

# “The red one!”: On learning to refer to things based on discriminative properties

Angeliki Lazaridou and Nghia The Pham and Marco Baroni

University of Trento

{angeliki.lazaridou|thepham.nghia|marco.baroni}@unitn.it

## Abstract

As a first step towards agents learning to communicate about their visual environment, we propose a system that, given visual representations of a referent (CAT) and a context (SOFA), identifies their *discriminative attributes*, i.e., properties that distinguish them (*has\_tail*). Moreover, although supervision is only provided in terms of discriminativeness of attributes for pairs, the model learns to assign plausible attributes to specific objects (SOFA-*has\_cushion*). Finally, we present a preliminary experiment confirming the referential success of the predicted discriminative attributes.

## 1 Introduction

There has recently been renewed interest in developing systems capable of genuine language understanding (Hermann et al., 2015; Hill et al., 2015). In this perspective, it is important to think of an appropriate general framework for teaching language to machines. Since we use language primarily for communication, a reasonable approach is to develop systems within a genuine communicative setup (Steels, 2003; Mikolov et al., 2015). Our long-term goal is thus to develop communities of computational agents that learn how to use language efficiently in order to achieve communicative success (Vogel et al., 2013; Foerster et al., 2016).

Within this general picture, one fundamental aspect of meaning where communication is indeed crucial is the act of *reference* (Searle, 1969; Abbott, 2010), the ability to successfully talk to others about things in the external world. A specific instantiation of reference studied in this paper is that of referring expression generation (Dale and

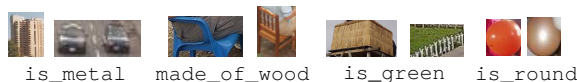


Figure 1: Discriminative attributes predicted by our model. Can you identify the intended referent? See Section 6 for more information

Reiter, 1995; Mitchell et al., 2010; Kazemzadeh et al., 2014). A *necessary* condition for achieving successful reference is that referring expressions (REs) accurately distinguish the intended referent from any other object in the context (Dale and Haddock, 1991). Along these lines, we present here a model that, given an intended referent and a context object, predicts the attributes that discriminate between the two. Some examples of the behaviour of the model are presented in Figure 1.

Importantly, and distinguishing our work from earlier literature on generating REs (Krahmer and Van Deemter, 2012): (i) the input objects are represented by natural images, so that the agent must learn to extract relevant attributes from realistic data; and (ii) no direct supervision on the attributes of a single object is provided: the training signal concerns their discriminativeness for object pairs (that is, during learning, the agent might be told that *has\_tail* is discriminative for  $\langle \text{CAT}, \text{SOFA} \rangle$ , but not that it is an attribute of cats). We use this “pragmatic” signal since it could later be replaced by a measure of success in actual communication between two agents (e.g., whether a second agent was able to pick the correct referent given a RE).

## 2 Discriminative Attribute Dataset

We generated the Discriminative Attribute Dataset, consisting of pairs of (intended) *referents* and *contexts*, with respect to which the referents should be identified by their distinctive attributes.



(referent, context)	visual instances	discriminative attributes
(CAT, SOFA)		has_tail, has_cushion, ...
(CAT, APPLE)		has_legs, is_green, ...

Table 1: Example training data

Our starting point is the Visual Attributes for Concepts Dataset (ViSA) (Silberer et al., 2013), which contains *per-concept* (as opposed to *per-image*) attributes for 500 concrete concepts (CAT, SOFA, MILK) spanning across different categories (MAMMALS, FURNITURE), annotated with 636 general attributes. We disregarded ambiguous concepts (e.g., *bat*), thus reducing our working set of concepts  $C$  to 462 and the number of attributes  $V$  to 573, as we eliminated any attribute that did not occur with concepts in  $C$ . We extracted on average 100 images annotated with each of these concepts from ImageNet (Deng et al., 2009). Finally, each image  $i$  of concept  $c$  was associated to a *visual instance vector*, by feeding the image to the VGG-19 ConvNet (Simonyan and Zisserman, 2014), as implemented in the MatConvNet toolkit (Vedaldi and Lenc, 2015), and extracting the second-to-last fully-connected (fc) layer as its 4096-dimensional visual representation  $\mathbf{v}_{c_i}$ .

We split the target concepts into *training*, *validation* and *test* sets, containing 80%, 10% and 10% of the concepts in each category, respectively. This ensures that (i) the intersection between train and test concepts is empty, thus allowing us to test the generalization of the model across different objects, but (ii) there are instances of all categories in each set, so that it is possible to generalize across training and testing objects. Finally we build all possible combinations of concepts in the training split to form pairs of referents and contexts  $\langle c_r, c_c \rangle$  and obtain their (binary) attribute vectors  $\mathbf{p}_{c_r}$  and  $\mathbf{p}_{c_c}$  from ViSA, resulting in 70K training pairs. From the latter, we derive, for each pair, a concept-level “discriminative-ness” vector by computing the symmetric difference  $\mathbf{d}_{c_r, c_c} = (\mathbf{p}_{c_r} - \mathbf{p}_{c_c}) \cup (\mathbf{p}_{c_c} - \mathbf{p}_{c_r})$ . The latter will contain 1s for discriminative attributes, 0s elsewhere. On average, each pair is associated with 20 discriminative attributes. The final training data are triples of the form  $\langle c_r, c_c, \mathbf{d}_{c_r, c_c} \rangle$  (the model *never* observes the attribute vectors of specific concepts), to be associated with visual in-

stances of the two concepts. Table 1 presents some examples.

Note that ViSA contain concept-level attributes, but images contain specific instances of concepts for which a general attribute might not hold. This introduces a small amount of noise. For example, *is\_green* would in general be a discriminative attribute for apples and cats, but it is not for the second sample in Table 1. Using datasets with *per-image* attribute annotations would solve this issue. However, those currently available only cover specific classes of concepts (e.g., only clothes, or animals, or scenes, etc.). Thus, taken separately, they are not general enough for our purposes, and we cannot merge them, since their concepts live in different attribute spaces.

### 3 Discriminative Attribute Network

The proposed *Discriminative Attribute Network* (DAN) learns to predict the discriminative attributes of referent object  $c_r$  and context  $c_c$  without direct supervision at the attribute level, but relying only on discriminativeness information (e.g., for the objects in the first row of Table 1, the gold vector would contain 1 for *has\_tail*, but 0 for both *is\_green* and *has\_legs*). Still, the model is implicitly encouraged to embed objects into a consistent attribute space, to generalize across the discriminativeness vectors of different training pairs, so it also effectively learns to annotate objects with visual attributes.

Figure 2 presents a schematic view of DAN, focusing on a single attribute. The model is presented with two concepts  $\langle \text{CAT}, \text{SOFA} \rangle$ , and randomly samples a visual instance of each. The instance visual vectors  $\mathbf{v}$  (i.e., ConvNet second-to-last fc layers) are mapped into *attribute* vectors of dimensionality  $|V|$  (cardinality of all available attributes), using weights  $\mathbf{M}^a \in \mathbf{R}^{4096 \times |V|}$  shared between the two concepts. Intuitively, this layer should learn whether an attribute is active for a specific object, as this is crucial for determining whether the attribute is discriminative for an object pair. In Section 5, we present experimental evidence corroborating this hypothesis.

In order to capture the *pairwise* interactions between attribute vectors, the model proceeds by concatenating the two units associated with the *same* visual attribute  $v$  across the two objects (e.g., the units encoding information about *has\_tail*) and pass them as input to the *discriminative* layer.

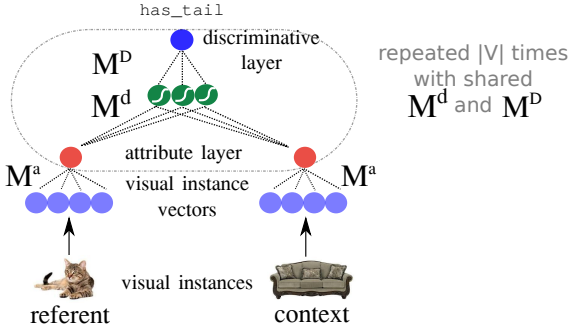


Figure 2: Schematic representation of DAN. For simplicity, the prediction process is only illustrated for `has_tail`

The discriminative layer processes the two units by applying a linear transformation with weights  $M^d \in \mathbf{R}^{2 \times h}$ , followed by a sigmoid activation function, finally deriving a single value by another linear transformation with weights  $M^D \in \mathbf{R}^{h \times 1}$ . The output  $\hat{d}_v$  encodes the predicted degree of discriminativeness of attribute  $v$  for the specific reference-context pair. The same process is applied to all attributes  $v \in V$ , to derive the estimated discriminativeness vector  $\hat{\mathbf{d}}$ , using the same shared weights  $M^d$  and  $M^D$  for each attribute.

To learn the parameters  $\theta$  of the model (i.e.  $M^a$ ,  $M^d$  and  $M^D$ ), given training data  $\langle c_r, c_c, \mathbf{d}_{c_r, c_c} \rangle$ , we minimize MSE between the gold vector  $\mathbf{d}_{c_r, c_c}$  and model-estimated  $\hat{\mathbf{d}}_{c_r, c_c}$ . We trained the model with rmsprop and with a batch size of 32. All hyperparameters (including the hidden size  $h$  which was set to 60) were tuned to maximize performance on the validation set.

#### 4 Predicting Discriminativeness

We evaluate the ability of the model to predict attributes that discriminate the intended referent from the context. Precisely, we ask the model to return all *discriminative attributes* for a pair, independently of whether they are positive for the referent or for the context (given images of a cat and a building, both `+is_furry` and `-made_of_bricks` are discriminative of the cat).

**Test stimuli** We derive our test stimuli from the VisA test split (see Section 2), containing 2000 pairs. Unlike in training, where the model was presented with specific *visual instances* (i.e., single images), for evaluation we use *visual concepts* (CAT, BED), which we derive by *averaging* the vectors of all images associated to an object (i.e., deriving CAT from all images of cats), due to lack

Model	Precision	Recall	F1
DAN	0.66	0.49	0.56
attribute+sym. difference	0.64	0.48	0.55
no attribute layer	0.63	0.33	0.43
Random baseline	0.16	0.16	0.16

Table 2: Predicting discriminative features

of gold information on *per-image* attributes.

**Results** We compare DAN against a random baseline based on per-attribute discriminativeness probabilities estimated from the training data and an ablation model without attribute layer. We test moreover a model that is trained with supervision to predict attributes and then deterministically computes the discriminative attributes. Specifically, we implemented a neural network with one hidden layer, which takes as input a visual instance, and it is trained to predict its gold attribute vector, casting the problem as logistic regression, thus relying on supervision at the attribute level. Then, given two paired images, we let the model generate their predicted attribute vectors and compute the discriminative attributes by taking the symmetric difference of the predicted attribute vectors as we do for DAN. For the DAN and its ablation, we use a 0.5 threshold to deem an attribute discriminative, without tuning.

The results in Table 2 confirm that, with appropriate supervision, DAN performs discriminativeness prediction reasonably well – indeed, as well as the model with similar parameter capacity requiring direct supervision on an attribute-by-attribute basis, followed by the symmetric difference calculation. Interestingly, allowing the model to embed visual representations into an intermediate attribute space has a strong positive effect on performance. Intuitively, since DAN is evaluated on novel concepts, the mediating attribute layer provides more high-level semantic information helping generalization, at the expense of extra parameters compared to the ablation without attribute layer.

#### 5 Predicting Attributes

Attribute learning is typically studied in supervised setups (Ferrari and Zisserman, 2007; Farhadi et al., 2009; Russakovsky and Fei-Fei, 2010). Our model learns to embed visual objects in an attribute space through indirect supervision about attribute discriminativeness for specific  $\langle \text{referent}, \text{context} \rangle$  pairs. Attributes are *never*

explicitly associated to a specific concept during training. The question arises of whether discriminativeness pushes the model to learn plausible concept attributes. Note that the idea that the semantics of attributes arises from their distinctive function within a communication system is fully in line with the classic structuralist view of linguistic meaning (Geeraerts, 2009).

To test our hypothesis, we feed DAN the same test stimuli (visual concept vectors) as in the previous experiment, but now look at activations in the attribute layer. Since these activations are real numbers whereas gold values (i.e., the visual attributes in the ViSA dataset) are binary, we use the validation set to learn the threshold to deem an attribute active, and set it to 0.5 without tuning. Note that no further training and no extra supervision other than the discriminativeness signal are needed to perform attribute prediction. The resulting binary attribute vector  $\hat{\mathbf{p}}_c$  for concept  $c$  is compared against the corresponding gold attribute vector  $\mathbf{p}_c$ .

**Results** We compare DAN to the random baseline and to an explicit *attribute classifier* similar to the one used in the previous experiment, i.e., a one-hidden-layer neural network trained with logistic regression to predict the attributes. We report moreover the best F1 score of Silberer et al. (2013), who learn a SVM for each visual attribute based on HOG visual features. Unlike in our setup, in theirs, images for the same concept are used both for training and to derive visual attributes (our setup is “zero-shot” at the concept level, i.e., we predict attributes of concepts not seen in training). Thus, despite the fact that they used presumably less accurate pre-CNN visual features, the setup is much easier for them, and we take their performance to be an upper bound on ours.

DAN reaches, and indeed surpasses, the performance of the model with direct supervision at the attribute level, confirming the power of discriminativeness as a driving force in building semantic representations. The comparison with Silberer’s model suggests that there is room for improvement, although the noise inherent in concept-level annotation imposes a relatively low bound on realistic performance.

Model	Precision	Recall	F1
DAN	0.58	0.64	0.61
direct supervision	0.56	0.60	0.58
Silberer et. al.	0.70	0.70	0.70
Random baseline	0.13	0.12	0.12

Table 3: Predicting concept attributes

## 6 Evaluating Referential Success

We finally ran a pilot study testing whether DAN’s ability to predict discriminative attributes at the concept level translates into producing features that would be useful in constructing successful referential expressions for specific object instances.

**Test stimuli** Our starting point is the ReferIt dataset (Kazemzadeh et al., 2014), consisting of REs denoting objects (delimited by bounding boxes) in natural images. We filter out any  $\langle \text{RE}, \text{bounding box} \rangle$  pair whose RE does not overlap with our attribute set  $V$  and annotate the remaining ones with the overlapping attribute, deriving data of the form  $\langle \text{RE}, \text{bounding box}, \text{attribute} \rangle$ . For each intended referent of this type, we sample as context another  $\langle \text{RE}, \text{bounding box} \rangle$  pair such that (i) the context RE does not contain the referent *attribute*, so that the latter is a likely discriminative feature; (ii) referent and context come from different images, so that their bounding boxes do not accidentally overlap; (iii) there is maximum word overlap between referent and contexts REs, creating a realistic referential ambiguity setup (e.g., two cars, two objects in similar environments). Finally we sample maximally 20  $\langle \text{referent}, \text{context} \rangle$  pairs per *attribute*, resulting in 790 test items. For each referent and context we extract CNN visual vectors from their bounding boxes, and feed them to DAN to obtain their discriminative attributes. Note that we used the ViSA-trained DAN for this experiment as well.

**Results** For 12% of the test  $\langle \text{referent}, \text{context} \rangle$  pairs, the discriminative *attribute* is contained in the set of discriminative attributes predicted by DAN. A random baseline estimated from the distribution of attributes in the ViSA dataset would score 15% recall. This baseline however on average predicts 20 discriminative attributes, whereas DAN activates, only 4. Thus, the baseline has a trivial recall advantage.

In order to evaluate whether in general the discriminative attributes activated by DAN would lead to referential success, we further sampled a

subset of 100 ⟨referent, context⟩ test pairs. We presented them separately to two subjects (one a co-author of this study) together with the attribute that the model activated with the largest score (see Figure 1 for examples). Subjects were asked to identify the intended referent based on the attribute. If both agreed on the same referent, we achieved referential success, since the model-predicted attribute sufficed to coherently discriminate between the two images. Encouragingly, the subjects agreed on 78% of the pairs ( $p < 0.001$  when comparing against chance guessing, according to a 2-tailed binomial test). In cases of disagreement, the predicted attribute was either too generic or very salient in both objects, a behaviour observed especially in same-category pairs (e.g., `is_round` in Figure 1).

## 7 Conclusion

We presented DAN, a model that, given a referent and a context, learns to predict their discriminative features, while also inferring visual attributes of concepts as a by-product of its training regime. While the predicted discriminative attributes can result in referential success, DAN is currently lacking all other properties of reference (Grice, 1975) (salience, linguistic and pragmatic felicity, etc). We are currently working towards adding communication (thus simulating a speaker-listener scenario (Golland et al., 2010)) and natural language to the picture.

## Acknowledgments

This work was supported by ERC 2011 Starting Independent Research Grant n. 283554 (COMPOSES). We gratefully acknowledge the support of NVIDIA Corporation with the donation of the GPUs used for this research.

## References

Barbara Abbott. 2010. *Reference*. Oxford University Press, Oxford, UK.

Robert Dale and Nicholas Haddock. 1991. Content determination in the generation of referring expressions. *Computational Intelligence*, 7(4):252–265.

Robert Dale and Ehud Reiter. 1995. Computational interpretations of the gricean maxims in the generation of referring expressions. *Cognitive science*, 19(2):233–263.

Jia Deng, Wei Dong, Richard Socher, Lia-Ji Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *Proceedings of CVPR*, pages 248–255, Miami Beach, FL.

Ali Farhadi, Ian Endres, Derek Hoiem, and David Forsyth. 2009. Describing objects by their attributes. In *Proceedings of CVPR*, pages 1778–1785, Miami Beach, FL.

Vittorio Ferrari and Andrew Zisserman. 2007. Learning visual attributes. In *Proceedings of NIPS*, pages 433–440, Vancouver, Canada.

Jakob N. Foerster, Yannis M. Assael, Nando de Freitas, and Shimon Whiteson. 2016. Learning to communicate to solve riddles with deep distributed recurrent q-networks. Technical Report arXiv:1602.02672.

Dirk Geeraerts. 2009. *Theories of lexical semantics*. Oxford University Press, Oxford, UK.

Dave Golland, Percy Liang, and Dan Klein. 2010. A game-theoretic approach to generating spatial descriptions. In *Proceedings of the 2010 conference on empirical methods in natural language processing*, pages 410–419. Association for Computational Linguistics.

Herbert P Grice. 1975. *Logic and conversation*. Syntax and Semantics.

Karl Moritz Hermann, Tomáš Kočiský, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. In *Advances in Neural Information Processing Systems (NIPS)*.

Felix Hill, Antoine Bordes, Sumit Chopra, and Jason Weston. 2015. The Goldilocks principle: Reading children’s books with explicit memory representations. <http://arxiv.org/abs/1511.02301>.

Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara L Berg. 2014. Referitgame: Referring to objects in photographs of natural scenes. In *EMNLP*, pages 787–798.

Emiel Kraahmer and Kees Van Deemter. 2012. Computational generation of referring expressions: A survey. *Computational Linguistics*, 38(1):173–218.

Tomas Mikolov, Armand Joulin, and Marco Baroni. 2015. A roadmap towards machine intelligence. *arXiv preprint arXiv:1511.08130*.

Margaret Mitchell, Kees van Deemter, and Ehud Reiter. 2010. Natural reference to objects in a visual domain. In *Proceedings of the 6th international natural language generation conference*, pages 95–104. Association for Computational Linguistics.

Olga Russakovsky and Li Fei-Fei. 2010. Attribute learning in large-scale datasets. In *Proceedings of ECCV*, pages 1–14.

- John R. Searle. 1969. *Speech Acts: An Essay in the Philosophy of Language*. Cambridge University Press.
- Carina Silberer, Vittorio Ferrari, and Mirella Lapata. 2013. Models of semantic representation with visual attributes. In *Proceedings of ACL*, pages 572–582, Sofia, Bulgaria.
- Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*.
- Luc Steels. 2003. Social language learning. In Mario Tokoro and Luc Steels, editors, *The Future of Learning*, pages 133–162. IOS, Amsterdam.
- Andrea Vedaldi and Karel Lenc. 2015. *MatConvNet – Convolutional Neural Networks for MATLAB*. Proceeding of the ACM Int. Conf. on Multimedia.
- Adam Vogel, Max Bodoia, Christopher Potts, and Daniel Jurafsky. 2013. Emergence of gricean maxims from multi-agent decision theory. In *HLT-NAACL*, pages 1072–1081.