

Language Grounding with 3D Objects

Anonymous Author(s)

Affiliation

Address

email

Abstract:

Seemingly simple natural language requests to a robot are generally underspecified, for example *Can you bring me the wireless mouse?* When viewing mice on the shelf, the number of buttons or presence of a wire may not be visible from certain angles or positions. Flat images of candidate mice may not provide the discriminative information needed for *wireless*. The world, and objects in it, are not flat images but complex 3D shapes. If a human requests an object based on any of its basic properties, such as color, shape, or texture, robots should perform the necessary exploration to accomplish the task. In particular, while substantial effort and progress has been made on understanding explicitly visual attributes like color and category, comparatively little progress has been made on understanding language about shapes and contours. In this work, we introduce a novel reasoning task that targets both visual and non-visual language about 3D objects. Our new benchmark **ShapeNet Annotated with Referring Expressions (SNARE)** requires a model to choose which of two objects is being referenced by a natural language description. We introduce several CLIP-based [1] models for distinguishing objects and demonstrate that while recent advances in jointly modeling vision and language are useful for robotic language understanding, it is still the case that these models are weaker at understanding the 3D nature of objects – properties which play a key role in manipulation. In particular, we find that adding view estimation to language grounding models improves accuracy on both SNARE and when identifying objects referred to in language on a robot platform.

Keywords: Benchmark, Language Grounding, Vision, 3D

1 Introduction

Joint language and vision models are often trained on image captions which have a bias towards canonically “visual” attributes of objects, such as color, rather than functional ones like shape. Image captions omit properties that require understanding objects as 3D, rather than flat, concepts. In this work, we show how the effect of this disconnect is that while robotics research has benefited greatly from advances in computer vision, vision-and-language models cannot always be directly applied to robotics. People use geometric and physical properties of objects when describing them. For example, a robot carrying out a request to *Bring me the mug with the wide handle* needs to be able to spot the *wide handle*, even if the mug rests on a countertop with the handle oriented out of view, but a robot can rotate such an object to investigate in further detail. Generally, identifying objects based on natural language is more challenging from non-canonical viewpoints, which are common for robots in home, office, and industrial environments.

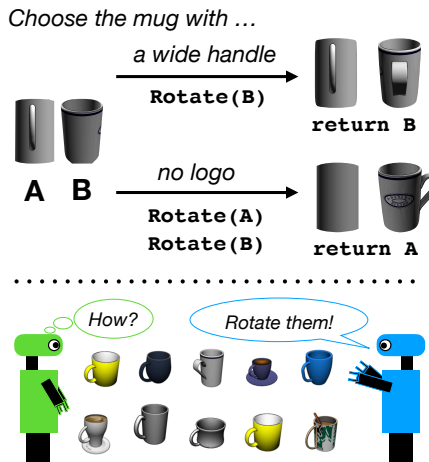


Figure 1: Deciding between 3D objects described in natural language can require rotating objects to see multiple faces.

44 The codification and commodification of ResNet [2] architectures pretrained on ImageNet [3] has
45 yielded immense progress in off-the-shelf computer vision techniques. Such representations provide
46 a strong foundation for visual reasoning. With the introduction of pretrained Transformer [4] archi-
47 tectures like BERT [5], natural language processing has similar off-the-shelf tools for representing
48 language. Modern visual language grounding is done through massive, internet-scale pretraining
49 to combine these strengths: aligning images with their captions and content descriptions in natural
50 language using Transformer models [6, 1, 7, 8].

51 However, internet images suffer from a type of *reporting bias*. These images are captured by sighted
52 humans, from human-centric vantage points [9]. There is a domain shift between such canonical
53 depictions of objects in images and objects as seen through a robot camera [10, 11], as studied in
54 computer vision between such images and those captured by the blind [12]. For example, images of
55 *mugs* online almost exclusively feature a full or three-fourths view of the handle, while a mug in the
56 wild can easily be oriented such that the handle is occluded from the camera by the mug body.

57 A person requesting a mug with a wide handle has a *mental model* of the referent object that includes
58 its full 3D spatial information; humans do not imagine objects only in 2 dimensions. Further, someone
59 trying to retrieve said mug will differentiate it from other mugs on the countertop by viewing the
60 mugs from different angles to inspect the handles.

61 Our goal is to similarly imbue robots with three dimensional notions of language grounding, by
62 encouraging them to *rotate* an object when selecting the referent of a language description. We
63 introduce the **ShapeNet Annotated with Referring Expressions (SNARE)** dataset, which provides
64 discriminative natural language descriptions of 3D ShapeNet [13] object models. SNARE includes
65 both *visual* descriptions that focus on colors and object categories and *blindfolded* descriptions that
66 focus on shapes and parts (Section 3). We train agents to select the correct object given its description,
67 and find that estimating the current view and performing rotation on objects to obtain an additional
68 vantage point improves accuracy (Figure 1).

69 In particular, we introduce **Language Grounding through Object Rotation (LAGOR)** which performs
70 view estimation as an auxiliary loss to encourage 3D object understanding (Section 4) during the
71 SNARE task. We find that using multiple views to select the correct object improves over single-view
72 object selection, using the large-scale vision-and-language CLIP [1] model as a backbone to score
73 how well an image and language MATCH. We show that enabling agents to rotate objects to a new
74 view, while performing auxiliary view estimation, can result in higher accuracy at selecting referent
75 objects than consuming panoramic views of objects while simultaneously being more realistic for a
76 physical robot platform.

77 Our key contributions are:

- 78 • SNARE, a benchmark dataset for identifying 3D object models given natural language referring
79 expressions including *visual* and *blindfolded* descriptions;
- 80 • Baseline models for SNARE demonstrating both zero-shot and fine-tuned performance of state-
81 of-the-art vision-and-language models; and
- 82 • LAGOR, an initial SNARE model that estimates current view, performs rotations, then makes a
83 referent prediction, which we demonstrate on a physical robot platform.

84 2 Related Work

85 Using natural language to work with robot partners is a long-standing goal in robotics [14]. We
86 argue that internet-scale, pretrained vision-and-language models offer a powerful starting point for
87 human-robot collaboration. Unlike 2D internet images, physical objects can be picked up and moved
88 by robot agents. Agents can perform information-seeking behaviors on objects to improve their
89 world [15, 16, 17] and language understanding [18, 19]. Such world interaction is inextricable from
90 language grounding [20], motivating language annotations for higher-fidelity referents than static
91 images. We introduce SNARE, comprised of language referring expressions for 3D objects, and the
92 LAGOR model to combine language and multiple object views to achieve better 3D understanding.
93 SNARE extends lines of work tying language to static images, 3D object models, and even physical
94 objects. We demonstrate that LAGOR generalizes to physical objects manipulated by a tabletop
95 robot §5.2. LAGOR is inspired by models that perform information-seeking actions informed by
96 learned world models to better achieve goals.

97 **Image-based Object Identification** Object classification [3, 2] is the first step towards image cap-
 98 tioning and visual question answering [21]. Particular object instances can be found with region seg-
 99 mentation models such as MaskRCNN [22], enabling object referent tasks such as GuessWhat!?! [23].
 100 Combining visual recognition with Transformer-based language understanding to jointly attend to
 101 language and visual tokens leads to improved downstream performance on many vision-and-language
 102 tasks [6, 24, 8]. Joint embedding approaches that learn a shared subspace for representing language
 103 and vision tokens [25] also achieve state-of-the-art performance when trained at scale [1]. We show
 104 that these large-scale, pretrained models fall short of human performance on the SNARE task, and
 105 that synthesizing 2D views from multiple vantage points improves object identification performance.

106 **Language Grounding in 3D** Prior works have associated single-word attributes with 3D object
 107 models based on latent representations of 3D meshes [26]. To learn spatial language, vision-and-
 108 language navigation (VLN) [27] models infer navigation actions from instructions and visual obser-
 109 vations in 3D simulated worlds. Such tasks can be extended to include simulated world interactions
 110 such as picking and placing objects [28] and using appliances and tools [29]. Models for these tasks
 111 extend Transformer representations to include action taking [30, 31] and invoke object classification
 112 methods to create a semantic understanding of the world [32]. SNARE poses a complementary
 113 challenge, providing data for selecting referent objects in the presence of distractors by taking into
 114 account the multiple views possible of objects in 3D space.

115 The work most similar to ours is Shape-
 116 Glot [33], a collection of referring expres-
 117 sions for ShapeNet objects to discriminate
 118 between two distractors. While Shape-
 119 Glot is similar in spirit, it contains train-
 120 ing data only for `chair` objects, and tests
 121 on four additional categories. By contrast,
 122 SNARE spans 262 distinct object cate-
 123 gories (Table 1). ShapeGlot focuses on a
 124 model’s abilities to learn particular parts of
 125 objects, such as chair arms and legs, while
 126 SNARE aims to ground language more
 127 generally to 3D features of objects span-
 128 ning color, shape, category, and parts. Further, ShapeGlot models take in a pointcloud representation
 129 of objects along with visual features, where our LAGOR model does not assume access to the 3D
 130 model of an object and instead operates solely on different 2D views achieved through object rotation.

	Data	Fold	# Cats	# Objs	# Ref Exps
ShapeGlot	Chairs		1	4,511	78,789
	Other		4	200	400
	Total		5	4,711	79,189
SNARE	Train		207	6,153	39,104
	Val		7	371	2,304
	Test		48	1,357	8,751
	Total		262	7,881	50,159

Table 1: Fold summaries in SNARE, with the Shape-
 Glot (SG) benchmark size for contrast.

131 **Physical Object Identification** Modeling the connections between language to physical objects
 132 enables robots to identify objects for manipulation by color, shape, and category [34, 35]. The
 133 physical forces and sounds objects make during manipulation actions can also be associated with
 134 words such as *rattling* and *heavy* for multimodal understanding beyond vision [18, 36]. Prior work
 135 has gathered language annotations for the YCB Benchmark object set [37] to explore how language
 136 descriptions provide priors on object affordances [38]. In our tabletop robot experiments, we use
 137 camera views of novel objects to evaluate zero shot transfer of LAGOR to the real world with minimal
 138 object rotations to achieve language-aligned camera views to select the correct referent object.

139 **Information-Seeking Actions** Embodiment, whether in simulation or the physical world, affords
 140 agents the opportunity to seek out new information to help with a given task. By predicting the new
 141 information that can be gotten from different actions, agents can *prospect* over potential futures for
 142 planning [39, 40, 41]. In the aforementioned VLN task, predicting *what* will be seen along different
 143 potential routes facilitates more efficient navigation [42]. Our LAGOR model estimates the current
 144 object view before performing a rotation to obtain a new view, for example rotating mugs to view
 145 handles head-on when charged to find a *wide handle* (Figure 1).

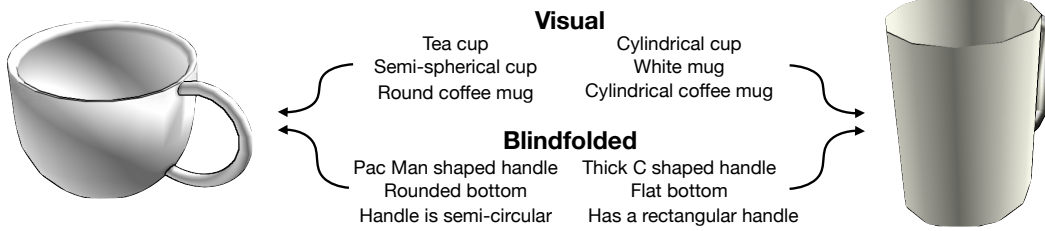


Figure 2: Example object pair and object referring expressions in SNARE. Given a referring expression and one (or more) views of two contrasting objects, a model must decide which object is being referenced by the language. Language annotations are collected in two forms: **Visual** and **Blindfolded**. The latter expressions tasked annotators with describing objects tactilely.

146 3 ShapeNet Annotated with Referring Expressions (SNARE)

147 We introduce SNARE, as a new benchmark for grounding natural language referring expressions to
 148 distinguish 3D objects. The annotations are collected to complement the ACRONYM¹ [43] grasping
 149 dataset and include language that targets both visual and tactile attributes of objects. Our goal is
 150 to enable future research to ground both from multi-view vision, as in this work, and directly from
 151 robotic grasp contact point data (Section 6).

152 To construct SNARE, we select a subset of 7,897 ACRONYM [43] object models from ShapeNet-
 153 Sem [44, 13]. We obtain over 50 thousand natural language referring expressions in English for
 154 these object models via Amazon Mechanical Turk (AMT).² Figure 2 gives an example of referring
 155 expressions collected for two mug objects used to distinguish between the two.

156 To elicit referring expressions with high specificity, we frame the annotation as a discriminative task.
 157 AMT workers were presented with two object models side-by-side from the same ShapeNet category,
 158 for example two different `OutdoorTable` object meshes, and asked to complete sentences like

- 159 • The way to tell Object A from Object B is that Object A looks like a(n) ____
- 160 • Blindfolded, the way to tell Object A from Object B is that Object B is a(n) ____

161 with referring expressions. We displayed the object models as GIFs rendered to give a 360 rotating
 162 view of each object. The resulting SNARE task asks models to take in one such referring expression
 163 and decide whether it applies to Object A or Object B.

164 Priming in the annotation prompts invokes visual features (*... looks like...*) and non-visual features
 165 (*Blindfolded, ...*) in the referring expressions. By pairing objects from the same ShapeNet category,
 166 referring expressions must go beyond categorical information like *brown dresser* to differentiate one
 167 object from the other with higher specificity, similar to ShapeGlot’s choice to use distractor objects
 168 from a single training category [33]. We collected six referring expressions per object—three primed
 169 to be visual and three primed to be blindfolded. Each referring expression was vetted through a
 170 secondary task on AMT where, given the language expression, workers had to correctly select the
 171 referent object. Every referring expression in SNARE was correctly associated to its referent object
 172 by a majority of such annotators.

173 SNARE referring expressions average 4.27 words, with blindfolded expressions longer (4.95) than
 174 visual expressions (3.63). Visual expressions focus more on color and object category words, such as
 175 *table*, than blindfolded expressions. By contrast, blindfolded expressions use more shape and part
 176 words, such as *rectangle* and *legs*. We estimated that blindfolded expressions use more shape words
 177 (14% vs 5% of all words) while visual use more color words (22% vs 1% of all words), by traversing
 178 the WordNet [45] hierarchy for each word and noting whether it is a hyponym of *color* or *shape*.

179 Each SNARE instance is a tuple of (referring expression, Object A, Object B), and models must
 180 select which of Object A or B the referent of the natural language expression. We split these data
 181 instances into train, validation, and test folds by ShapeNet category. We ensure that closely related
 182 categories such as `2Shelves` and `3Shelves` or `DiningTable` and `AccentTable` are within

¹Yes, that dataset title is “ACRONYM” [43].

²Section 7.1 contains additional details about the AMT study.

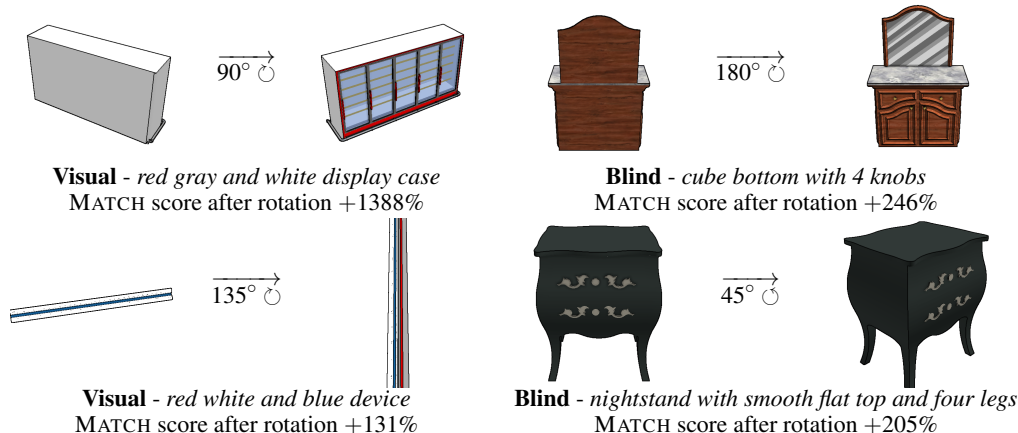


Figure 3: Example MATCH score increases after performing object rotation. Rotations involve exposing colors (top left, bottom left), canonical faces (top right), and parts (bottom right). A robot encountering objects in the wild will find initial object orientations that do not line up with natural language referring expressions for those objects, as reflected in SNARE.

183 the same fold. For example, all shelf-related and bed-related categories are sorted into the train and
 184 test folds, respectively, so that inter-category information does not leak across folds. More details can
 185 be found Section 7.3. Table 1 summarizes the overall data statistics.

186 4 Methods

187 There has been a recent proliferation of off-the-shelf, increasingly large-scale pretrained vision and
 188 language alignment models applicable for language grounding in robotics [1, 6, 7, 8, 46]. We set out
 189 to answer two questions about such models, using SNARE as a testbed:

- 190 • Can existing language and vision models ground language in the SNARE benchmark?
- 191 • Do SNARE models generalize to a robot object selection task better than off-the-shelf models?

192 To answer the first question, we first train a MATCH module that learns a predictive head on top
 193 of a CLIP [1] backbone, and evaluate this module on single and multiple views of 3D objects. To
 194 answer the second question, we introduce the LAGOR model, which uses view estimation as an
 195 auxiliary loss while predicting language expression referent objects for SNARE, and evaluate on
 196 both SNARE §5.1 and a physical robot platform §5.2. LAGOR examines an initial 3D object view
 197 and an additional, post-rotation view, striking a balance between single-view models and those that
 198 attempt to capture each object’s entire 3D structure.

199 4.1 Language-View Match (MATCH) Module

200 The MATCH module takes in a language expression L and single view of an object, V_i , and produces a
 201 match score $s(L, V_i)$. To decide whether Object A or Object B is the referent of a language expression
 202 on SNARE, we interpret $\operatorname{argmax}_{O \in \{A, B\}} s(L, V_{i,O})$ as the model’s prediction.

203 CLIP [1] serves as a backbone on which we add additional, learnable layers for the SNARE
 204 task. We use CLIP’s transformer-based sentence encoder to extract language features $\mathbf{l} : \mathbb{R}^{1 \times 512}$.
 205 We use CLIP ViT-B/32 to extract visual features for image $V_{i,A}$, $\mathbf{v}_{i,A} : \mathbb{R}^{1 \times 512}$, and image $V_{i,B}$,
 206 $\mathbf{v}_{i,B} : \mathbb{R}^{1 \times 512}$. These modality-specific feature vectors are concatenated: $[\mathbf{v}_{i,A}; \mathbf{l}] : \mathbb{R}^{1 \times 1024}$ and
 207 $[\mathbf{v}_{i,B}; \mathbf{l}] : \mathbb{R}^{1 \times 1024}$ and independently run through a learnable, multi-layer perceptron that gradually
 208 reduces the dimensionality from 1024 to 512, then 256, and finally to a single-dimensional score
 209 s . This final score comparison training mirrors the multiple-choice formulation used in existing
 210 unimodal [5] and multimodal transformers [7]. We keep the CLIP encoders frozen during training.

211 ShapeNet objects are 3D meshes, incompatible with the 2D input expected by off-the-shelf vision and
 212 language models. We sample eight rendered viewpoints around each object at 45 degree increments,

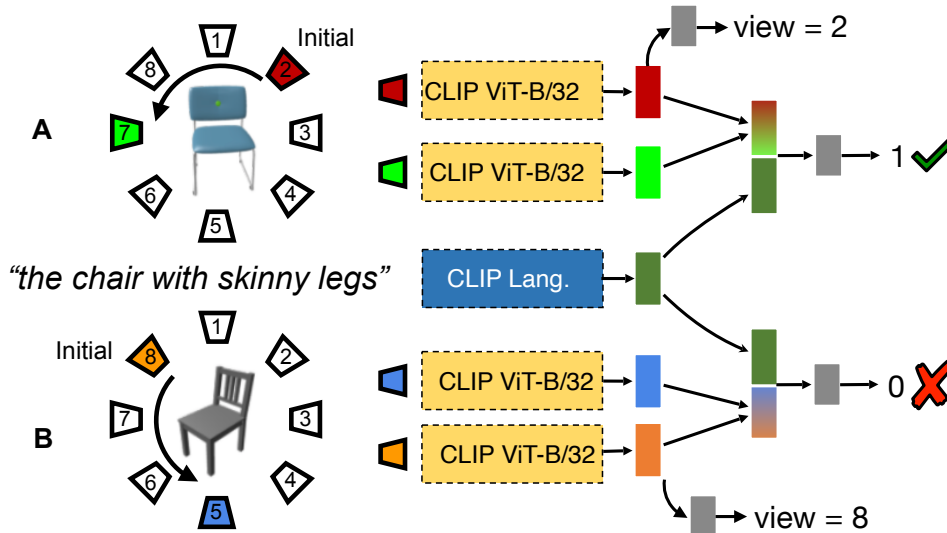


Figure 4: **LAGOR model** relies on the pretrained CLIP [1] architecture as a backbone encoder across multiple views compared to the encoded referring expression. In addition, the model sees substantial improvements from an auxiliary loss that predicts the orientation of the initial view.

213 and compare off-the-shelf CLIP and MATCH module accuracies when considering single or multiple
 214 of these views. To aggregate multiple views, we maxpool over view embedding vectors $\mathbf{v}_{i,A} \dots \mathbf{v}_{j,A}$.
 215 MATCH is trained with cross-entropy-loss, \mathcal{L}_s , to predict a binary label of whether the referring
 216 expression matches the object represented in the view image. MATCH is trained for 50 epochs on
 217 SNARE, with object views chosen at random during each step.

218 **Baselines.** We compare our fine-tuned MATCH against a zero-shot CLIP classifier. Instead of fine-
 219 tuning CLIP, we use the cosine distance between visual and language features to pick the referred
 220 object. That is, we select $\text{argmax}_{O \in \{A,B\}} \mathbf{1} \cdot \mathbf{v}_{i,O}$. We also evaluate against a trained ViLBERT baseline
 221 that consumes multiple object views at once, but find that it does not perform as well as the MATCH
 222 module; more details about the ViLBERT baseline can be found in Section 7.4.

223 4.2 Language Grounding through Object Rotation (LAGOR)

224 A robot tasked with retrieving an object given a natural language expression should not have to gather
 225 a 360 degree view of each candidate referent object to make a decision. We find that taking two
 226 random views of an object as input, as opposed to just one, nearly closes the gap with MATCH module
 227 performance on a 360 view of objects in SNARE (Section 5.1). Further, we have the intuition that
 228 estimating the *current* viewing angle can enable a model to develop a global reference frame for
 229 language grounding against 3D objects (Figure 3). Thus, we propose LAGOR, a model that takes
 230 in two views of each candidate object before selecting the language referent, and performs *view*
 231 *estimation* on candidate objects as an auxiliary prediction (Figure 4).

232 In addition to predicting the object referent using a pre-trained MATCH module that considers two
 233 views, LAGOR predicts the current view of each object using a cross-entropy loss, \mathcal{L}_v , against a
 234 1-hot vector representing a discrete set of 8 views at 45 degree offsets. LAGOR learns a multi-layer
 235 perceptron that takes in a visual embedding V_i and reduces dimensionality from 512 to 256 to 128 to
 236 64 to a vector of 8 logits representing the 8 discrete object views. The MATCH loss, \mathcal{L}_s , is combined
 237 with the view estimation loss for a final loss function $\mathcal{L} = \mathcal{L}_v + 0.2 * \mathcal{L}_s$. LAGOR is trained for 50
 238 epochs on SNARE, with object views chosen at random during each step. An interesting direction for
 239 future work is predicting the *optimal*, rather than random, next view to rotate to when disambiguating
 240 a query, though that could introduce a cyclic dependency with the MATCH module.

241 4.3 Robot Demonstration

242 We evaluate LAGOR trained on SNARE on a robot platform, taking two random views of objects
 243 and evaluating the robot’s ability to select which is the correct referent of a language expression

Model	Views	Validation			Test		
		Visual \subset	Blind \subset	All	Visual \subset	Blind \subset	All
ViLBERT	360	89.5	76.6	83.1	80.2	73.0	76.6
CLIP	360	84.5	66.1	75.3	80.0	61.4	70.9
MATCH	360	90.6	79.3	85.0	85.9	71.3	78.7
CLIP	Single	79.5	65.2	72.3	73.9	60.4	67.3
MATCH	Single	89.4	75.6	82.5	84.1	69.6	77.0
CLIP	Two	81.7	65.5	73.6	76.2	61.0	68.8
MATCH	Two	91.2	75.1	83.2	85.8	70.9	78.5
LAGOR	Two	91.5	81.2	86.3	86.6	72.0	79.4
Human (U)	360	94.0	90.6	92.3	93.4	88.9	91.2
Human (M)	360	100.0	100.0	100.0	100.0	100.0	100.0

Table 2: **Accuracy on the SNARE benchmark.** LAGOR outperforms off-the-shelf vision and language models, often even compared to those considering a 360 representation of the object, by adding a level of 3D object understanding in the form of a view estimation auxiliary loss. For multi-view models, we maxpool over visual embedding vectors from multiple views. CLIP rows are zero-shot performance without fine-tuning on SNARE, while MATCH rows use SNARE training data. Human accuracy is conservatively calculated as the number of SNARE instances for which the correct referent object was identified by *every* voting annotator (**Unanimous**); for all SNARE instances a majority of voting annotators correctly selected the referent (**Majority**).

244 (Section 5.2). We compare LAGOR to the off-the-shelf CLIP model that considers a single view of
 245 the object only, before performing any view-gathering rotations.

246 In particular, we set two objects on the workspace of a Franka Emika Panda³ with a wrist-mounted
 247 Intel Realsense D414.⁴ We capture an initial view of the scene, then segment the objects from the
 248 workspace using the Unseen Object Clustering [47] algorithm. We feed the segmented objects as the
 249 initial views to LAGOR and CLIP. We move the arm to a second vantage point above the workspace
 250 and capture another image of the objects, repeat segmentation, and treat these segmentations as the
 251 second views for LAGOR.

252 These experiments introduce a number of new sources of noise not present in the simulated setting,
 253 for example changes in camera angle, shifting lighting during rotation, and automatic segmentation
 254 to account for cluttered scenes. In Section 5.2, we include a pair of examples where it is clear how
 255 each of these variables affects the final results. For example, segmentation leads to the creation of
 256 voids inside of objects not present in SNARE.

257 5 Results

258 While CLIP [1] and ViLBERT [6] are competitive baselines for many language and vision tasks, in
 259 this section we show that they select correct 3D object referents in SNARE substantially less often
 260 than the MATCH and LAGOR models. Then, we deploy CLIP and LAGOR on a robot tasked with
 261 selecting objects conditioned on natural language referring expressions, and find that view estimation
 262 and gathering additional object viewpoints improves referent identification.

263 5.1 LAGOR Performance on SNARE

264 Table 2 gives LAGOR accuracy compared to other two-view alternatives, as well as single view
 265 and 360 view models. We begin with comparisons to existing models trained (ViLBERT) and
 266 zero-shot (CLIP) as advertised and recommended in their respective papers. There are three key
 267 takeaways. First, LAGOR outperforms other two-view models, and the performance gap against
 268 MATCH with two views is solely due to the view estimation auxiliary loss. Thus, estimating initial

³<https://www.franka.de/>

⁴<https://www.intelrealsense.com/depth-camera-d415/>

269 object view appears to imbue LAGOR with some level of 3D object understanding that aids in the
 270 interpretation of how language corresponds to 3D object models. Second, LAGOR outperforms
 271 all 360 view models on nearly every metric, indicating that gathering *more* object views is not as
 272 helpful as understanding the 3D *relationship* between views. Third, and most importantly, while
 273 no models achieve human level accuracy, the performance difference is particularly striking on the
 274 blindfolded referring expression subset of data, which lags 5-15% behind visual referring expressions
 275 across models. That gap supports our intuition that blindfolded referring expressions capture a
 276 complementary and challenging linguistic space currently understudied for vision-language models
 277 but key to everyday manipulation. Thereby, SNARE opens interesting avenues for future work
 278 exploring the shape and grasp-points of objects.

279 5.2 Robot Results

280 We run single view, zero-shot CLIP against LAGOR on
 281 8 pairs of objects with language descriptions. CLIP is
 282 able to select the correct object 3/8 times (38% accuracy),
 283 while LAGOR is correct 4/8 times (50% accuracy). Fig-
 284 ure 5 highlights a success and a failure case for LAGOR.
 285 LAGOR succeeds over CLIP on shape and part descrip-
 286 tions, such as *thin metal handle* and *the one with a lid*. A
 287 task at which CLIP succeeds but LAGOR fails is optical
 288 character recognition (OCR); CLIP is known to “read” la-
 289 bels, enabling it to identify *the juice carton*, labeled with
 290 “juice”, while LAGOR’s training on SNARE caused this
 291 ability to be forgotten. Additional examples and discus-
 292 sions are given in Section 7.5.

293 6 Conclusions

294 We introduce ShapeNet Annotated with Referring
 295 Expressions (SNARE), a challenge task to ground lan-
 296 guage to 3D object models. We show that fine-tuned,
 297 massively pretrained vision and language models fall short
 298 of human performance at identifying object referents of natural language expressions by a wide
 299 margin (Table 2). While models improve as more 2D views of 3D objects are available, a robot
 300 agent can only take in one at a time. We introduce Language Grounding through Object Rotation
 301 (LAGOR), a model that performs a rotation operation on a candidate referent object to achieve a more
 302 language-aligned view (Figure 1). We find that LAGOR is able to identify referent objects using
 303 only an initial and final view more accurately than models utilizing 360 views. LAGOR achieves
 304 this result by performing view estimation as an auxiliary loss, enabling learning latent 3D structure.
 305 Finally, we transfer the trained LAGOR model to a physical robot and demonstrate its ability to select
 306 language-aligned object views and better select referent objects after performing informed rotations.

307 In the future, a rotation *policy* could be learned to examine up to N views across M candidate
 308 objects subject to a referring expression, rather than performing a random rotation to a novel view.
 309 By assigning rotation actions expected discriminative values using a trained rotation module, and
 310 penalties based on the time it takes a particular hardware to achieve an object rotation, one could
 311 learn a POMDP policy aggregating object views seen so far to decide whether to next obtain a new
 312 view or make a guess about the referent object. Obtaining different object views could be done either
 313 by picking up and rotating the object or by moving the robot base to see the object from another angle
 314 (as in this paper). Further, since 3D meshes are available for each object, it may be possible to extract
 315 shape-level information using a PointNet [48], similar to experiments attempted by ShapeGlot [33].
 316 Because we target physical robot applications, we may be able to tie language to sets of *graspable*
 317 *points* encoding gripper orientation and position [43] for each object.



Figure 5: (Top) Success: LAGOR correctly identifies the thin handle of the red mug. (Bottom) Failure: LAGOR fails to build a good representation of the more abstract concept “open”.

318 **References**

319 [1] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,
320 P. Mishkin, J. Clark, G. Krueger, and I. Sutskever. Learning transferable visual models from
321 natural language supervision. *arXiv*, 2021.

322 [2] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *Computer*
323 *Vision and Pattern Recognition (CVPR)*, 2016.

324 [3] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale
325 Hierarchical Image Database. In *Computer Vision and Pattern Recognition (CVPR)*, 2009.

326 [4] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and
327 I. Polosukhin. Attention is all you need. In *Neural Information Processing Systems (NeurIPS)*,
328 2017.

329 [5] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. BERT: Pre-training of deep bidirectional
330 transformers for language understanding. In *North American Chapter of the Association for*
331 *Computational Linguistics (NAACL)*, 2019.

332 [6] J. Lu, D. Batra, D. Parikh, and S. Lee. ViLBERT: Pretraining task-agnostic visiolinguistic
333 representations for vision-and-language tasks. In *Neural Information Processing Systems*
334 *(NeurIPS)*, pages 13–23, 2019.

335 [7] Y.-C. Chen, L. Li, L. Yu, A. E. Kholy, F. Ahmed, Z. Gan, Y. Cheng, and J. Liu. UNITER:
336 Universal image-text representation learning. In *European Conference on Computer Vision*
337 *(ECCV)*, 2020.

338 [8] L. H. Li, M. Yatskar, D. Yin, C.-J. Hsieh, and K.-W. Chang. VisualBERT: A simple and
339 performant baseline for vision and language. *arXiv*, 2019.

340 [9] A. Torralba and A. A. Efros. Unbiased look at dataset bias. In *Computer Vision and Pattern*
341 *Recognition (CVPR)*, 2011.

342 [10] T. Kanade and M. Hebert. First-person vision. *Proceedings of the IEEE*, 100(8):2442–2453,
343 2012.

344 [11] P. Anderson, A. Shrivastava, J. Truong, A. Majumdar, D. Parikh, D. Batra, and S. Lee. Sim-
345 to-real transfer for vision-and-language navigation. In *Conference on Robot Learning (CoRL)*,
346 2020.

347 [12] D. Gurari, Q. Li, A. J. Stangl, A. Guo, C. Lin, K. Grauman, J. Luo, and J. P. Bigham. Vizwiz
348 grand challenge: Answering visual questions from blind people. In *Computer Vision and Pattern*
349 *Recognition (CVPR)*, June 2018.

350 [13] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva,
351 S. Song, H. Su, J. Xiao, L. Yi, and F. Yu. ShapeNet: An Information-Rich 3D Model Repository.
352 Technical Report arXiv:1512.03012 [cs.GR], Stanford University — Princeton University —
353 Toyota Technological Institute at Chicago, 2015.

354 [14] S. Tellex, N. Gopalan, H. Kress-Gazit, and C. Matuszek. Robots that use language. *The Annual*
355 *Review of Control, Robotics, and Autonomous Systems*, 15, 2020.

356 [15] K. Takahashi, T. Ogata, J. Nakanishi, G. Cheng, and S. Sugano. Dynamic motion learning
357 for multi-dof flexible-joint robots using active-passive motor babbling through deep learning.
358 *Advanced Robotics*, 31(18):1002–1015, 2017.

359 [16] A. Mousavian, C. Eppner, and D. Fox. 6-dof graspnet: Variational grasp generation for object
360 manipulation. In *International Conference on Computer Vision (ICCV)*, pages 2901–2910,
361 2019.

362 [17] T. Nagarajan and K. Grauman. Learning affordance landscapes for interaction exploration in
363 3D environments. In *Neural Information Processing Systems (NeurIPS)*, 2020.

- 364 [18] J. Thomason, A. Padmakumar, J. Sinapov, N. Walker, Y. Jiang, H. Yedidsion, J. Hart, P. Stone,
365 and R. J. Mooney. Jointly improving parsing and perception for natural language commands
366 through human-robot dialog. *The Journal of Artificial Intelligence Research (JAIR)*, 67, 2020.
- 367 [19] C. Lynch and P. Sermanet. Grounding language in play. In *arXiv*, 2020.
- 368 [20] Y. Bisk, A. Holtzman, J. Thomason, J. Andreas, Y. Bengio, J. Chai, M. Lapata, A. Lazaridou,
369 J. May, A. Nisnevich, et al. Experience grounds language. *arXiv preprint arXiv:2004.10151*,
370 2020.
- 371 [21] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang. Bottom-up and
372 top-down attention for image captioning and visual question answering. In *Computer Vision
373 and Pattern Recognition (CVPR)*, 2018.
- 374 [22] K. He, G. Gkioxari, P. Dollár, and R. B. Girshick. Mask R-CNN. *International Conference on
375 Computer Vision (ICCV)*, 2017.
- 376 [23] H. De Vries, F. Strub, S. Chandar, O. Pietquin, H. Larochelle, and A. Courville. Guess-
377 what?! visual object discovery through multi-modal dialogue. In *Computer Vision and Pattern
378 Recognition (CVPR)*, 2017.
- 379 [24] J. Lu, V. Goswami, M. Rohrbach, D. Parikh, and S. Lee. 12-in-1: Multi-task vision and language
380 representation learning. In *Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- 381 [25] C. W. Leong and R. Mihalcea. Going beyond text: A hybrid image-text approach for measuring
382 word relatedness. In *International Joint Conference on Natural Language Processing (IJCNLP)*,
383 2011.
- 384 [26] V. Cohen, B. Burchfiel, T. Nguyen, N. Gopalan, S. Tellex, and G. Konidaris. Grounding
385 language attributes to objects using bayesian eigenobjects. In *Intelligent Robots and Systems
386 (IROS)*, 2019.
- 387 [27] P. Anderson, Q. Wu, D. Teney, J. Bruce, M. Johnson, N. Sünderhauf, I. Reid, S. Gould, and
388 A. Van Den Hengel. Vision-and-language navigation: Interpreting visually-grounded navigation
389 instructions in real environments. In *Computer Vision and Pattern Recognition (CVPR)*, 2018.
- 390 [28] J. Abramson, A. Ahuja, A. Brussee, F. Carnevale, M. Cassin, S. Clark, A. Dudzik, P. Georgiev,
391 A. Guy, T. Harley, F. Hill, A. Hung, Z. Kenton, J. Landon, T. Lillicrap, K. Mathewson, A. Muldal,
392 A. Santoro, N. Savinov, V. Varma, G. Wayne, N. Wong, C. Yan, and R. Zhu. Imitating interactive
393 intelligence. *arXiv*, 2020.
- 394 [29] M. Shridhar, J. Thomason, D. Gordon, Y. Bisk, W. Han, R. Mottaghi, L. Zettlemoyer, and
395 D. Fox. ALFRED: A Benchmark for Interpreting Grounded Instructions for Everyday Tasks. In
396 *Computer Vision and Pattern Recognition (CVPR)*, 2020.
- 397 [30] A. Majumdar, A. Shrivastava, S. Lee, P. Anderson, D. Parikh, and D. Batra. Improving vision-
398 and-language navigation with image-text pairs from the web. In *European Conference on
399 Computer Vision (ECCV)*, 2020.
- 400 [31] Y. Hong, Q. Wu, Y. Qi, C. Rodriguez-Opazo, and S. Gould. A recurrent vision-and-language
401 BERT for navigation. *arXiv*, 2020.
- 402 [32] K. P. Singh, S. Bhambri, B. Kim, R. Mottaghi, and J. Choi. MOCA: A modular object-centric
403 approach for interactive instruction following. *arXiv*, 2020.
- 404 [33] P. Achlioptas, J. Fan, R. X. Hawkins, N. D. Goodman, and L. J. Guibas. ShapeGlot: Learning
405 language for shape differentiation. In *International Conference on Computer Vision (ICCV)*,
406 2019.
- 407 [34] M. Shridhar, D. Mittal, and D. Hsu. INGRESS: Interactive visual grounding of referring
408 expressions. *The International Journal of Robotics Research (IJRR)*, 39(2-3):217–232, 2020.
- 409 [35] T. Kollar, S. Tellex, M. R. Walter, A. S. Huang, A. Bachrach, S. Hemachandra, E. Brunskill,
410 A. Banerjee, D. Roy, S. Teller, and N. Roy. Generalized grounding graphs: A probabilistic
411 framework for understanding grounded commands. *arXiv*, 2017.

- 412 [36] J. Thomason, J. Sinapov, M. Svetlik, P. Stone, and R. J. Mooney. Learning multi-modal
413 grounded linguistic semantics by playing “I spy”. In *International Joint Conference on Artificial*
414 *Intelligence (IJCAI)*, 2016.
- 415 [37] B. Calli, A. Singh, J. Bruce, A. Walsman, K. Konolige, S. Srinivasa, P. Abbeel, and A. M.
416 Dollar. Yale-cmu-berkeley dataset for robotic manipulation research. *The International Journal*
417 *of Robotics Research (IJRR)*, 36(3):261–268, 2017.
- 418 [38] R. Scalise, J. Thomason, Y. Bisk, and S. Srinivasa. Improving robot success detection using
419 static object data. In *Intelligent Robots and Systems (IROS)*, 2019.
- 420 [39] X. Zhang, J. Sinapov, and S. Zhang. Planning multimodal exploratory actions for online robot
421 attribute learning. In *Robotics: Science and Systems (RSS)*, 2021.
- 422 [40] S. Nair, S. Savarese, and C. Finn. Goal-aware prediction: Learning to model what matters. In
423 *International Conference on Machine Learning (ICML)*, 2020.
- 424 [41] C. Paxton, Y. Bisk, J. Thomason, A. Byravan, and D. Fox. Prospection: Interpretable plans from
425 language by predicting the future. In *International Conference on Robotics and Automation*
426 *(ICRA)*. IEEE, 2019.
- 427 [42] J. Y. Koh, H. Lee, Y. Yang, J. Baldridge, and P. Anderson. Pathdreamer: A world model for
428 indoor navigation. *arXiv*, 2021.
- 429 [43] C. Eppner, A. Mousavian, and D. Fox. ACRONYM: A large-scale grasp dataset based on
430 simulation. In *2021 IEEE Int. Conf. on Robotics and Automation, ICRA*, 2020.
- 431 [44] A. X. Chang, M. Savva, and P. Hanrahan. Semantically-enriched 3d models for common-sense
432 knowledge. In *FPIC Workshop, Computer Vision and Pattern Recognition (CVPR)*, 2015.
- 433 [45] G. A. Miller. Wordnet: A lexical database for english. *Communications of the ACM*, 38(11):
434 39–41, 1995.
- 435 [46] P. Zhang, X. Li, X. Hu, J. Yang, L. Zhang, L. Wang, Y. Choi, and J. Gao. VinVL: Making visual
436 representations matter in vision-language models. *Computer Vision and Pattern Recognition*
437 *(CVPR)*, 2021.
- 438 [47] Y. Xiang, C. Xie, A. Mousavian, and D. Fox. Learning rgb-d feature embeddings for unseen
439 object instance segmentation. In *Conference on Robot Learning (CoRL)*, 2020.
- 440 [48] C. R. Qi, L. Yi, H. Su, and L. J. Guibas. Pointnet++: Deep hierarchical feature learning on
441 point sets in a metric space. In *Neural Information Processing Systems (NeurIPS)*, 2017.
- 442 [49] J. Pennington, R. Socher, and C. Manning. Glove: Global vectors for word representation. In
443 *Empirical Methods in Natural Language Processing (EMNLP)*, 2014.
- 444 [50] P. Sharma, N. Ding, S. Goodman, and R. Soricut. Conceptual captions: A cleaned, hypernymed,
445 image alt-text dataset for automatic image captioning. In *Association for Computational*
446 *Linguistics (ACL)*, 2018.
- 447 [51] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. BERT: Pre-training of deep bidirectional
448 transformers for language understanding. In *North American Chapter of the Association for*
449 *Computational Linguistics (NAACL)*, 2019.