

Classification by Attention: Scene Graph Classification with Prior Knowledge

Sahand Sharifzadeh,¹ Sina Moayed Baharlou,^{1*} Volker Tresp^{1,2}

¹Ludwig Maximilian University of Munich

²Siemens AG

sharifzadeh@dbs.ifi.lmu.de

Abstract

A major challenge in scene graph classification is that the appearance of objects and relations can be significantly different from one image to another. Previous works have addressed this by relational reasoning over all objects in an image or incorporating prior knowledge into classification. Unlike previous works, we do not consider separate models for perception and prior knowledge. Instead, we take a multi-task learning approach by introducing *schema* representations and implementing the classification as an attention layer between image-based representations and the schemata. This allows for the prior knowledge to emerge and propagate within the perception model. By enforcing the model also to represent the prior, we achieve a strong inductive bias. We show that our model can accurately generate commonsense knowledge and that the iterative injection of this knowledge to scene representations, as a top-down mechanism, leads to significantly higher classification performance. Additionally, our model can be fine-tuned on external knowledge given as triples. When combined with self-supervised learning and with 1% of annotated images only, this gives more than 3% improvement in object classification, 26% in scene graph classification, and 36% in predicate prediction accuracy.

Introduction

Classifying objects and their relations in images, also known as scene graph classification, is a fundamental task in scene understanding and can play an essential role in applications such as recommender systems, visual question answering and decision making. Scene graph (SG) classification methods typically have a perception model that takes an image as input and generates a graph that describes the given image as a collection of (head, predicate, tail). One of the main challenges that current models face is diverse appearances of objects and relations across different images. This can be due to variations in lighting conditions, viewpoints, object poses, occlusions, etc. For example, the *Bowl* in Figure 1 is highly occluded and has very few image-based features. Therefore, a typical perception model fails to classify it. One approach to tackle this problem is to collect supportive evidence from the neighbors before classifying an entity. This

*S.M. Baharlou contributed to this project while he was a visiting researcher at the Ludwig Maximilian University of Munich. Copyright © 2021, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

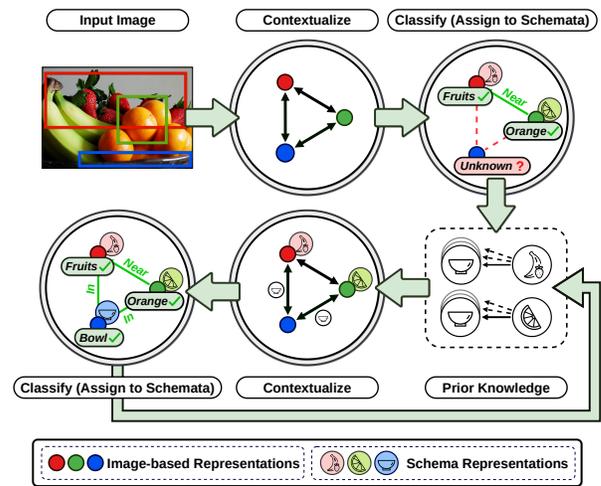


Figure 1: An example of scene graph classification where the *Bowl* lacks sufficient visual input. The top right is the initially predicted graph from the visual inputs only. The bottom left is the prediction of our model after considering both image-based representations and the prior knowledge about *Fruits* and *Oranges* (schemata). The long arrow near the bottom indicates recursion.

can be done, for example, by message passing between all the image-based object representations in an image, using graph convolutional neural networks (GCN) (Kipf and Welling 2016) or LSTMs (Hochreiter and Schmidhuber 1997). The main issue with this approach is the combinatorial explosion of all possible *image-based* neighbor representations¹.

A current theory in cognitive psychology states that humans solve this challenge by reasoning over the pre-existing representations of neighboring objects instead of relying on the perceptual inputs only (Piaget 1923); philosophers often argue that humans have a form of mental representation for objects and concepts (Kant 1787). These representations do not depend on a given image but are rather *symbol-based*. There are different opinions on how these representations come to be. Piaget called these representations *schema* (plu-

¹For a more detailed probabilistic analysis of this issue, refer to the section *GCN vs. Prior Model: A matter of inductive biases*.

ral *schemata*) and suggested that we acquire them in our earlier perceptions. When an object is being perceived, the mind assigns it to a schema in a process called *assimilation*. By relational reasoning over schemata, assimilation helps to predict the facts surrounding the observation (Arbib 1992)².

Nevertheless, learning and utilizing prior knowledge is still a significant challenge in AI research. Recently, Zellers et al. (2018) and Chen et al. (2019c,b) proposed to correct the SG prediction errors by comparing them to the co-occurrence statistics of triples in the training dataset or an external source. The statistics can, for example, suggest that it is common to see (Fruits, in, Bowls). Furthermore, instead of relying on simple co-occurrence statistics, one can create a prior model with knowledge graph embeddings (KGE) (Nickel et al. 2016) that can generalize beyond the given triples (Baier, Ma, and Tresp 2017, 2018; Hou et al. 2019). KGEs typically consist of an *embedding matrix*, that assigns a vector representation to each entity, and an *interaction function* that predicts the probability of a triple given the embeddings of its head, predicate, and tail. This allows KGE models to generalize to unseen relations. For example, Man and Woman will be given similar embeddings since they appear in similar relations during the training. As a result, if the model observes (Woman, rides, Horse), it can generalize, for example, to (Man, rides, Camel).

However, in the described approaches, unlike Piaget’s schemata, the perception and the prior are treated naively as independent components; they are trained separately and from different inputs (either from images or from triples), and their predictions are typically fused in the probability space, e.g., by multiplication. Other than requiring redundant parameters and computations, this makes the prior model agnostic to the image-based knowledge and the perception model agnostic to the prior knowledge. For example, the collection of triples might contain (Woman, rides, Horse) but have no triples regarding a Donkey. While the images can represent the visual similarities of a Horse to a Donkey, the triples lack this information. If we train a prior model purely based on the triples, the model fails to generalize. We can avoid this by training the prior model from a combination of triples and images. As for another example, in Figure 1, the prior model might suggest (Fruits, in, Bowl) but it might also suggest (Fruits, areOn, Trees). To decide between the two, one should still consider the visual cues from the given image.

To address these shortcomings, we entangle the perception and prior in a single model with shared parameters trained by multi-task learning. Therefore, instead of training a separate embedding matrix for a prior model, we exploit the perception model’s classification layer; when we train a classification layer on top of contextualized image-based representations, the classification weights capture both relational and image-based class embeddings (Refer to Figure 2). Unfortunately, the classification’s common realization as a fully connected layer does not allow us to feed these *network weights* to an interaction function. To this end, we employ a more general formulation of classification as an attention

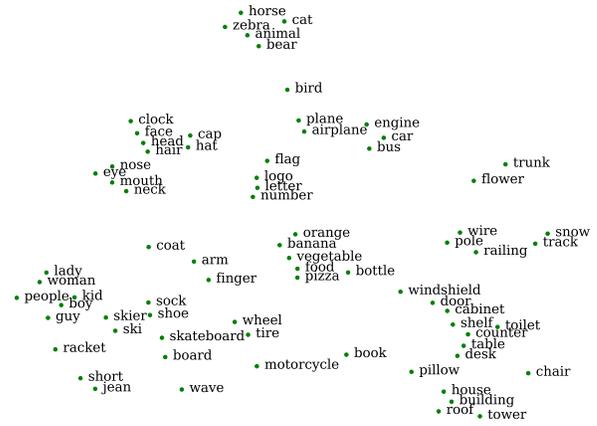


Figure 2: t-SNE visualization of the object classification weights that have been trained on top of contextualized image-based representations. Entities that appear similar to each other, or participate in similar relations, have a closer semantic affinity. This enables link prediction similar to Knowledge Graph Embeddings, and leads to generalization.

layer instead. In this layer, the extracted image-based and contextually enriched representations attend to trainable schema embeddings of all classes such that (a) the attention coefficients are the classification scores (we enforce this by applying a classification loss on the attention outputs), and (b) the attention values carry the prior knowledge that is injected into the image-based representations (a top-down mechanism).

Furthermore, instead of training a separate interaction function for the prior model, we exploit the message passing function that we already have available in the perception model; after fusing the schemata and the image-based object-representations, we contextualize and classify the representations again. Other than computational efficiency, this has the advantage that the image-based object representations and the schemata are combined in the embedding space rather than the probability space.

We train the schemata using the Visual Genome (Krishna et al. 2017) dataset. We show that our model can accurately generate the captured commonsense knowledge and that iterative injection of this knowledge, as a top-down signal, leads to significantly higher classification accuracy. Additionally, we draw from the recent advancements in self-supervised learning and show that the schemata can be trained with only a fraction of labeled images. This allows us to fine-tune the perception model without any additional images; instead, we can use a knowledge base of hand-crafted or external triples and train with their *mental images* (schemata). As a result, compared to the self-supervised baseline, and with 1% of the training data, our model achieves more than 3% improvement in object classification, 26% in scene graph classification, and 36% in predicate prediction; an accuracy that is almost equal to when using 100% of the labeled images.

Related Works

While the concept of schemata can be applied to any form of perceptual processing, and there are recent deep learning

²We leave out the discussion on the schemata in Kant’s view.

models of action schemata (Kansky et al. 2017), we focus on the figurative schemata in the visual scene understanding domain. There has also been a body of related research on relational reasoning outside the scene graph domain (Wu, Lenz, and Saxena 2014; Deng et al. 2014; Hu et al. 2016, 2017; Santoro et al. 2017; Sabour, Frosst, and Hinton 2017). Nevertheless, research in this field was largely accelerated after the release of Visual Relation Detection (VRD) (Lu et al. 2016) and the Visual Genome (Krishna et al. 2017) datasets. Baier, Ma, and Tresp (2017, 2018) proposed the first KG-based model of prior knowledge that improves SG classification. VTransE (Zhang et al. 2017) proposed to capture relations by applying the KGE model of TransE (Bordes et al. 2013) on the visual embeddings. Yu et al. (2017) employed a teacher-student model to distill external language knowledge. Iterative Message Passing (Xu et al. 2017), Neural Motifs (Zellers et al. 2018) (NM), and Graph R-CNN (Yang et al. 2018) used RNNs and graph convolutions to propagate image context. Tang et al. (2019) exploited dynamic tree structures and Chen et al. (2019a) proposed a method based on multi-agent policy gradients. Sharifzadeh et al. (2019) employed the predicted pseudo depth maps in addition to the 2D information. In general, scene graph classification methods are closely related to KGE models (Nickel, Tresp, and Krieger 2011; Nickel et al. 2016). For an extensive discussion on the connection between perception, KGEs, and cognition, refer to (Tresp, Sharifzadeh, and Konopatzki 2019; Tresp et al. 2020). The link prediction in KGEs arises from the compositionality of the trained embeddings. Some other forms of compositionality in neural networks are discussed in (Montufar et al. 2014). In this work, we introduce assimilation, which strengthens the representations within the neural network’s causal structure, addressing an issue raised by Fodor, Pylyshyn et al. (1988). Some of the issues that we address in this work have also been recently discussed by Bengio (2017); Goyal et al. (2019); Mittal et al. (2020).

Methods

In summary, after an initial classification step, we combine the image-based representations with the schemata of predicted classes. We then collect supportive evidence from the neighbors before re-classifying each entity (Ref. Figure 1). In what follows, bold lower case letters denote vectors, bold upper case letters denote matrices, and the letters denote scalar quantities or random variables. Subscripts and superscripts denote variables and calligraphic upper case letters for sets.

Definitions

Let us consider a given image \mathbf{I} and a set of n bounding boxes $\mathcal{B} = \{\mathbf{b}_i\}_{i=1}^n$, $\mathbf{b}_i = [b_i^x, b_i^y, b_i^w, b_i^h]$, such that $[b_i^x, b_i^y]$ are the coordinates of \mathbf{b}_i and $[b_i^w, b_i^h]$ are its width and height. We build a **Scene Representation Graph**, $SRG = \{\mathcal{V}, \mathcal{E}\}$ as a structured presentation of the objects and predicates in \mathbf{I} . $\mathcal{X}^o = \{\mathbf{x}_i^o\}_{i=1}^n$, $\mathbf{x}_i^o \in \mathbb{R}^d$ denote the features of object nodes and $\mathcal{X}^p = \{\mathbf{x}_i^p\}_{i=1}^m$, $\mathbf{x}_i^p \in \mathbb{R}^d$ denote the features of predicate nodes³. Each \mathbf{x}_i^o is initialized by a pooled

³Similar to (Yang et al. 2018; Koncel-Kedziorski et al. 2019), we consider each object node as direct neighbors with its predicate

image-based object representation, extracted by applying VGG16 (Simonyan and Zisserman 2014) or ResNet-50 (He et al. 2016) on the image contents of \mathbf{b}_i . Each \mathbf{x}_i^p is initialized by applying a two layered fully connected network on the relational position vector \mathbf{t} between a head i and a tail j where $\mathbf{t} = [t_x, t_y, t_w, t_h]$, $t_x = (b_i^x - b_j^x)/b_i^w$, $t_y = (b_i^y - b_j^y)/b_j^h$, $t_w = \log(b_i^w/b_j^w)$, $t_h = \log(b_i^h/b_j^h)$. The implementation details of the networks are provided in the Supplementary.

Scene graph classification is the mapping of each node in scene representation graph to a label where each object node is from the label set \mathcal{C}^o and each predicate node from \mathcal{C}^p . The resulting labeled graph is a set of triples referred to as the **Scene Graph**. We also define a **Probabilistic Knowledge Graph (PKG)** as a graph where the weight of a triple is the expected value of observing that relation given the head and tail classes and regardless of any given images⁴. Later we will show that our model can accurately generate the PKG, i.e., the commonsense that is captured from perceptions during training.

In what comes next, \mathbf{x}_i^o and \mathbf{x}_i^p are treated identically except for classification with respect to \mathcal{C}^o or \mathcal{C}^p . Therefore, for a better readability, we only write \mathbf{x}_i .

Contextualized Scene Representation Graph

We obtain contextualized object representations \mathbf{z}_i by applying a graph convolutional neural network, on SRG . We also refer to this module as our *interaction function*. We use a Graph Transformer as a variant of the Graph Network Block (Battaglia et al. 2018; Koncel-Kedziorski et al. 2019) with multi-headed attentions as

$$\mathbf{m}_i^{\mathcal{N}(i)} = \frac{1}{K} \sum_{k=1}^K \sum_{j \in \mathcal{N}(i)} \alpha_{ij}^{(l,k)} \mathbf{W}^{(l,k)} \mathbf{z}_j^{(l,t)} \quad (1)$$

$$\mathbf{z}'_i^{(l)} = LN(\mathbf{z}_i^{(l,t)} + \mathbf{m}_i^{\mathcal{N}_{in}(i)} + \mathbf{m}_i^{\mathcal{N}_{out}(i)}) \quad (2)$$

$$\mathbf{z}_i^{(l+1,t)} = LN(\mathbf{z}'_i^{(l)} + f(\mathbf{z}'_i^{(l)})), \quad (3)$$

where $\mathbf{z}_i^{(l,t)}$ is the embedding of node i in the l -th graph convolution layer and t -th assimilation. In the first layer $\mathbf{z}_i^{(0,t)} = \mathbf{x}_i$. LN is the layer norm (Ba, Kiros, and Hinton 2016), K is the number of attentional heads and $\mathbf{W}^{(l,k)}$ is the weight matrix of the k -th head in layer l . $\mathcal{N}(i)$ represent the set of neighbors, which are either incoming $\mathcal{N}_{in}(i)$ or outgoing $\mathcal{N}_{out}(i)$. $f(\cdot)$ is a two layered feed-forward neural network with Leaky ReLU non-linearities between each layer. $\alpha_{ij}^{(l,k)}$ denotes the attention coefficients in each head and is defined as

$$e_{ij}^{(l,k)} = \sigma(\mathbf{h}^{(l,k)} \cdot [\mathbf{z}_i^{(l)} \parallel \mathbf{W}^{(l,k)} \mathbf{z}_j^{(l)}]) \quad (4)$$

$$\alpha_{ij}^{(l,k)} = \frac{\exp(e_{ij}^{(l,k)})}{\sum_{q \in \mathcal{N}(i)} \exp(e_{iq}^{(l,k)})} \quad (5)$$

nodes and each predicate node as direct neighbors with its head and tail object nodes.

⁴Note that while typical knowledge graphs such as Freebase are based on object instances, given the nature of our image dataset, we focus on classes.

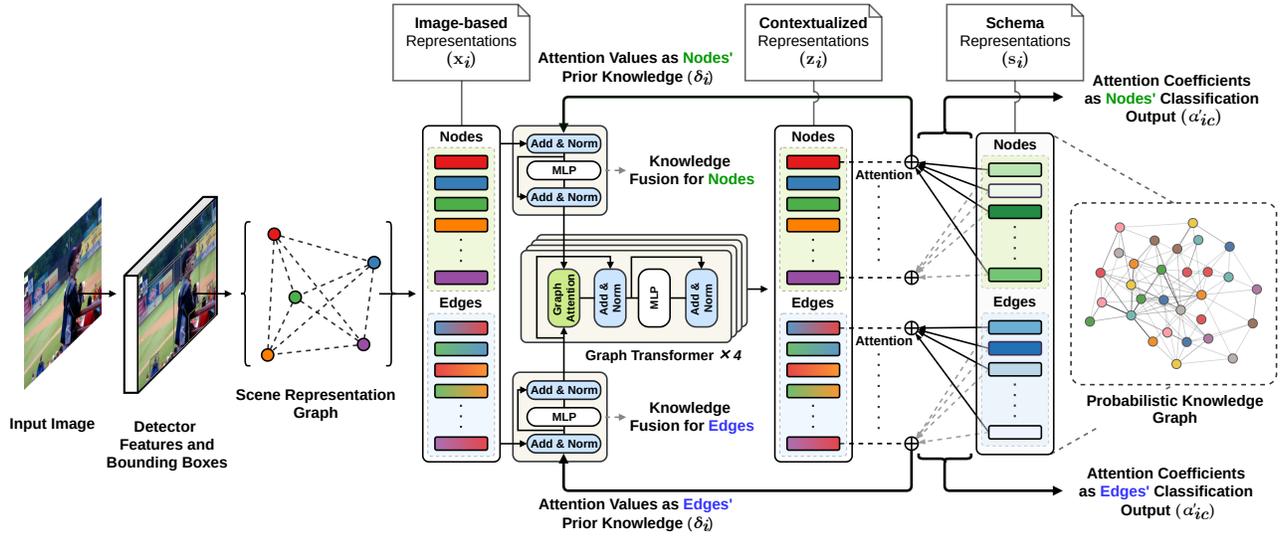


Figure 3: We formulate the classification as attention layer between object and schema representations. Contextualizing image-based object representations before classification encourages the schemata to learn *image-based relational* prior knowledge. As a result, the attention values that are injected from the schemata to scene representations and then propagated. In this way, they enrich the image-based representations with prior knowledge. Additionally, the interactions between schemata can reconstruct the probabilistic knowledge graph (right).

with $\mathbf{h}^{(l,k)}$ as a learnable weight vector and \parallel denoting concatenation. σ is the Leaky ReLU with the slope of 0.2.

Schemata

We define the schema of a class c as an embedding vector \mathbf{s}_c . We realize object and predicate classification by an attention layer between the contextualized representations and the schemata such that the classification outputs α'_{ic} are computed as the attention coefficients between \mathbf{z}_i and \mathbf{s}_c as

$$\alpha'_{ic} = \text{softmax}(a(\mathbf{z}_i^{(L,t)}, \mathbf{s}_c)) \quad (6)$$

where, $a(\cdot)$ is the attention function that we implement as the dot-product between the input vectors, and $\mathbf{z}_i^{(L,t)}$ is the output from the last (L -th) layer of the Graph Transformer. The attention values δ_i capture the schemata messages as

$$\delta_i = \sum_{c \in \mathcal{C}} \alpha'_{ic} \mathbf{s}_c \quad (7)$$

and we inject them back to update the scene representations as

$$\mathbf{u}_i = LN(\mathbf{x}_i + \delta_i) \quad (8)$$

$$\mathbf{z}_i^{(0,t+1)} = LN(\mathbf{u}_i + g(\mathbf{u}_i)) \quad (9)$$

where $g(\cdot)$ is a two-layered feed-forward network with Leaky ReLU non-linearities. Note that we compute \mathbf{u}_i by fusing the attention values with the *original image features* \mathbf{x}_i . Therefore, the outputs from previous Graph Transformer layers will not be accumulated, and the original image-based features will not vanish.

We define *assimilation* as the set of computations from $\mathbf{z}_i^{(L,t)}$ to $\mathbf{z}_i^{(L,t+1)}$. This includes the initial classification step

(Eq. 6), fusion of schemata with image-based vectors (Eq. 9) and the application of the interaction function on the updated embeddings (Eq. 3). We expect to get refined object representations after the assimilation. Therefore, we assimilate several times such that after each update of the classification results, the priors are also updated accordingly. During training, and for each step of assimilation, we employ a supervised attention loss, i.e. categorical cross entropy, between the one-hot encoded ground truth labels and α'_{ic} . This indicates a multi-task learning strategy where one task (for the first assimilation) is to optimize for $P(y_q | x_1, \dots, x_\theta)$, with x_q as a random variable representing the image-based features of q , y_q as the label, and $\theta = m + n$. The other set of tasks is to optimize for $P(y_q^{t+1} | x_1^t, \dots, x_\theta^t, y_1^t, \dots, y_\theta^t)$. We refer to the first task as **IC**, for **Image-based Classification** and to the second set of tasks as **ICP** for **Image-based Classification with Prior knowledge**. We train the second task using teacher forcing and by setting the labels to their ground truth values. Therefore, in order to prevent collapse, we set the edge schemata to zero. This resembles link prediction, such that we denote an incomplete input graph. Note that even when no images are available, we can still train for the *ICP* from a collection of external or hand-crafted triples by directly assigning $\mathbf{z}_i^{(0,t+1)} = \delta_i$ such that $\alpha'_{ic} = \text{onehot}(c_i)$.

GCN vs. Prior Model: A matter of inductive biases

Typical GCNs, such as the Graph Transformer, take the features derived from each bounding box as input, apply non-linear transformations and propagate them to the neighbors in the following layers. Each GCN layer consists of fully connected neural networks. Therefore, *theoretically* they can also model and propagate prior knowledge that is *not* visible

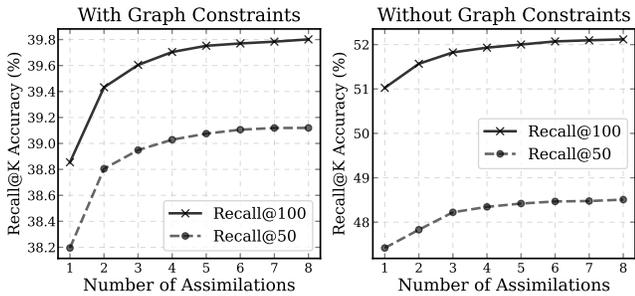


Figure 4: The results of our ablation studies. We study the effect of each assimilation in scene graph classification. Note that the model has been trained for only 4 assimilations yet it can generalize.

in bounding boxes. However, experimental results of previous works (and also this work) confirm that explicit modeling and propagation of prior knowledge (*ICP*) can still improve the classification accuracy. Why is that the case?

Let us consider the following. According to the universal approximation theorem (Csáji et al. 2001), when we solve for *IC* as $P(y_k|x_1, \dots, x_o)$, our model might learn to capture a desired form of $P(y_k|x_1, \dots, x_o, y_1, \dots, y_o)$. However, in practice, the learning algorithm does not always find the best function. Therefore, we require appropriate inductive biases to guide us through the learning process. As Caruana (1997) puts: “*Multitask Learning is an approach to inductive transfer that improves generalization by using the domain information contained in the training signals of related tasks as an inductive bias. It does this by learning tasks in parallel while using a shared representation; what is learned for each task can help other tasks be learned better*”. For example, in the encoder-decoder models for machine translation, e.g. Transformers (Vaswani et al. 2017), the prediction is often explicitly conditioned not just on the encoded inputs but also on the decoded outputs from the previous tokens. Therefore, the decoding in each step can be interpreted as computing $P(y_k|x_1, \dots, x_o, y_1, \dots, y_{k-1})$. Note that the previous predictions such as y_1 , cannot benefit from the future predictions $\{y_2, \dots, y_o\}$. However, in our model, we provide an explicit bias towards utilizing predictions in *all indices*. In fact, our model can be interpreted as an encoder-decoder network, where the decoder consists of multiple decoders. Therefore, the decoding depends not just on the encoded image features but also on the previously decoded outputs. In other words, by injecting schema embeddings, as embeddings that are trained over *all images*, we impose the bias to propagate *what is not visible in the bounding box*. As will be shown later, we can train for *ICP* and *IC* even with smaller splits of annotated images, which can lead to competitive results with fewer labels. Additionally, assimilation enables us to quantify the propagated prior knowledge. This interpretability is another advantage that GCNs alone do not have.

Settings We train our models on the common split of Visual Genome (Krishna et al. 2017) dataset containing images labeled with their scene graphs (Xu et al. 2017). This split

Method		SGCls		PredCls		Mean
		@50	@100	@50	@100	
Unconstrained	IMP+ (Xu et al. 2017)	12.1	16.9	20.3	28.9	19.5
	FREQ (Zellers et al. 2018)	13.5	19.6	24.8	37.3	23.8
	SMN (Zellers et al. 2018)	15.4	20.6	27.5	37.9	25.3
	KERN(Chen et al. 2019c)	19.8	26.2	36.3	49.0	32.8
	Schemata	21.4	28.8	40.1	54.9	36.3
Constrained	IMP (Xu et al. 2017)	3.1	3.8	6.1	8.0	5.2
	IMP+ (Xu et al. 2017)	5.8	6.0	9.8	10.5	8.0
	FREQ (Zellers et al. 2018)	6.8	7.8	13.3	15.8	10.9
	SMN (Zellers et al. 2018)	7.1	7.6	13.3	14.4	10.6
	KERN(Chen et al. 2019c)	9.4	10.0	17.7	19.2	14.0
	VCTree(Tang et al. 2019)	10.1	10.8	17.9	19.4	14.5
	Schemata	10.1	10.9	19.1	20.7	15.2
Schemata - PKG	-	-	8.2	9.4	-	

Table 1: Comparison of the mR@50 and mR@100, with and without graph constraints for SGCls and PredCls.

takes the most frequent 150 object and 50 predicate classes in total, with an average of 11.5 objects and 6.2 predicates in each image. We report the experimental results on the test set, under two standard classification settings of predicate classification (**PredCls**): predicting predicate labels given a ground truth set of object boxes and object labels, and scene graph classification (**SGCls**): predicting object and predicate labels, given the set of object boxes. Another popular setting is the scene graph detection (SGDet), where the network should also detect the bounding boxes. Since the focus of our study is not on improving the object detector backbone and our improvements in SGDet were similar to the improvements in SGCls, we do not report them here. For those results, please refer to our official code repository. We report all the results under *constrained* and *unconstrained* setups (Yu et al. 2017). In the unconstrained setup, we allow for multiple predicate labels, whereas in the constrained setup, we only take the top-1 predicted predicate label.

Metrics We use Recall@K (**R@K**) as the standard metric. R@K computes the mean prediction accuracy in each image given the top *K* predictions. In VG, the distribution of labeled relations is highly imbalanced. Therefore, we additionally report Macro Recall (Sharifzadeh et al. 2019; Chen et al. 2019c) (**mR@K**) to reflect the improvements in the long tail of the distribution. In this setting, the overall recall is computed by taking the mean over recall per predicate.

Experiments The goal of our experiments is (A) to study whether injecting prior knowledge into scene representations can improve the classification and (B) to study the common-sense knowledge that is captured in our model. In what follows, *backbone* refers to VGG16/ResNet-50 that generates the *SRG*, and *main model* refers to part of the network that applies contextualization and assimilation. The backbone can be trained from a set of labeled images (in a supervised manner), unlabeled images (in a self-supervised manner), or a combination of the two. The main model can be trained from

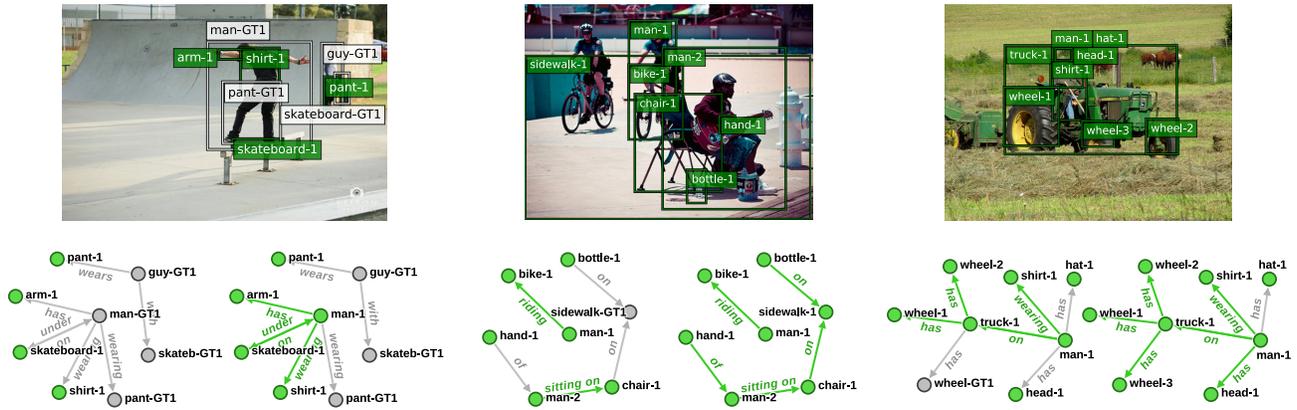


Figure 5: Qualitative examples of improved scene graph classification results (Recall@50) through assimilations of our model. From left to right is after each assimilation. Green and gray colors indicate true positives and false negatives concluded by the model. For example consider the middle image, where the sidewalk was initially misclassified as a street. After seeing a biker in the image and a man sitting on a chair, a reasonable inference is that this should be a sidewalk.

Training	SGCls R@100			PredCls R@100			Object Classification		
	1%	10%	100%	1%	10%	100%	1%	10%	100%
Sup - IC	1.84 ± 0.26	13.90 ± 0.97	33.6	40.61 ± 0.84	52.51 ± 1.19	62.0	14.38 ± 0.57	38.45 ± 1.21	64.2
Self-Sup - IC	12.12 ± 0.47	26.14 ± 0.77	36.8	48.10 ± 0.54	58.14 ± 0.35	63.4	40.75 ± 0.48	56.97 ± 0.76	68.0
Self-Sup - IC & ICP	15.36 ± 0.38	27.37 ± 0.47	37.1	65.68 ± 0.12	65.42 ± 0.19	65.7	42.09 ± 0.65	58.60 ± 0.56	68.4

Table 2: Comparison of R@100 for SGCls, PredCls and Object Classification tasks on smaller splits of the VG dataset.

a set of labeled images (the *IC* task), a prior knowledge base (*ICP*) or a combination of the two. For (A), we conduct the following studies:

1. We train both the backbone and the main model from all the labeled images and for both tasks. We use the VGG-16 backbone as trained by Zellers et al. (2018). This allows us to compare the results with the related works directly. We evaluate the classification accuracy for 8 assimilations (until the changes are not significant anymore). Table 1 compares the performance of our model to the state-of-the-art under mR@K (for the R@K results refer to the supplementary). As shown, our model exceeds the others on average and under most settings. supplementary. Figure 4 shows our ablation study, indicating that the accuracy is improved after each assimilation.
2. To qualitatively examine these results, we present some of the images and their scene graphs after two assimilations, in Figure 5. For example in the right image, while the wheel is almost fully occluded, we can still classify it once we classify other objects and employ commonsense (e.g., trucks have wheels). Another interesting example is the middle image, where the sidewalk is initially misclassified as a street. After seeing a biker in the image and a man sitting on a chair, a reasonable inference is that this should be a sidewalk! Similarly, in the left image, the man is facing away from the camera, and his pose makes it hard to classify him unless we utilize our prior knowledge about the arm, pants, shirt, and skateboard.
3. Figure 7 shows the improvements per each predicate

class. The results indicate that most improvements occur in under-represented classes. This means that we have achieved a generalization performance that is beyond the simple reflection of the dataset’s statistical bias.

4. To understand the importance of prior knowledge compared to having a large set of labeled images, we conduct the following study: we uniformly sample two splits with 1% and 10% of VG. The images in each split are considered as *labeled*. We ignore the labels of the remaining images and consider them as unlabeled⁵. Instead, we treat the set of ignored labels as a form of external/hand-crafted knowledge in the form of triples. For each split, we train the full model (I) with a backbone that has been trained in a supervised fashion with the respective split and no pre-training, and the main model that has been trained for *IC* (without commonsense) with the respective split, (II) with a backbone that has been pre-trained on ImageNet (Deng et al. 2009) and fine-tuned on the Visual Genome (in a self-supervised fashion with BYOL (Grill et al. 2020)) and fine-tuned on the respective split of the visual genome (in a supervised fashion) and the main model that has been trained for *IC* with the respective split, and (III) Similar to 2, except that we include the *ICP* and train the main model by assimilating the entire prior knowledge base including

⁵Note that these splits are different from the recently proposed few-shot learning set by Chen et al. (2019d). In (Chen et al. 2019d), the goal is to study the few-shot learning of *predicates* only. However, we explore a more competitive setting, where only a fraction of both *objects* and *predicates* are labeled.

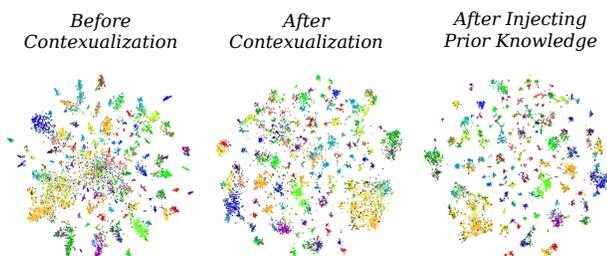


Figure 6: t-SNE visualization of object representations.

the external triples. We discard their image-based features (x_i) for the triples outside a split. Also, to treat all triples equally when injecting the prior knowledge, we discard *all* image-based features and directly feed the δ_i to the graph transformer. To prevent collapse, we randomly drop some of the δ_i s. Since BYOL is based on ResNet-50, for a fair comparison, we train all models in this experiment with ResNet-50 (including another model that we train with 100% of the data). In the Scene Graph Classification community, the results are often reported under an arbitrary random seed, and previous works have not reported the summary statistics over several runs before. To allow for a fair comparison of our model to those works (on the 100% set), we followed the same procedure in the study **A1**. However, to encourage a statistically more stable comparison of future models in this experiment, we report the summary statistics (arithmetic mean and standard deviation) over five random fractions (1% and 10%) of VG training set⁶. As shown in Table 2, utilizing prior knowledge allows to achieve almost the same predicate prediction accuracy with 1% of the data only. Also, we largely improve object classification and scene graph classification.

Evaluation

For **B** we consider the following studies:

1. We visualize the semantic affinity of *schema representations* by employing t-SNE (Maaten and Hinton 2008). As we can see in Figure 2, the schema representations of entities that are *visually* or *relationally* similar are the closest to each other.
2. We inspect the semantic affinity of *object representations* by employing t-SNE (I) before contextualization, (II) after contextualization and (III) after injecting prior knowledge. The results are represented in Figure 6. Each color represents a different object class. This investigation confirms that object representations will get into more separable clusters after injecting prior knowledge.
3. Finally, we evaluate our model’s accuracy in link prediction. The goal is to quantitatively evaluate our model’s understanding of *relational* commonsense, i.e., relational structure of the probabilistic knowledge graph. Similar to a KGE link prediction, we predict the predicate given *head* and *tail* of a relation. In other words, we feed our model

⁶The splits are available at: <https://github.com/sharifza/schemata>

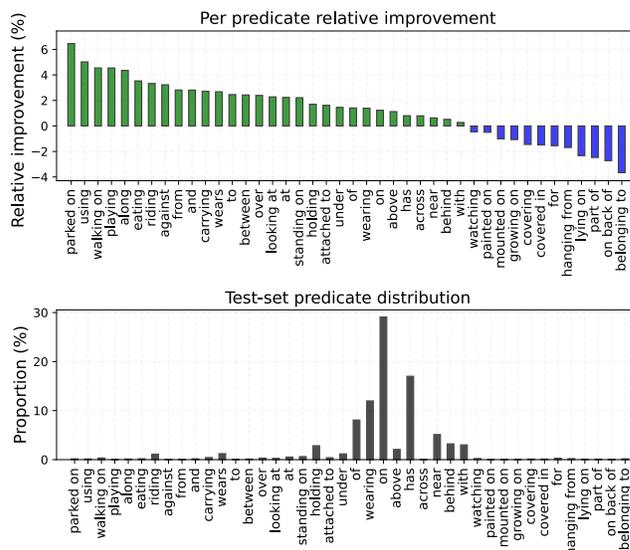


Figure 7: The top shows the per-predicate classification accuracy improvement after injecting prior knowledge, in SGCLs R@100. The bottom shows the distribution of sample proportion for the predicates in the VG.

with the schema of head and tail, together with a *zero-vector* for the image-based representations. As we can see in Table 1, in *Schemata - PKG*, even if we do not provide any image-based information, our model can still *guess* the expected predicates similar to a KGE model. While this guess is not as accurate as when we present it with an image, the accuracy is still remarkable.

Conclusion

We discussed schemata as mental representations that enable compositionality and reasoning. To model schemata in a deep learning framework, we introduced them as representations that encode image-based and relational prior knowledge of objects and predicates in each class. By defining classification as an attention layer instead of a fully connected layer, we introduced an inductive bias that enabled the propagation of prior knowledge. Our experiments on the Visual Genome dataset confirmed the effectiveness of assimilation through qualitative and quantitative measures. Our model achieved higher accuracy under most settings and could also accurately predict the commonsense knowledge. Additionally, we showed that our model could be fine-tuned from external sources of knowledge in the form of triples. When combined with pre-trained schemata in a self-supervised setting, this leads to a predicate prediction accuracy that is almost equal to the full model. Also, it gives significant improvements in the scene graph and object classification tasks. We hope that this work will open new research directions in utilizing commonsense to learn from little annotations.

Acknowledgments

We would like to thank Max Berrendorf, Dario Konopatzki, Shaya Akbarinejad, Lisa Machata, Shabnam Sadegh, and the anonymous reviewers for the fruitful discussions and helpful feedback on the manuscript. This work has been partially funded by the German Federal Ministry of Education and Research (BMBF) under Grant No. 01IS18036A.

References

- Arbib, M. A. 1992. Schema theory. *The encyclopedia of artificial intelligence 2*: 1427–1443.
- Ba, J. L.; Kiros, J. R.; and Hinton, G. E. 2016. Layer normalization. *arXiv preprint arXiv:1607.06450*.
- Baier, S.; Ma, Y.; and Tresp, V. 2017. Improving visual relationship detection using semantic modeling of scene descriptions. In *International Semantic Web Conference*, 53–68. Springer.
- Baier, S.; Ma, Y.; and Tresp, V. 2018. Improving information extraction from images with learned semantic models. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, 5214–5218. AAAI Press.
- Battaglia, P. W.; Hamrick, J. B.; Bapst, V.; Sanchez-Gonzalez, A.; Zambaldi, V.; Malinowski, M.; Tacchetti, A.; Raposo, D.; Santoro, A.; Faulkner, R.; et al. 2018. Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*.
- Bengio, Y. 2017. The consciousness prior. *arXiv preprint arXiv:1709.08568*.
- Bordes, A.; Usunier, N.; Garcia-Duran, A.; Weston, J.; and Yakhnenko, O. 2013. Translating embeddings for modeling multi-relational data. In *Advances in neural information processing systems*, 2787–2795.
- Caruana, R. 1997. Multitask learning. *Machine learning* 28(1): 41–75.
- Chen, L.; Zhang, H.; Xiao, J.; He, X.; Pu, S.; and Chang, S.-F. 2019a. Counterfactual critic multi-agent training for scene graph generation. In *Proceedings of the IEEE International Conference on Computer Vision*, 4613–4623.
- Chen, T.; Xu, M.; Hui, X.; Wu, H.; and Lin, L. 2019b. Learning semantic-specific graph representation for multi-label image recognition. In *Proceedings of the IEEE International Conference on Computer Vision*, 522–531.
- Chen, T.; Yu, W.; Chen, R.; and Lin, L. 2019c. Knowledge-embedded routing network for scene graph generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 6163–6171.
- Chen, V. S.; Varma, P.; Krishna, R.; Bernstein, M.; Re, C.; and Fei-Fei, L. 2019d. Scene graph prediction with limited labels. In *Proceedings of the IEEE International Conference on Computer Vision*, 2580–2590.
- Csáji, B. C.; et al. 2001. Approximation with artificial neural networks. *Faculty of Sciences, Eötvös Loránd University, Hungary* 24(48): 7.
- Deng, J.; Ding, N.; Jia, Y.; Frome, A.; Murphy, K.; Bengio, S.; Li, Y.; Neven, H.; and Adam, H. 2014. Large-scale object classification using label relation graphs. In *European conference on computer vision*, 48–64. Springer.
- Deng, J.; Dong, W.; Socher, R.; Li, L.-J.; Li, K.; and Fei-Fei, L. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, 248–255. Ieee.
- Fodor, J. A.; Pylyshyn, Z. W.; et al. 1988. Connectionism and cognitive architecture: A critical analysis. *Cognition* 28(1-2): 3–71.
- Goyal, A.; Lamb, A.; Hoffmann, J.; Sodhani, S.; Levine, S.; Bengio, Y.; and Schölkopf, B. 2019. Recurrent independent mechanisms. *arXiv preprint arXiv:1909.10893*.
- Grill, J.-B.; Strub, F.; Althé, F.; Tallec, C.; Richemond, P. H.; Buchatskaya, E.; Doersch, C.; Pires, B. A.; Guo, Z. D.; Azar, M. G.; et al. 2020. Bootstrap your own latent: A new approach to self-supervised learning. *arXiv preprint arXiv:2006.07733*.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 770–778.
- Hochreiter, S.; and Schmidhuber, J. 1997. Long short-term memory. *Neural computation* 9(8): 1735–1780.
- Hou, J.; Wu, X.; Qi, Y.; Zhao, W.; Luo, J.; and Jia, Y. 2019. Relational Reasoning using Prior Knowledge for Visual Captioning. *arXiv preprint arXiv:1906.01290*.
- Hu, H.; Deng, Z.; Zhou, G.-T.; Sha, F.; and Mori, G. 2017. Labelbank: Revisiting global perspectives for semantic segmentation. *arXiv preprint arXiv:1703.09891*.
- Hu, H.; Zhou, G.-T.; Deng, Z.; Liao, Z.; and Mori, G. 2016. Learning structured inference neural networks with label relations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2960–2968.
- Kansky, K.; Silver, T.; Mély, D. A.; Eldawy, M.; Lázaro-Gredilla, M.; Lou, X.; Dorfman, N.; Sidor, S.; Phoenix, S.; and George, D. 2017. Schema networks: Zero-shot transfer with a generative causal model of intuitive physics. *arXiv preprint arXiv:1706.04317*.
- Kant, I. 1787. *Kritik der reinen Vernunft:[Hauptband]*. Walter de Gruyter.
- Kipf, T. N.; and Welling, M. 2016. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*.
- Koncel-Kedziorski, R.; Bekal, D.; Luan, Y.; Lapata, M.; and Hajishirzi, H. 2019. Text generation from knowledge graphs with graph transformers. *arXiv preprint arXiv:1904.02342*.
- Krishna, R.; Zhu, Y.; Groth, O.; Johnson, J.; Hata, K.; Kravitz, J.; Chen, S.; Kalantidis, Y.; Li, L.-J.; Shamma, D. A.; et al. 2017. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International Journal of Computer Vision* 123(1): 32–73.

- Lu, C.; Krishna, R.; Bernstein, M.; and Fei-Fei, L. 2016. Visual relationship detection with language priors. In *European Conference on Computer Vision*, 852–869. Springer.
- Maaten, L. v. d.; and Hinton, G. 2008. Visualizing data using t-SNE. *Journal of machine learning research* 9(Nov): 2579–2605.
- Mittal, S.; Lamb, A.; Goyal, A.; Voleti, V.; Shanahan, M.; Lajoie, G.; Mozer, M.; and Bengio, Y. 2020. Learning to combine top-down and bottom-up signals in recurrent neural networks with attention over modules. In *International Conference on Machine Learning*, 6972–6986. PMLR.
- Montufar, G. F.; Pascanu, R.; Cho, K.; and Bengio, Y. 2014. On the number of linear regions of deep neural networks. In *Advances in neural information processing systems*, 2924–2932.
- Nickel, M.; Murphy, K.; Tresp, V.; and Gabrilovich, E. 2016. A review of relational machine learning for knowledge graphs. *Proceedings of the IEEE* 104(1): 11–33.
- Nickel, M.; Tresp, V.; and Kriegel, H.-P. 2011. A three-way model for collective learning on multi-relational data. In *Icml*, volume 11, 809–816.
- Piaget, J. 1923. *Langage et pensée chez l'enfant*. Delachaux et Niestlé.
- Sabour, S.; Frosst, N.; and Hinton, G. E. 2017. Dynamic routing between capsules. *arXiv preprint arXiv:1710.09829* .
- Santoro, A.; Raposo, D.; Barrett, D. G.; Malinowski, M.; Pascanu, R.; Battaglia, P.; and Lillicrap, T. 2017. A simple neural network module for relational reasoning. In *Advances in neural information processing systems*, 4967–4976.
- Sharifzadeh, S.; Moayed Baharlou, S.; Berrendorf, M.; Koner, R.; and Tresp, V. 2019. Improving Visual Relation Detection using Depth Maps. *arXiv preprint arXiv:1905.00966* .
- Simonyan, K.; and Zisserman, A. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* .
- Tang, K.; Zhang, H.; Wu, B.; Luo, W.; and Liu, W. 2019. Learning to compose dynamic tree structures for visual contexts. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 6619–6628.
- Tresp, V.; Sharifzadeh, S.; and Konopatzki, D. 2019. A Model for Perception and Memory. *Conference on Cognitive Computational Neuroscience* .
- Tresp, V.; Sharifzadeh, S.; Konopatzki, D.; and Ma, Y. 2020. The Tensor Brain: Semantic Decoding for Perception and Memory. *arXiv preprint arXiv:2001.11027* .
- Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, Ł.; and Polosukhin, I. 2017. Attention is all you need. In *Advances in neural information processing systems*, 5998–6008.
- Wu, C.; Lenz, I.; and Saxena, A. 2014. Hierarchical Semantic Labeling for Task-Relevant RGB-D Perception. In *Robotics: Science and systems*.
- Xu, D.; Zhu, Y.; Choy, C. B.; and Fei-Fei, L. 2017. Scene graph generation by iterative message passing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 5410–5419.
- Yang, J.; Lu, J.; Lee, S.; Batra, D.; and Parikh, D. 2018. Graph r-cnn for scene graph generation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 670–685.
- Yu, R.; Li, A.; Morariu, V. I.; and Davis, L. S. 2017. Visual relationship detection with internal and external linguistic knowledge distillation. In *IEEE International Conference on Computer Vision (ICCV)*.
- Zellers, R.; Yatskar, M.; Thomson, S.; and Choi, Y. 2018. Neural motifs: Scene graph parsing with global context. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 5831–5840.
- Zhang, H.; Kyaw, Z.; Chang, S.; and Chua, T. 2017. Visual Translation Embedding Network for Visual Relation Detection. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, 3107–3115. IEEE Computer Society. ISBN 978-1-5386-0457-1. doi:10.1109/CVPR.2017.331. URL <https://doi.org/10.1109/CVPR.2017.331>.