang, "ican: Instance-centric attention network for human-object interaction detection," *arXiv preprint arXiv:1808.10437*, 2018.

[42] Y. Song, W. Li, L. Zhang, J. Yang, E. Kiciman, H. Palangi, J. Gao, C.-C. J. Kuo, and P. Zhang, "Novel Human-Object Interaction Detection via Adversarial Domain Generalization," *arXiv:2005.11406 [cs]*, May 2020. arXiv: 2005.11406.

[43] T. Wang, R. M. Anwer, M. H. Khan, F. S. Khan, Y. Pang, L. Shao, and J. Laaksonen, "Deep Contextual Attention for Human-Object Interaction Detection," *arXiv:1910.07721 [cs]*, Oct. 2019. arXiv: 1910.07721.

[44] T. Gupta, A. Schwing, and D. Hoiem, "No-Frills Human-Object Interaction Detection: Factorization, Layout Encodings, and Training Techniques," in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 9677–9685, 2019.

[45] Y.-L. Li, S. Zhou, X. Huang, L. Xu, Z. Ma, H.-S. Fang, Y. Wang, and C. Lu, "Transferable Interactiveness Knowledge for Human-Object Interaction Detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3585–3594, 2019.

[46] P. Zhou and M. Chi, "Relation Parsing Neural Network for Human-Object Interaction Detection," in

*Proceedings of the IEEE International Conference on Computer Vision*, pp. 843–851, 2019.

[47] B. Wan, D. Zhou, Y. Liu, R. Li, and X. He, "Pose-aware Multi-level Feature Network for Human Object Interaction Detection," *arXiv:1909.08453 [cs]*, Sept. 2019. arXiv: 1909.08453.

[48] T. Wang, T. Yang, M. Danelljan, F. S. Khan, X. Zhang, and J. Sun, "Learning Human-Object Interaction Detection Using Interaction Points," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, (Seattle, WA, USA), pp. 4115–4124, IEEE, June 2020.

[49] Y. Liao, S. Liu, F. Wang, Y. Chen, C. Qian, and J. Feng, "PPDM: Parallel Point Detection and Matching for Real-Time Human-Object Interaction Detection," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, (Seattle, WA, USA), pp. 479–487, IEEE, June 2020.

[50] K. He, G. Gkioxari, P. Dollr, and R. Girshick, "Mask r-cnn," in *Computer Vision (ICCV), 2017 IEEE International Conference on*, pp. 2980–2988, IEEE, 2017.

## A. Supplementary Material

This supplementary material contains:

1. Details about the hardware and software infrastructure used to implement the method;

2. Several ablation studies (Appendix A.2) and additional experiments (Appendix A.3);

3. Details about the HOI Detection experiment described in Section 5.5 (Appendix A.4);

4. A visual comparison of HICO's action representation spaces obtained by different means (Appendix A.5) to expand on what we show in Section 5.4.

### A.1. Infrastructure Details

The model has been implemented in Python 3.6 using PyTorch v0.4.1. Experiments have been run on a single NVIDIA GeForce GTX TITAN X GPU on a server with an Intel(R) Core(TM) i7-5930K CPU and 64GB of RAM running CentOS Linux 7. The code will be made available upon publication.

### A.2. Ablation experiments

In this section we describe several ablation experiments. All the reported results are computed on HICO.

| Method | With | | Without | |
|---|---|---|---|---|
| | All | Unseen | All | Unseen |
| $\lambda = \rho = 0$ | 13.79 | 6.93 | 13.11 | 6.75 |
| $\lambda = 1$ | 16.02 | 10.08 | 15.14 | 9.85 |
| $\rho = 10$ | 14.02 | 7.16 | 13.32 | 6.94 |
| $\lambda = 1, \rho = 10$ | 16.02 | 10.20 | 15.14 | 9.95 |

Table 6: Ablation: internal representation for seen classes.

| Method | Objects only | | Objects and actions | |
|---|---|---|---|---|
| | All | Unseen | All | Unseen |
| $\lambda = \rho = 0$ | 13.79 | 6.93 | 13.31 | 6.20 |

Table 7: Ablation: word embeddings.

#### A.2.1 Alternative representation

As described in Section 4.2, we train an alternative representation for seen classes that we call *internal representation*, denoted by $\mathbf{Z}_{INT}^q$ for $q \in \{O, A\}$. The rationale for this is that the semi-supervised training used to train representation of unseen classes via the GCN might introduce noise in the representation of seen classes, which could be trained in a fully-supervised fashion thank to the availability of instance labels. Results shown in Table 6 corroborate our hypothesis: models that learn the internal representation for seen classes perform consistently better than the corresponding ones which only learn GCN representations. All results are statistically significant at the 95% confidence interval.

#### A.2.2 Word embeddings

In Section 4.2 we argued that word embeddings are not well-suited to provide affordance information about actions, because word embeddings relate words based on co-occurrence in a sentence. Therefore, while they can capture affordance-based similarity for objects (e.g., in the sentence "I eat an apple and a banana", objects *apple* and *banana* co-occur because they both afford action *eating*), this effect is weaker for actions (e.g., in the sentence "People were eating and drinking" the actions co-occur not because they are afforded by the same objects, but rather because they can be performed in the same context). To empirically verify this intuition, we perform an ablation experiment whose result are shown in Table 7. Differences are statistically significant at the 99% confidence interval. These results justify why we do not add a component based on word embeddings in Equation (2).

### A.3. Sensitivity Experiments

In this section we evaluate how sensitive our model is to the available information, in particular to the amount of

| Method | All | Unseen |
|---|---|---|
| $\lambda = 1, \mu = 1$ (Table 1) | 16.02 | 10.08 |
| $\lambda = 1, \mu = 0.6$ | 20.91 | 12.72 |
| $\lambda = 1, \mu = 0.3$ | **27.19** | **18.21** |

Table 8: Sensitivity of the model to the amount of unseen labels.

| Method | All | Unseen |
|---|---|---|
| $\lambda = 1, \nu = 1$ (Table 1) | **16.02** | **10.08** |
| $\lambda = 1, \nu = 0.8$ | 15.43 | 9.34 |
| $\lambda = 1, \nu = 0.6$ | 14.50 | 8.00 |

Table 9: Sensitivity of the model to the completeness of the affordance graph.

unseen labels and completeness of the affordance graph.

### A.3.1 Amount of unseen labels

We perform an experiment to evaluate how much the amount of unseen labels impacts performance. Specifically, we define a hyperparameter $\mu \in [0, 1]$ as the ratio of unseen classes with respect to the experiment reported in Table 1. Thus, a $\mu = 1$ corresponds to the reported experiment (31 unseen object classes and 63 unseen action classes), $\mu = 0.6$ means to keep around 60% of the unseen classes (for a total of 19/38 unseen object/action classes, while the remaining 40% are added to the seen classes) and a value of $\mu = 0.3$ means to only keep around 30% of the unseen classes. Results are shown in Table 8 and show that the amount of unseen labels greatly affects performance.

### A.3.2 Completeness of the affordance graph

Our approach makes extensive use of the affordance graph. It is natural to assume that a sparser graph leads to worse results, since it contains less information, thus we design an experiment to verify this assumption. In particular, we define a hyperparameter $\nu \in [0, 1]$ as the proportion of edges of the affordance graph with respect to one used in the experiments reported in Table 1. For instance, $\nu = 0.8$ means to sample around 80% of the edges to keep and remove the remaining 20%. Results can be viewed in Table 9 and they show that the more sparse the affordance graph, the lower the performance.

### A.4. Details about the HOI Detection Experiment

In this section we describe the settings that we used for our HOI Detection experiment on HICO-DET and COCO-a (Section 5.5), which differ from the settings of the HOI Recognition experiments (Section 5.3).

### A.4.1 Experimental Setup

The focus of this experiments is Zero-Shot HOI Detection when there are unseen actions. On HICO-DET our training set contains the same unseen actions as the recognition experiment (∼50% of the total, as described in Section 5.3), while on COCO-a there are 114 unseen actions, corresponding to 1048 unseen interactions. In both cases there are no unseen object classes, therefore the object branch is removed: since we do not perform zero-shot on objects, we can rely purely on the scores provided by a pre-trained object detector (we use Mask R-CNN [50] with ResNet-50 [21] as backbone). Note that this is possible because we use a model pre-trained on COCO [26], which has the same object categories as HICO-DET and COCO-a.

### A.4.2 Architectural Changes

Contrary to HICO, HICO-DET and COCO-a contain localised information: each interaction in an image refers to a specific person and object, and a bounding box for each is provided. Therefore, we adapted our model to deal with image regions instead of whole images. The object detector provides visual features for every person, every object and every region that represents a possible interaction (i.e., the tightest region that contains both a person and an object). This means that, for each example $i$, we have three visual feature vectors: $\boldsymbol{h}_i^{(h)}$, $\boldsymbol{h}_i^{(o)}$ and $\boldsymbol{h}_i^{(a)}$ for human, object and action respectively. We compute the interaction representation as

$$\boldsymbol{x}_i = f_1([\boldsymbol{h}_i^{(h)}, \boldsymbol{s}_i^{(h)}]) + f_1([\boldsymbol{h}_i^{(o)}, \boldsymbol{s}_i^{(o)}]) + f_1(\boldsymbol{h}_i^{(a)}) , \quad (10)$$

where $f_1$ is defined as usual as an MLP, $[\cdot, \cdot]$ indicates concatenation and $\boldsymbol{s}^{(\cdot)}$ are object classification score vectors returned by the object detector.

### A.4.3 Sampling Interactions

During training, we keep all detected object bounding boxes and add the ground-truth ones that do not have any match, i.e., there is no detected box whose intersection-over-union (IoU) is greater than 0.5. We keep as positive interaction examples all human-object pairs whose subject and object are correctly classified and overlap with the subject/object (respectively) of a ground-truth interaction (again, the threshold for IoU is 0.5). Among the pairs that are not positive interactions, we sample negative ones, at a rate of 3 negatives per positive (this is a widely used ratio, see for example [6]). At inference time, we only keep human candidates with a confidence score greater than 0.7 and threshold object ones at 0.3. Every possible human-object pair in the image is considered as a candidate interaction and classified by the model.

### A.4.4 Changes to the Training Procedure

When using the regularisation loss $\mathcal{L}_{REG}$ on HICO-DET, we found it beneficial to only enable it (that is, set $\rho > 0$) after the first 5 epochs. This allows the model to learn class representations first, and only later regularise them.

The model is trained with minibatch Stochastic Gradient Descent (SGD) with a learning rate of $0.001$ and weight decay coefficient of $5 \cdot 10^{-4}$. We train our model for a maximum of 10 epochs (due to the high amount of training samples: more than 1.2M interactions, compared to the $\sim$30k training images for HICO) and a batch size of 64, 75% of which is constituted by negative samples as previously mentioned.
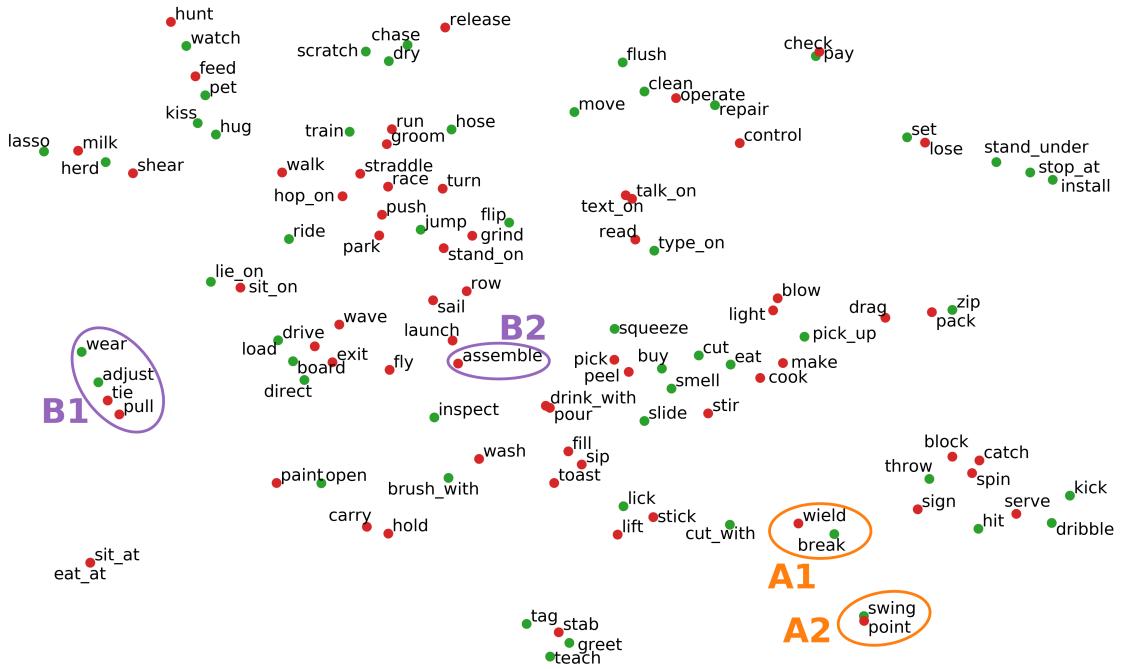
## A.5. Visualisation of Representation Spaces

In this section we show how the representation spaces vary depending on whether our regularisation loss $\mathcal{L}_{REG}$ is used (Appendix A.5.1).

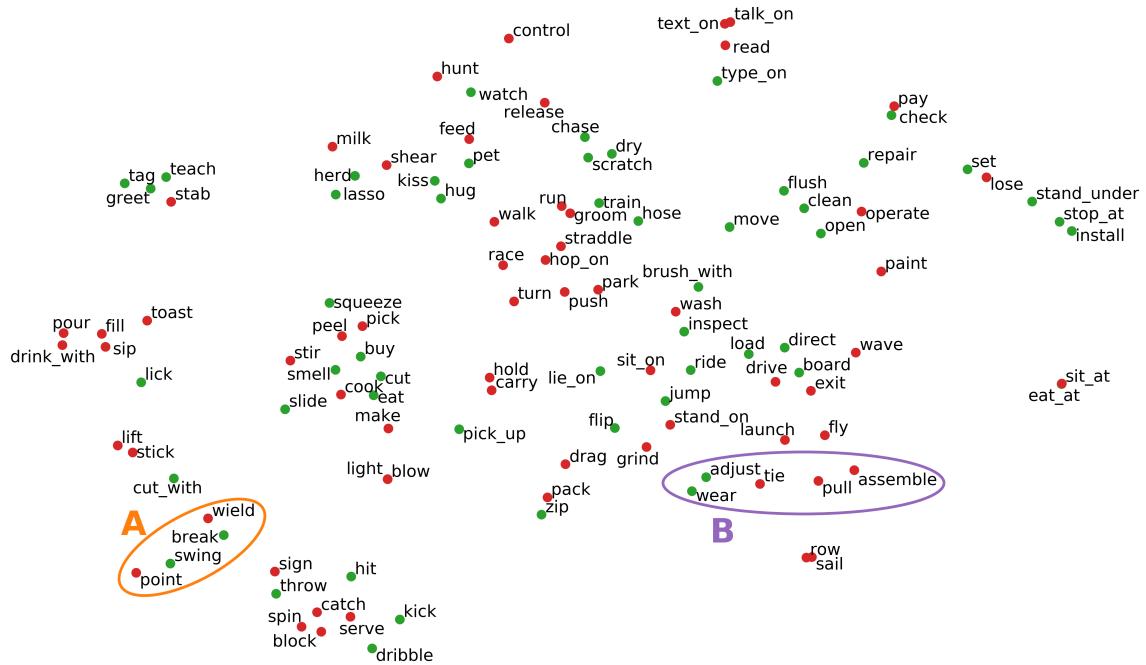### A.5.1 Effect of Regularisation

Figure 5 shows action representations obtained with our model on HICO without (Figure 5a) and with (Figure 5b) regularisation. It can be seen that the base model is already quite effective at grouping actions by affordances, as representations are computed through a GCN over the affordance graph. Adding the proposed regularisation further promotes clustering based on functional similarity, i.e., it tends to group together actions based on what objects they can be performed on. For instance, let us consider the group $wield$, $break$, $point$ and $swing$. Three of them ($wield$, $break$ and $swing$) can be performed on a $baseball\_bat$, while two ($swing$ and $point$) can be performed on a $remote\_control$. In the unregularised model (Figure 5a) $swing$ is correctly clustered with $point$ (A2), but quite distant from $wield$ and $break$ (A1), whereas the regularisation brings the two groups closer to each other, effectively merging them (A). This effect is magnified for $pull$ and $assemble$, which can both be performed on a $kite$. In the base model, $pull$ is only grouped with actions that can be performed on a $tie$ (B1), but the regularisation helps in bringing the two clusters (B1 and B2) together because they share a common action (B).

### A.5.2 Are Affordances Captured by Word Embeddings?

We show a comparison between the learnt representation space and the word embedding space in Figure 6. The figure shows that actions are not clustered by affordance in the word embedding space, further confirming the efficacy of our approach.

(a) Unregularised model

(b) Regularised model

Figure 5: Comparison between unregularised (a) and regularised (b) model. The latter is better at grouping functionally similar actions. Best seen in colour.

(a) Learnt representation
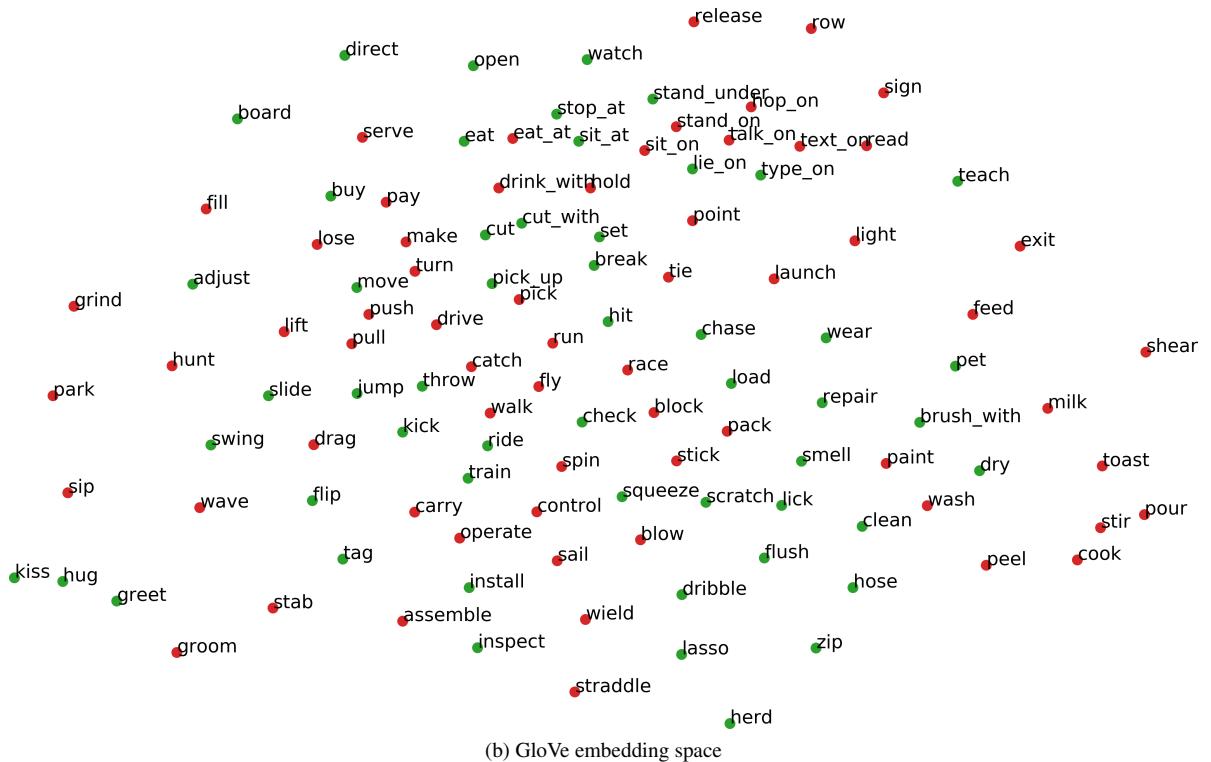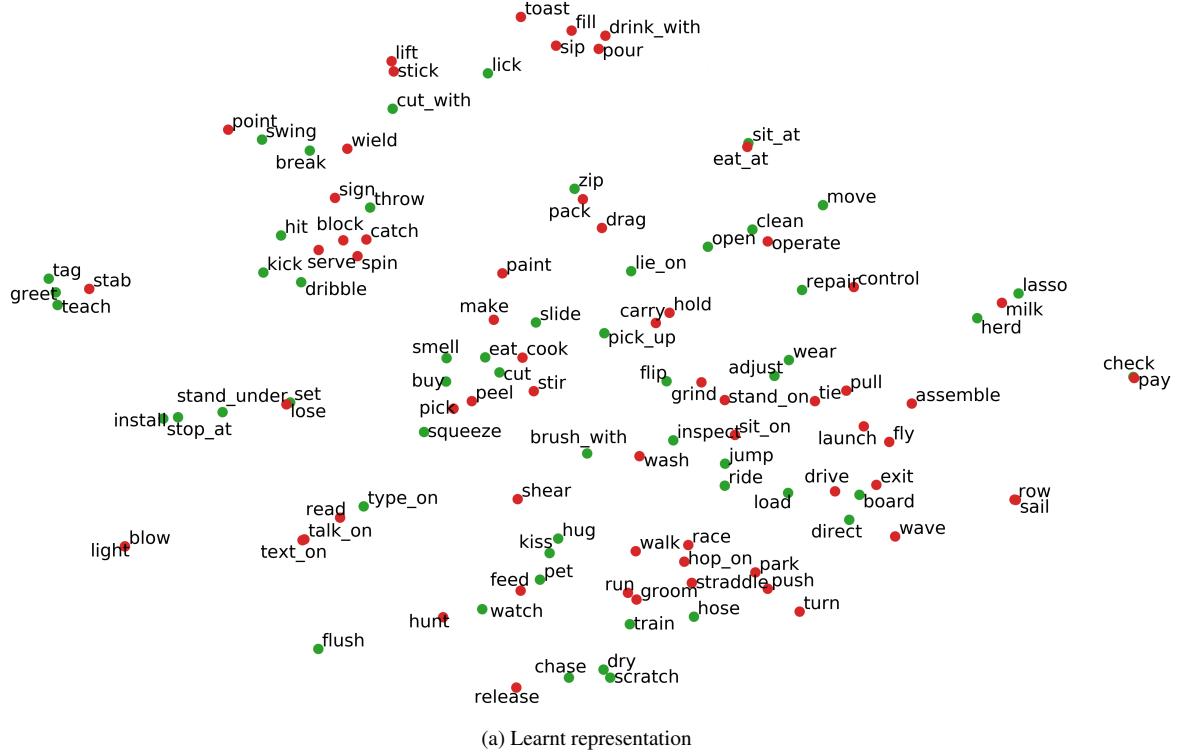


(b) GloVe embedding space

Figure 6: Comparison between the learnt representation space (a) and the pre-trained word embedding space (b). The proposed approach is effective at grouping actions by affordances, while GloVe embeddings do not capture this relationship. Best seen in colour.