

Combining distributional and referential information for naming objects through cross-modal mapping and direct word prediction

Anonymous ACL submission

Abstract

We compare three recent models of referential word meaning that link visual object representations to lexical representations in a distributional vector space, either directly through cross-modal mapping or indirectly through visual predictors for individual words. We use these models to predict object names as they could be used in naturalistic referring expressions. We find that cross-modal mapping generally produces semantically appropriate and mutually highly similar object names in its top- n list, but sometimes fails to make desired distinctions. Visual word predictors, on the other hand, can react to more subtle visual distinctions and select specific terms, but sometimes stray taxonomically very far from the correct one. Combination of the approaches improves over the individual predictions in a standard naming task. All approaches can be extended to the zero-shot naming case, where the correct name is one for which no instances were seen during training; again they show complementary strengths and weaknesses, depending on the setup and the lexical relation of the unattested object name to known ones.

1 Introduction

Expressions referring to objects in visual scenes typically include a word naming the *type* of the object: E.g., “house” in Figure 1 (a), or, as a very general type, “thingy” in Figure 1 (d). Determining such a name is a crucial step for referring expression generation (REG) systems, as many other decisions, concerning e.g. the selection of attributes, follow from it (Dale and Reiter,

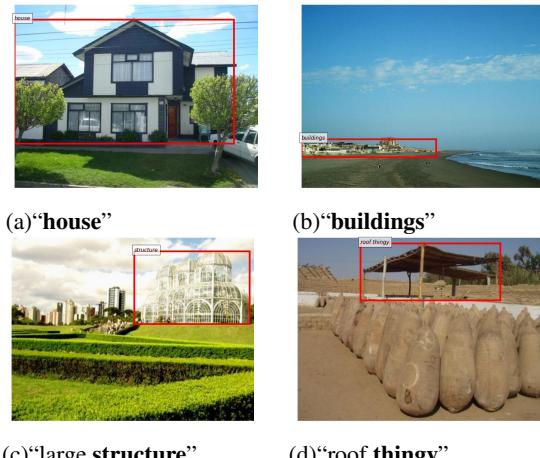


Figure 1: Examples of object names in the REFERIT corpus referring to instances of buildings

1995; Krahmer and Van Deemter, 2012). For a long time, however, research on REG mostly assumed the availability of symbolic representations of referent and scene, and sidestepped questions about how speakers actually choose these names, due to the lack of models capable of capturing what a word like *house* refers to in the real world.

Recent advances in image processing promise to fill this gap, with state-of-the-art computer vision systems being able to classify images into thousands of different categories (eg. Szegedy et al. (2015)). However, classification is not naming (Ordonez et al., 2016). Classification schemes are typically designed to be “flat”, with labels being on the same ontological level and, ideally, having disjunct extensions. In contrast, humans seem to be more flexible as to the chosen level of generality. Depending on the prototypicality of the object to name, and possibly other visual properties, a general name might be more or less appropriate. For instance, a robin can be named “bird”, but a penguin is better referred to as “penguin” (Rosch, 1978); along the same lines, the rather unusual building in Figure 1 that is not easy to otherwise

100 categorise was named “*structure*”.

101 Other work at the intersection of image and
 102 language processing has investigated models that
 103 learn to directly associate visual objects with
 104 a representation of word meaning, for example
 105 through cross-modal transfer into distributional
 106 vector spaces. Under the assumption that such se-
 107 mantic spaces represent, in some form at least, tax-
 108 onomic knowledge, this makes labels on different
 109 levels of specificity available for a given object.
 110 Moreover, if the mapping is sufficiently general, it
 111 should be able to map objects to an appropriate la-
 112 bel, even if during training of the mapping this la-
 113 bel has not been seen (*zero-shot learning*). While
 114 indeed performing with some promise on this task
 115 (Lazaridou et al., 2014), this approach does not
 116 generally outperform standard object classifica-
 117 tion with known categories (Frome et al., 2013;
 118 Norouzi et al., 2013).

119 This paper pursues the hypothesis that an accu-
 120 rate model of referential word meaning does not
 121 need to fully integrate visual and lexical knowl-
 122 edge (e.g. as expressed in a distributional vector
 123 space), but at the same time, has to go beyond
 124 treating words as independent labels. We extend
 125 upon work on learning models of referential word
 126 use from corpora of images paired with referring
 127 expressions (Schlangen et al., 2016; Anonymous,
 128 in press) that treats words as individual predictors
 129 capturing referential appropriateness. We explore
 130 different ways of linking these predictors to distri-
 131 butional knowledge, during application and during
 132 training. We find that these improve over direct
 133 cross-modal mapping and direct visual classifica-
 134 tion in a standard and a zero-shot setup of an ob-
 135 ject naming task, as they allow for a more flexi-
 136 ble combination of lexical and visual information
 137 when modeling referential meaning.

138 2 Related Work

139 **Grounding and Reference** An early example
 140 for work in REG that goes beyond Dale and Re-
 141 iter (1995)’s dominant symbolic paradigm is Deb
 142 Roy’s work from the early 2000s (Roy et al.,
 143 2002; Roy, 2002, 2005). More recently, research
 144 on REG, which has traditionally been done on
 145 small toy data sets, is being scaled up to real-
 146 world images (Kazemzadeh et al., 2014; Gkatzia
 147 et al., 2015; Zarrieß and Schlangen, 2016; Mao
 148 et al., 2015). In this paper, we focus on a particu-
 149 lar problem posed by REG on real-world images,

150 namely generating the appropriate head noun for
 151 a given object. Similarly, Ordonez et al. (2016)
 152 have studied the problem of deriving appropriate
 153 object names, or so-called entry-level categories,
 154 from the output of an object recognizer. Their ap-
 155 proach focusses links abstract object categories in
 156 ImageNet to actual words via various translation
 157 procedures. We are interested in learning referen-
 158 tial appropriateness and extensional word mean-
 159 ings directly from actual human referring expres-
 160 sions (REs) paired with objects in images, using an
 161 existing object recognizer for feature extraction.

**162 Multi-modal and cross-modal distributional se-
 163 mantics** Distributional semantic models are a
 164 well-known method for capturing lexical word
 165 meaning in a variety of tasks (Turney and Pan-
 166 tel, 2010; Mikolov et al., 2013; Erk, 2016). Re-
 167 cent work on multi-modal distributional vector
 168 spaces (Feng and Lapata, 2010; Silberer and La-
 169 pata, 2014; Kiela and Bottou, 2014; Lazaridou
 170 et al., 2015b; Kottur et al., 2016) has aimed at cap-
 171 turing semantic similarity even more accurately by
 172 integrating distributional and perceptual features
 173 associated with words (mostly taken from images)
 174 into a single representation. More related to our
 175 work are cross-modal mapping models (Socher
 176 et al., 2013; Frome et al., 2013; Norouzi et al.,
 177 2013; Lazaridou et al., 2014), that learn to trans-
 178 fer a representation of an object or image in the
 179 visual space to a vector in a distributional space.
 180 When tested on standard object recognition tasks,
 181 transfer, however, comes at a price. Frome et al.
 182 (2013) and Norouzi et al. (2013) both find that
 183 it slightly degrades performance as compared to
 184 a plain object classification using standard accu-
 185 racy metrics (called flat “hit @k metric” in their
 186 paper). Interestingly though, Frome et al. (2013)
 187 report better performance using “hierarchical pre-
 188 cision”, which essentially means that transfer pre-
 189 dictors words that are ontologically closer to the gold
 190 label and makes “semantically more reasonable er-
 191 rors”. To the best of our knowledge, this pattern
 192 has not been systematically investigated any fur-
 193 ther. Another known problem with cross-modal
 194 transfer is that it seems to generalize less well than
 195 expected, i.e. tends to reproduce word vectors ob-
 196 served during training (Lazaridou et al., 2015a). In
 197 this work, we present a model that exploits distri-
 198 butional knowledge for learning referential word
 199 meaning as well, but explore and compare differ-
 ent ways of combining visual and lexical aspects
 of referential word meaning.

200 3 Task and Data

201 We define *object naming* as follows: Given an object x in an image, the task is to predict a word
 202 w that could be used as the head noun of a realistic referring expression. (Cf. discussion above:
 203 “*bird*” when naming a robin, but “*penguin*” when
 204 naming a penguin.) To get at this, we develop our
 205 approach using a corpus of referring expressions
 206 produced by human users under natural, interactive
 207 conditions (Kazemzadeh et al., 2014), and
 208 train and test on the corresponding head nouns in
 209 these REs. This is similar to picture naming setups
 210 used in psycholinguistic research (cf. Levelt et al.
 211 (1991)) and based on the simplifying assumption
 212 that the name used for referring to an object can be
 213 determined successfully without looking at other
 214 objects in the image.

215 We now summarise the details of our setup:

216 **Corpus** We train and test on the REFERIT corpus
 217 (Kazemzadeh et al., 2014), which is based
 218 on the SAIAPR image collection (Grubinger et al.,
 219 2006) (99.5k image regions; 120K REs). We follow
 220 (Schlangen et al., 2016) and select words with
 221 a minimum frequency of 40 in these two data sets,
 222 which gives us a vocabulary of 793 words.

223 **Names** For most of our experiments, we only
 224 use a subset of this vocabulary, namely the set of
 225 object names. As the REs contain nouns that
 226 cannot be considered to be names (*background*, *bottom*,
 227 etc.), we extract from the semantically
 228 annotated portion of the REFERIT corpus a list of
 229 names which correspond to ‘entry-level’ nouns in
 230 terms of (Kazemzadeh et al., 2014). This gives
 231 us a list of 159 names. Thus, our experiments are
 232 on a smaller scale as compared to (Ordonez et al.,
 233 2016). Nevertheless, the data is challenging, as the
 234 corpus contains references to objects that fall outside
 235 of the object labeling scheme that available
 236 object recognition systems are typically optimized
 237 for, cf. Hu et al. (2015)’s discussion on “stuff”
 238 entities such as “*sky*” or “*grass*” in the REFERIT data.
 239 For testing, we remove relational REs (containing
 240 a relational preposition such as ‘left of X’), be-
 241 cause here we cannot be sure that the head noun
 242 of the target is fully informative; we also remove
 243 REs with more than one head noun from our list
 244 (i.e. these are mostly relational expressions as well
 245 such as ‘girl laughing at boy’). We pair each im-
 246 age region from the test set with its corresponding
 247 names from the remaining REs.

248 **Image and Word Embeddings** Following
 249 Schlangen et al. (2016), we derive representations
 250 of our visual inputs with a convolutional neural
 251 network, ‘GoogleNet’ (Szegedy et al., 2015),
 252 which was trained on the ImageNet corpus (Deng
 253 et al., 2009), and extract the final fully-connected
 254 layer before the classification layer, to give us a
 255 1024 dimensional representation of the region.
 256 We add 7 features that encode information about
 257 the region relative to the image, thus representing
 258 each object as a vector of 1031 features. As dis-
 259 tributional word vectors, we use the `word2vec`
 260 representations provided by Baroni et al. (2014)
 261 (trained with CBOW, 5-word context window, 10
 262 negative samples, 400 dimensions).

263 4 Three Models of Interfacing Visual and 264 Distributional Information

265 4.1 Direct Cross-Modal Mapping

266 Following e.g. Lazaridou et al. (2014), referential
 267 meaning can be represented as a translation func-
 268 tion that projects visual representations of objects
 269 to linguistic representations of words in a distri-
 270 butional vector space. Thus, in contrast to standard
 271 object recognition systems or the other models we
 272 will use here, cross-modal mapping does not treat
 273 words as individual labels or classifiers, but learns
 274 to directly predict continuous representations of
 275 words in a vector space, such as the space defined
 276 by the `word2vec` embeddings that we use in this
 277 work. This model will be called TRANSFER below.

278 During training, we pair each object with the
 279 distributional embedding of its name, and use
 280 standard Ridge regression for learning the trans-
 281 formation. Lazaridou et al. (2014) and Lazaridou
 282 et al. (2015a) test a range of technical tweaks and
 283 different algorithms for cross-modal mapping. For
 284 ease of comparison with other models, we stick
 285 with simple Ridge Regression in this work.

286 For decoding, we map an object into the dis-
 287 tributional space, and retrieve the nearest neigh-
 288 bors of the predicted vector using cosine simila-
 289 rity. In theory, the model should generalize easily
 290 to words that it has not observed in a pair with an
 291 object during training as it can map an object any-
 292 where in the distributional space.

293 4.2 Lexical Mapping Through Individual 294 Word Classifiers

295 Another approach is to keep visual and distribu-
 296 tional information separate, by training a separate
 297 visual classifier for each word w in the vocabu-

300 lary. Predictions can then be mapped into distributional
 301 space during application time via the vectors
 302 of the predicted words. Here, we use Schlangen
 303 et al. (2016)’s WAC model, building the training
 304 set for each word w as follows: all visual objects
 305 in a corpus that have been referred to as w are
 306 used as positive instances, the remaining objects
 307 as negative instances. Thus, the classifiers learn
 308 to predict referential appropriateness for individual
 309 words based on the visual features of the objects
 310 they refer to, in isolation of other words.

311 During decoding, we apply all word classifiers
 312 from the model’s vocabulary to the given object,
 313 and take the argmax over the individual word
 314 probabilities. The model can be used to predict
 315 names directly, without links into a distributional
 316 space.

317 In order to extend the model’s vocabulary for
 318 zero-shot learning, we follow Norouzi et al. (2013)
 319 and associate the top n words with their corre-
 320 sponding distributional vector and compute the
 321 convex combination of these vectors. Then, in par-
 322 allel to cross-modal mapping, we retrieve the near-
 323 est neighbors of the combined embedding from the
 324 distributional space. Thus, with this model, we use
 325 two different modes of decoding: one that projects
 326 into distributional space, one that only applies the
 327 available word classifiers.

328 4.3 Word Prediction via Cross-Modal 329 Similarity Mapping

331 Finally, we implement an approach that combines
 332 ideas from cross-modal mapping with the WAC
 333 model: we train individual predictors for each
 334 word in the vocabulary, but, during training, we
 335 exploit lexical similarity relations encoded in a
 336 distributional space. Instead of treating a word as a
 337 binary classifier, we annotate its training instances
 338 with a fine-grained similarity signal according to
 339 their object names. When building the training set
 340 for such a word predictor w , instead of simply di-
 341 viding objects into w and $\neg w$ instances, we label
 342 each object with a real-valued similarity obtained
 343 from cosine similarity between w and v in a dis-
 344 tributional vector space, where v is the word that
 345 was used to refer to the object. Thus, we task the
 346 model with jointly learning similarities and refer-
 347 ential appropriateness, by training it with Ridge
 348 regression on a continuous output space. Object
 349 instances where $v = w$ (i.e., the positive instances
 in the binary setup) have maximal similarity; the

350 remaining instances have a lower value which is
 351 more or less close to maximal similarity. This is
 352 the SIM-WAP model, recently proposed in (Anony-
 353 mous).

354 Importantly, and going beyond (Anonymous),
 355 this model allows for an innovative treatment of
 356 words that only exist in a distributional space
 357 (without being paired with visual referents in the
 358 image corpus): as the predictors are trained on a
 359 continuous output space, no genuine positive in-
 360 stances of a word’s referent are needed. When
 361 training a predictor for such a word w , we use
 362 all available objects from our corpus and anno-
 363 tate them with the expected lexical similarity
 364 between w and the actual object names v , which for
 365 all objects will be below the maximal value that
 366 marks genuine positive instances. During decod-
 367 ing, this model does not need to project its pre-
 368 dictions into a distributional space, but it simply
 369 applies all available predictors to the object, and
 370 takes the argmax over the predicted referential ap-
 371 propriateness scores.

372 5 Experiment 1: Naming Objects

373 This Section reports on experiments in a stan-
 374 dard setup of the object naming task where all
 375 object names are paired with visual instances of
 376 their referents during training. In a compara-
 377 ble task, i.e. object recognition with known ob-
 378 ject categories, cross-modal projection or trans-
 379 fer approaches have been reported to perform
 380 worse than standard object classification methods
 381 (Frome et al., 2013; Norouzi et al., 2013). This
 382 seems to suggest that lexical or at least distri-
 383 butional knowledge is detrimental when learning
 384 what a word refers to in the real world and that
 385 referential meaning should potentially be learned
 386 from visual object representation only.

387 5.1 Model comparison

388 **Setup** We use the train/test split of REFERIT data
 389 as in (Schlangen et al., 2016). We consider image
 390 regions with non-relational referring expressions
 391 that contain at least one of the 159 head nouns
 392 from the list of entry-level nouns (see section 3).
 393 This amounts to 6208 image regions for testing
 394 and 73K instances for training.

395 **Results** Table 1 shows accuracies in the object
 396 naming task for the TRANSFER, WAC and SIM-
 397 WAP models according to their accuracies in the
 398 top n , including two variants of WAC where its top

5 and top 10 predictions are project into the distributional space. Overall, the differences in accuracy between the models are small, but the various models that link their predictions to word representations in the distributional space all perform slightly worse than the plain WAC model, i.e. individual word classifiers trained on visual features only. This suggests that referential meanings for a word are learned less accurately when mapping from visual to distributional space, which replicates results reported in the literature on standard object recognition benchmarks.

	hit @k(%)		
	@1	@2	@5
transfer	48.34	60.49	74.89
wac	49.34	61.86	75.35
wac, project top5	48.73	61.10	74.07
wac, project top10	48.68	61.23	74.31
sim-wap	48.13	60.60	75.40

Table 1: Accuracies in object naming

5.2 Model combination

In order to get more insight into why the TRANSFER and SIM-WAP models produce slightly worse results than individual visual word classifiers, we now test to what extent the different models are complementary and combine them by aggregating over their naming predictions. If the models are complementary, their combination should lead to more confident and accurate naming decisions.

Setup We combine TRANSFER, SIM-WAP and WAC by aggregating the scores they predict for different object names for a given object. During testing, we apply all models to an image region and consider words ranked among the top 10. We first normalize the referential appropriateness scores in each top-10 list and then compute their sum. This aggregation scheme will give more weight to words that appear in the top 10 list of different models, and less weight to words that only get top-ranked by a single model. We test on the same data as in Section 5.1.

	hit @k(%)		
	1	5	10
sim-wap + transfer	49.10	61.78	75.81
sim-wap + wac	51.10	63.45	77.92
transfer + wac	51.13	63.76	77.84
wac + transfer + sim-wap	52.19	64.71	78.40

Table 2: Object naming acc., combined models

Results Table 2 shows that all model combinations improve over the results of their isolated models in Table 1, suggesting that WAC, TRANSFER and SIM-WAP indeed do capture complementary aspects of referential word meaning. On their own, the distributionally informed models are less tuned to specific word occurrences than the visual word classifiers in the WAC model, but they can add useful information which leads to a clear overall improvement. We take this as a promising finding, supporting our initial hypothesis that knowledge on lexical distributional meaning should and can be exploited when learning how to use words for reference.

	Av. cosine distance			
	among top k		gold - top k	
	5	10	5	10
transfer	0.68	0.73	0.72	0.75
wac	0.82	0.80	0.82	0.84
sim-wap	0.68	0.74	0.72	0.75

Table 3: Cosine distances between word2vec embeddings of nouns generated in the top k

5.3 Analysis

Figure 2 illustrates objects from our test set where the combination of TRANSFER, SIM-WAP and WAC predicts an accurate name, whereas the models in isolation do not. These examples give some interesting insight into why the models capture different aspects of referential word use and meaning.

Word Similarities Many of the examples in Figure 2 suggest that the object names ranked among the top 3 by the TRANSFER and SIM-WAP model are semantically similar to each other, whereas WAC generates object names on top that describe very different underlying object categories, such as *seal / rock* in Figure 2(a), *animal / lamp* in Figure 2(g) or *chair / shirt* in Figure 2(c). To quantify this general impression, Table 3 shows cosine distances among words in the top n generated by our models, using their word2vec embeddings. The average cosine distance between words in our vocabulary is 0.83. The transfer and sim-wap model rank words on top that are clearly more similar to each other than word pairs on average, whereas words ranked top by the wac model are more dissimilar. This parallels findings by Frome et al. (2013), discussed in Section 2. Additional evaluation metrics, such as success rates

500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
in a human evaluation (cf. Zarrieß and Schlangen
(2016)), would be an interesting direction for more
detailed investigation here.

528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
Word Use But even though the WAC classifiers
lack knowledge on lexical similarities, they seem
to able to detect relatively specific instances of
word use such as *hut* in Figure 2(b), *shirt* in 2(c) or
lamp in 2(h). Here, the combination with TRANS-
FER and SIM-WAP is helpful to give more weight
to the object name that is taxonomically correct
(sometimes pushing up words below the top-3 and
hence not shown in Figure 2). In Figure 1(e), SIM-
WAP and TRANSFER give more weight to typical
names for persons, whereas WAC top-ranks more
unusual names, reflecting that the person is diffi-
cult to identify visually. Another observation is
that the mapping models have difficulties deal-
ing with object names in singular and plural. As
these words have very similar representations in
the distributional space, they are often predicted as
likely variants among the top 10 by SIM-WAP and
TRANSFER, whereas the WAC model seems to pre-
dict inappropriate plural words less often among
the top 3. Such specific phenomena at the intersec-
tion of visual and semantic similarity have found
very little attention in the literature. We will in-
vestigate them further in our Experiments on zero-
shot naming in the following Section.

6 Zero-Shot Naming

550
551
552
553
554
555
556
557
Zero-shot learning is an attractive prospect for
REG from images, as it promises to overcome de-
pendence on pairings of visual instances and nat-
ural names being available for all names, if vis-
ual/referential data can be generalised from other
types of information. Previous work has looked
at the feasibility of zero-shot learning as a func-
tion of semantic similarity or ontological close-
ness between unknown and known categories, and
confirmed the intuition that the task is harder the
less close unknown categories are to known ones
(Frome et al., 2013; Norouzi et al., 2013).

558
559
560
561
562
563
564
565
566
567
568
569
570
Our experiments on object naming in Section 5
suggest that lexical similarities encoded in a dis-
tributional space might not always fully carry over
to referential meaning. This could constitute an
additional challenge for zero-shot learning, as dis-
tributional similarities might be misleading when
the model has to fully rely on them for learning
referential word meanings. Therefore, the fol-
lowing experiments investigate the performance of

550
551
552
553
554
555
556
557
our models in zero-shot naming as a function of
the lexical relation between unknown and known
object names, i.e. namely hypernyms and singu-
lar/plurals. Both relations are typically captured
by distributional models of word meaning in terms
of closeness in the vector space, but their visual
and referential relation is clearly different.

6.1 Vocabulary Splits and Testsets

558
559
560
561
562
563
564
565
566
567
568
569
570
Random As in previous work on zero-shot
learning, we consider zero-shot naming for words
of varying degrees of similarity in our vocabulary.
We randomly split our 159 names from Experi-
ment 1 into 10 subsets. We train the models on
90% of the nouns (and all their visual instances in
the image corpus) and test on the set of image re-
gions that are named with words which the model
did not observe during training. Results reported
in Table 4 on the random test set correspond to
averaged scores from cross-validation over the 10
splits.

571
572
573
574
575
576
577
Hypernyms We manually split the model’s vo-
cabulary into set of hypernyms (see Appendix A)
and the remaining nouns. We train the models on
those 84K image regions that where not named
with a hypernym, and test on 8895 image regions
that were named with a hypernym in the corpus.
We checked that for each of these hypernyms, the
vocabulary contains at least one or two names that
can be considered as hyponyms, i.e. the model
sees objects during training that are instances of
vehicle for example, but never encounters actual
uses of that name. This test set is particularly inter-
esting from an REG perspective, as objects named
with very general terms by human speakers are of-
ten difficult to describe with more common, but
more specific terms, as is illustrated by the uses of
structure and *thingy* in Figure 1.

578
579
580
581
582
583
584
585
586
587
588
589
Singulars/Plurals We pick 68 words from our
vocabulary that can be grouped into 34 singular-
plural noun pairs (see Appendix A). From each
pair, we randomly include the singular or plural
noun in the set of zero-shot nouns. Thus, we make
sure that the model encounters singular and plu-
ral names during training, but it never encounters
both variants of a name. This results in a more
even training/test split, i.e. we train on 23K image
regions and evaluate on 13825 instances.



Figure 2: Examples from object naming experiment where model combination is accurate

Zero-shot names	Model	full vocab				disjoint vocab	
		@ 1	@ 2	@ 5	@ 10	@ 1	@ 2
Random	transfer	0.05	2.38	16.57	35.71	41.49	62.34
	wac, project top10	0.00	4.42	21.16	39.17	38.03	58.07
	wac, project top5	0.00	4.39	21.63	40.01	37.46	57.36
	sim-wap	3.71	13.13	36.49	54.44	42.28	64.26
Hypernyms	transfer	0.07	1.25	7.75	29.93	59.88	73.88
	wac, project top10	0.00	3.01	15.55	36.99	50.51	66.33
	wac, project top5	0.00	2.78	16.75	38.13	47.73	64.38
	sim-wap	3.16	10.33	31.14	49.62	57.55	70.15
Singulars/Plurals	transfer	0.01	22.84	44.30	72.85	34.56	51.79
	wac, project top10	0.00	22.21	43.43	68.95	31.46	48.76
	wac, project top5	0.00	22.18	43.93	69.33	31.46	48.88
	sim-wap	15.39	34.73	56.62	77.32	37.24	54.02

Table 4: Accuracies in zero-shot object naming on different vocabulary splits

700 6.2 Evaluation

701 Some previous work on zero-shot image labeling
 702 assumes additional components that first identify
 703 whether an image should be labelled by a known
 704 or unknown word (Frome et al., 2013). We fol-
 705 low Lazaridou et al. (2014) and let the model de-
 706 cide whether to refer to an object by a known or
 707 unknown name. Related to that, distinct evalua-
 708 tion procedures have been used in the literature on
 709 zero-shot learning:

710 **Testing on full vocabulary** A realistic way to
 711 test zero-shot learning performance is to consider
 712 all words from a given vocabulary during testing,
 713 though the testset only contains instances of ob-
 714 jects that have been named with a ‘zero-shot word’
 715 (for which no visual instances were seen during
 716 training). Accuracies in this setup reflect how well
 717 the model is able to generalize, i.e. how often it
 718 decides to deviate from the words it was trained
 719 on, and (implicitly) predicts that the given object
 720 requires a “new” name. In case of the (i) hyper-
 721 nym and (ii) singular/plural test set, this accuracy
 722 also reflects to what extent the model is able to de-
 723 tect cases where (i) a more general or vague term
 724 is needed, where (ii) an unknown singular/plural
 725 counterpart of a known object type occurs.

726 **Testing on disjoint vocabulary** Alternatively,
 727 the model’s vocabulary can be restricted during
 728 testing to zero-shot words only, such that names
 729 encountered during training and testing are dis-
 730 joint, see e.g. (Lampert et al., 2009, 2013). This
 731 setup factors out the generalization problem, and
 732 assesses to what extent a model is able to cap-
 733 ture the referential meaning of a word that does
 734 not have instances in the training data.

735 6.3 Results

736 As compared to Experiment 1 where models
 737 achieved similar performance, differences are
 738 more pronounced in the zero-shot setup, as shown
 739 in Table 4. In particular, we find that the SIM-
 740 WAP model which induces individual predictors
 741 for words that have not been observed in the train-
 742 ing data is clearly more successful than TRANS-
 743 FER or WAC that project predictions into the dis-
 744 tributional space. When tested on the full vocabu-
 745 lary, we find that TRANSFER and WAC very rarely
 746 generate names whose referents were excluded
 747 from training, which is in line with observations
 748 made by Lazaridou et al. (2015a). The SIM-WAP

749 predictors generalize much better, in particular on
 750 the singular/plural testset.

751 An interesting exception is the good perfor-
 752 mance of the TRANSFER model on the hypernym
 753 test set, when evaluated with a disjoint vocabu-
 754 lary. This corroborates evidence from Experiment
 755 1, namely that the transfer model captures tax-
 756 onomic aspects of object names better than the
 757 other models. Projection via individual word clas-
 758 sifiers, on the other hand, seems to generalize bet-
 759 ter than TRANSFER, at least when looking at ac-
 760 curacies @2 ... @10. Thus, combining several
 761 vectors predicted by a model of referential word
 762 meaning can provide additional information, as
 763 compared to mapping an object to a single vec-
 764 tor in distributional space. More work is needed to
 765 establish how these approaches can be integrated
 766 more effectively.

767 7 Discussion and Conclusion

768 In this paper, we have investigated models of refer-
 769 ential word meaning, using different ways of com-
 770 bining visual information about a word’s referent
 771 and distributional knowledge about its lexical sim-
 772 ilarities. Previous cross-modal mapping models
 773 essentially force semantically similar objects to be
 774 mapped into the same area in the semantic space
 775 regardless of their actual visual similarity. We
 776 found that cross-modal mapping produces seman-
 777 tically appropriate and mutually highly similar ob-
 778 ject names in its top- n list, but does not preserve
 779 differences in referential word use (e.g. appropri-
 780 ateness of *person* vs. *woman*) especially within the
 781 same semantic field. We have shown that it is
 782 beneficial for performance in standard and zero-
 783 shot object naming to treat words as individual
 784 predictors that capture referential appropriateness
 785 and are only indirectly linked to a distributional
 786 space, either through lexical mapping during ap-
 787 plication or through cross-modal similarity map-
 788 ping during training. As we have tested these ap-
 789 proaches on a rather small vocabulary, which may
 790 limit generality of conclusions, future work will
 791 be devoted to scaling up these findings to larger
 792 test sets, as e.g. recently collected through con-
 793 versational agents (Das et al., 2016) that circumvent
 794 the need for human-human interaction data. Also
 795 from an REG perspective, various extensions of
 796 this approach are possible, such as the inclusion of
 797 contextual information during object naming and
 798 its combination with attribute selection.

800 References

801 Anonymous. in press.

802 Marco Baroni, Georgiana Dinu, and Germán
803 Kruszewski. 2014. Don't count, predict! a
804 systematic comparison of context-counting vs.
805 context-predicting semantic vectors. In *ACL (1)*.
806 pages 238–247.

807 Robert Dale and Ehud Reiter. 1995. Computational
808 interpretations of the gricean maxims in the genera-
809 tion of referring expressions. *Cognitive Science*
810 19(2):233–263.

811 Abhishek Das, Satwik Kottur, Khushi Gupta,
812 Avi Singh, Deshraj Yadav, José M. F.
813 Moura, Devi Parikh, and Dhruv Batra.
814 2016. Visual dialog. *CoRR* abs/1611.08669.
815 <http://arxiv.org/abs/1611.08669>.

816 Jia Deng, W. Dong, Richard Socher, L.-J. Li, K. Li, and
817 L. Fei-Fei. 2009. ImageNet: A Large-Scale Hierar-
818 chical Image Database. In *CVPR09*.

819 Katrin Erk. 2016. What do you know about
820 an alligator when you know the company it
821 keeps? *Semantics and Pragmatics* 9(17):1–63.
<https://doi.org/10.3765/sp.9.17>.

822 Yansong Feng and Mirella Lapata. 2010. Visual in-
823 formation in semantic representation. In *Human
824 Language Technologies: The 2010 Annual Confer-
825 ence of the North American Chapter of the Associa-
826 tion for Computational Linguistics*. Association for
827 Computational Linguistics, pages 91–99.

828 Andrea Frome, Greg S Corrado, Jon Shlens, Samy
829 Bengio, Jeff Dean, Marc Aurelio Ranzato, and
830 Tomas Mikolov. 2013. Devise: A deep visual-
831 semantic embedding model. In C. J. C. Burges,
832 L. Bottou, M. Welling, Z. Ghahramani, and K. Q.
833 Weinberger, editors, *Advances in Neural Infor-
834 mation Processing Systems 26*, Curran Associates, Inc.,
835 pages 2121–2129.

836 Dimitra Gkatzia, Verena Rieser, Phil Bartie, and
837 William Mackaness. 2015. From the virtual to the
838 real world: Referring to objects in real-world spatial
839 scenes. In *Proceedings of EMNLP 2015*. Associa-
840 tion for Computational Linguistics.

841 Michael Grubinger, Paul Clough, Henning Müller, and
842 Thomas Deselaers. 2006. The IAPR TC-12 bench-
843 mark: a new evaluation resource for visual informa-
844 tion systems. In *Proceedings of the International
845 Conference on Language Resources and Evaluation
846 (LREC 2006)*. Genoa, Italy, pages 13–23.

847 Ronghang Hu, Huazhe Xu, Marcus Rohrbach, Jiashi
848 Feng, Kate Saenko, and Trevor Darrell. 2015. Natu-
849 ral language object retrieval. *CoRR* abs/1511.04164.
<http://arxiv.org/abs/1511.04164>.

850 Sahar Kazemzadeh, Vicente Ordonez, Mark Matten,
851 and Tamara L Berg. 2014. ReferItGame: Referring
852

853 to Objects in Photographs of Natural Scenes. In *Pro-
854 ceedings of the Conference on Empirical Methods
855 in Natural Language Processing (EMNLP 2014)*.
856 Doha, Qatar, pages 787–798.

857 Douwe Kiela and Léon Bottou. 2014. Learning Image
858 Embeddings using Convolutional Neural Networks
859 for Improved Multi-Modal Semantics. In *Pro-
860 ceedings of the Conference on Empirical Methods in
861 Natural Language Processing (EMNLP-14)*.

862 Satwik Kottur, Ramakrishna Vedantam, José MF
863 Moura, and Devi Parikh. 2016. Visual word2vec
864 (vis-w2v): Learning visually grounded word embed-
865 dings using abstract scenes. In *Proceedings of the
866 IEEE Conference on Computer Vision and Pattern
867 Recognition*. pages 4985–4994.

868 Emiel Krahmer and Kees Van Deemter. 2012. Compu-
869 tational generation of referring expressions: A sur-
870 vey. *Computational Linguistics* 38(1):173–218.

871 Christoph H Lampert, Hannes Nickisch, and Stefan
872 Harmeling. 2009. Learning to detect unseen object
873 classes by between-class attribute transfer. In *IEEE
874 Computer Vision and Pattern Recognition*. IEEE,
875 pages 951–958.

876 Christoph H. Lampert, Hannes Nickisch, and Stefan
877 Harmeling. 2013. Attribute-based classification for
878 zero-shot visual object categorization. *IEEE Trans-
879 actions on Pattern Analysis and Machine Intelli-
880 gence* 36(3):453–465.

881 Angeliki Lazaridou, Elia Bruni, and Marco Baroni.
882 2014. Is this a wampimuk? Cross-modal map-
883 ping between distributional semantics and the visual
884 world. In *Proceedings of the 52nd Annual Meet-
885 ing of the Association for Computational Linguis-
886 tics (Volume 1: Long Papers)*. pages 1403–1414.

887 Angeliki Lazaridou, Georgiana Dinu, and Marco
888 Baroni. 2015a. Hubness and pollution: Delv-
889 ing into cross-space mapping for zero-shot learn-
890 ing. In *Proceedings of the 53rd Annual Meet-
891 ing of the Association for Computational Linguis-
892 tics and the 7th International Joint Con-
893 ference on Natural Language Processing (Vol-
894 ume 1: Long Papers)*. Association for Computa-
895 tional Linguistics, Beijing, China, pages 270–280.
<http://www.aclweb.org/anthology/P15-1027>.

896 Angeliki Lazaridou, Nghia The Pham, and Marco Baroni.
897 2015b. Combining language and vision with a
898 multimodal skip-gram model. In *Proceedings of the
899 2015 Conference of the North American Chapter of
900 the Association for Computational Linguistics: Hu-
901 man Language Technologies*. Association for Com-
902 putational Linguistics, Denver, Colorado, pages
903 153–163. <http://www.aclweb.org/anthology/N15-1016>.

904 Willem JM Levelt, Herbert Schriefers, Dirk Vor-
905 berg, Antje S Meyer, Thomas Pechmann, and Jaap
906 Havinga. 1991. The time course of lexical access in
907

900	speech production: A study of picture naming. <i>Psychological review</i> 98(1):122.	950
901		951
902	Junhua Mao, Jonathan Huang, Alexander Toshev, Oana Camburu, Alan L. Yuille, and Kevin Murphy. 2015. Generation and comprehension of unambiguous object descriptions. <i>ArXiv / CoRR</i> abs/1511.02283. http://arxiv.org/abs/1511.02283.	952
903		953
904		954
905		955
906		956
907	Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, <i>Advances in Neural Information Processing Systems 26</i> , pages 3111–3119.	957
908		958
909		959
910		960
911		961
912		962
913	Mohammad Norouzi, Tomas Mikolov, Samy Bengio, Yoram Singer, Jonathon Shlens, Andrea Frome, Greg S Corrado, and Jeffrey Dean. 2013. Zero-shot learning by convex combination of semantic embeddings.	963
914		964
915		965
916		966
917	Vicente Ordonez, Wei Liu, Jia Deng, Yejin Choi, Alexander C. Berg, and Tamara L. Berg. 2016. Learning to name objects. <i>Commun. ACM</i> 59(3):108–115.	967
918		968
919		969
920		970
921	Eleanor Rosch. 1978. Principles of Categorization. In Eleanor Rosch and Barbara B. Lloyd, editors, <i>Cognition and Categorization</i> , Lawrence Erlbaum, Hillsdale, N.J., USA, pages 27—48.	971
922		972
923		973
924		974
925	Deb Roy. 2005. Grounding words in perception and action: Computational insights. <i>Trends in Cognitive Science</i> 9(8):389–396.	975
926		976
927		977
928	Deb Roy, Peter Gorniak, Niloy Mukherjee, and Josh Juster. 2002. A trainable spoken language understanding system for visual object selection. In <i>Proceedings of the International Conference on Speech and Language Processing 2002 (ICSLP 2002)</i> . Colorado, USA.	978
929		979
930		980
931		981
932		982
933	Deb K. Roy. 2002. Learning visually-grounded words and syntax for a scene description task. <i>Computer Speech and Language</i> 16(3).	983
934		984
935		985
936	David Schlangen, Sina Zarriess, and Casey Kennington. 2016. Resolving references to objects in photographs using the words-as-classifiers model. In <i>Proceedings of the 54rd Annual Meeting of the Association for Computational Linguistics (ACL 2016)</i> .	986
937		987
938		988
939		989
940		990
941	Carina Silberer and Mirella Lapata. 2014. Learning grounded meaning representations with autoencoders. In <i>Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> . Association for Computational Linguistics, Baltimore, Maryland, pages 721–732. http://www.aclweb.org/anthology/P14-1068.	991
942		992
943		993
944		994
945		995
946		996
947	Richard Socher, Milind Ganjoo, Christopher D Manning, and Andrew Ng. 2013. Zero-shot learning through cross-modal transfer. In <i>Advances in neural information processing systems</i> . pages 935–943.	997
948		998
949		999