

# Building Language-Agnostic Grounded Language Learning Systems

Caroline Kery, Nisha Pillai, Cynthia Matuszek, and Francis Ferraro<sup>1</sup>

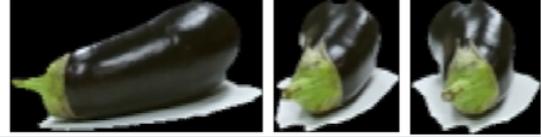
**Abstract**—Learning the meaning of *grounded* language—language that references a robot’s physical environment and perceptual data—is an important and increasingly widely studied problem in robotics and human-robot interaction. However, with a few exceptions, research in robotics has focused on learning groundings for a single natural language pertaining to rich perceptual data. We present experiments on taking an existing natural language grounding system designed for English and applying it to a novel multilingual corpus of descriptions of objects paired with RGB-D perceptual data. We demonstrate that this specific approach transfers well to different languages, but also present possible design constraints to consider for grounded language learning systems intended for robots that will function in a variety of linguistic settings.

## I. INTRODUCTION

As robots become less expensive, and more capable, it is becoming possible to imagine them being deployed in a variety of human-centric settings such as homes, schools, or workplaces. However, as robots become more accessible, it becomes more critical that they can be communicated with and controlled by non-specialists. Robotic assistants could be extremely helpful for a variety of people such as the elderly or the disabled, so for maximal accessibility it would be ideal to make the use of such robots as intuitive as possible. *Grounded language acquisition*, in which robots learn to understand language in the context of the sensed physical world around them, is a major focus of research for building robot systems that can interact and collaborate with human partners in a natural way. At the same time, studying language learning and interaction with a physically situated agent offers a mechanism for advancing natural language understanding [23], [12].

Although learning shared embeddings between sensor data and language is a rich and varied field of research in robotics, most of the work has focused on systems that operate on only a single language, frequently English. In the natural language processing community, there is work on building *multilingual* systems, which learn from [10], [13] or can be used with [34], [6] more than one language simultaneously. In this work, we address a complementary problem: Can systems that take advantage of physical percepts to learn meanings in a particular language be deployed in a setting where a different language is used?

In this work, we describe the application of an existing grounded language learning system that uses the words-as-classifiers model [21], [28], [24] to a novel corpus containing



This is an Italian eggplant. It is firm and dark purple when ripe.

Esta es una berenjena. La Berenjena se utiliza para preparar deliciosos platos.

यह एक बैंगन है। यह एक सब्जी का प्रकार है। इसका बना भर्ता भी स्वादिष्ट होता है।

Fig. 1. An example of data collected for grounded language learning. *Top*: Several images of an object in the dataset. *Below*: Descriptions of the object in unconstrained English, Spanish, and Hindi (respectively).

two additional languages. We discuss methods of acquiring suitable training data, the overall performance of the system on English, Spanish, and Hindi, and some characteristics of a grounded language process that should be considered when linguistic flexibility is desired. In practice, such flexibility should be considered highly desirable in order to avoid limiting the benefits of ubiquitous, collaborative robots to English-speaking settings.

The remainder of this work is organized as follows: In Section III-A, we describe our novel trilingual dataset, which is comprised of approximately 17,000 descriptions of objects collected with an RGB-D sensor. Each object has descriptions in unconstrained English, Spanish and Hindi. In Section III-B and Section III-C, we describe applying a grounded language acquisition system designed for English [25] to this corpus, including brief descriptions of the natural language processing tools used. Finally in Section IV, we demonstrate the need for human-provided training data, analyze specific sources of performance degradation across language pairs, and offer design suggestions for perceptually grounded language learning systems that are intended to be agnostic about the language being learned.

To the best of our knowledge, this is the first work that seeks to apply a grounded language learning system designed for one language to other languages in order to evaluate its transferability. Our main contributions are: This evaluation, along with a detailed analysis of sources of error; a dataset containing images and trilingual descriptions of those objects; and a set of suggestions and considerations for future research in this area. We will make our dataset available upon publication.

\*This material is based in part upon work supported by the National Science Foundation under Grant No. 1657469.

<sup>1</sup>All authors are with the University of Maryland, Baltimore County, 1000 Hilltop Circle, Baltimore, MD, 21250, USA.

kery1|npillai1|cmat|ferraro@umbc.edu

## II. RELATED WORK

Grounding natural language to the physical world [11] is highly relevant to robotics research, and researchers are addressing the question from a variety of perspectives. There have been a number of successful approaches in different areas, such as navigation [31], [20], understanding commands and directions [2], [1] or action words [5], grounding spatial relations and concepts [24], and referring expressions for objects in images [8], [36]. Other work has explored interactively grounding additional non-visual properties like sound and weight [32].

Our experiments focus on the problem of learning people’s preferred labels for color, shape, and semantic object labels from unconstrained descriptions of real-world objects. Understanding such attributes is critical for tasks such as grasping [18] and collaboration [14]. In this work, rather than learning connections between novel language and an existing formulation, we follow the approach of [21], extending the formal representation as novel words are encountered.

Many grounded learning systems associate words with the perceived world [30], [19]. In this paper, we demonstrate the importance of selecting natural language processing techniques carefully when improving the efficiency of such grounded learning systems. Language pre-processing techniques have been studied when measuring document clustering and retrieval [15], [3], but not in visual-linguistic grounding. This paper further examines visual classification tasks with lemmatized and/or stemmed tokens of various categories and proposes approaches to improve learning.

Some systems that ground language attributes apply complex, structured linguistic analysis techniques to utterances in order to examine how particular phrases might relate to other parts of the sentences [24], [21]. Applying these structured linguistic approaches to an arbitrary new language is a complex undertaking, sometimes requiring costly expert annotation. In this paper, we consider each token in isolation and do not use any additional language features from the descriptions when learning the meaning of the tokens, allowing the model to scale more easily across languages.

Most research has focused on learning English-based robotic language acquisition systems, although there are some non-English corpora that do exist that could potentially be used (e.g., [7], [33], among others). In contrast, we concentrate on developing language grounding models from descriptions of objects gathered via a low-cost robotic sensor. In our experiments, we extended the monolingual public image dataset of Pillai and Matuszek [25] by collecting (non-paired) Hindi and Spanish language descriptions from native speakers.

Frank *et al.* [9] evaluate the preference of descriptions generated by native language speakers over the descriptions translated from a different language, whereas we evaluate how well a system trained on either translated data or Amazon Mechanical Turk language descriptions collected from native speakers could then complete the object recognition task using AMT descriptions.

Our work is most similar to that of [13] and [10], who improve the performance of learning visual-semantic embeddings by training over multiple languages. However, our work differs in two important ways. First, rather than using descriptions of perceptual data from several languages matched to one another, we explore how a core learning architecture designed for a particular language transfers to a new language; we do not assume that *parallel* (aligned) multilingual descriptions are available. Second, we collect unconstrained language about particular objects as seen by a real sensor, rather than relying on caption data. This is consistent with the goal of deploying robots for use by non-specialists in different environments.

## III. APPROACH

There is a substantial body of research on learning mappings between a natural language and robot-usable representations. We follow [28] and [25] in using a *words-as-classifiers* model, in which a robot learns a mapping from each word in a language to classifiers trained on perceptual data over objects or actions. This process consists of three stages. First, perceptual data and descriptions of objects are collected in multiple languages, simulating inputs a robot might receive when being taught about a new object. Second, the language and percepts are pre-processed to extract features, and relevant language tokens (words) are extracted and paired with the features of the objects they describe, along with suitable negative examples [25]. Third, for each term, several binary classifiers are trained and evaluated against a held-out set of objects. Fig. 3 shows this architecture.

### A. Data Corpus

The system was trained and tested using a novel corpus consisting of RGB-D images paired with multilingual natural language descriptions of the objects in the images. This corpus extends the public dataset of [25], tripling the number of language descriptions. The dataset itself has 18 object categories (such as *carrot* and *arch*), with four instances (individual objects) in each category and an average of five images of each instance (see Fig. 2 for examples). Each image was taken from various angles with a Kinect2 depth sensor, yielding both RGB-D point clouds and regular color images. The goal in using these images was to collect images

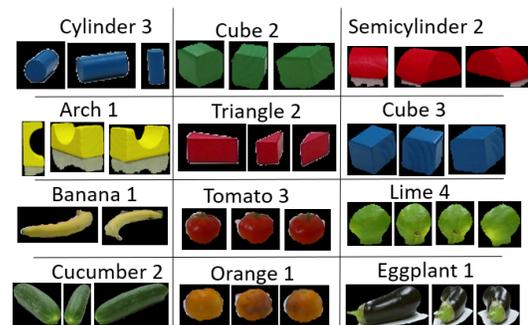


Fig. 2. Cropped Kinect2 images of twelve objects. All categories in the data set were toy blocks, fruit, or vegetables of various colors and shapes.

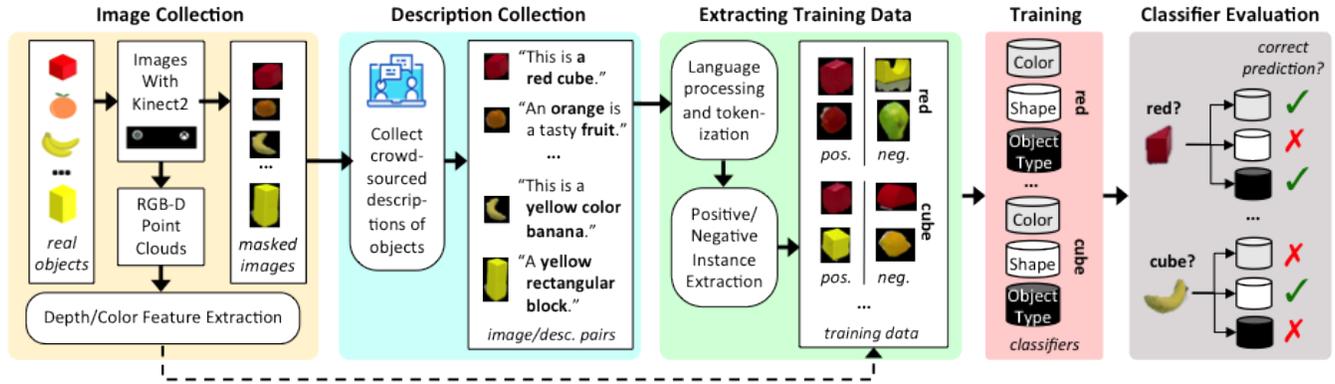


Fig. 3. This diagram shows the data flow of the grounded language acquisition system, from sensor and language data collection, through selection of training data and learning of classifiers, to eventual testing of those classifiers. Individual sections indicate the elements of the approach, as detailed in Section III. Note that while this diagram only shows classifiers and examples for “red” and “cube,” classifiers are learned for all tokens.

that accurately represent how a robot might perceive objects in its environment. The RGB images were masked so only the objects were visible to participants.

In existing and previous work, the English natural language descriptions of each instance were obtained using the Amazon Mechanical Turk crowd-sourcing platform. We chose to mimic this setup when collecting fluent speaker descriptions in Hindi and Spanish. These languages were chosen for their high number of speakers, as well as their varying dissimilarities to English. Our motivation for collecting new language data was to determine whether machine-translated training data is sufficient to handle language use by native speakers, or if it is necessary to involve speakers of that language to train an adequate model. As we discuss later in this paper, we found that the translated data was not sufficient for learning the language used by native speakers, resulting in an approximately 20% reduction in  $F_1$  score across all languages (see Fig. 5).

Workers were presented with cropped and masked images of objects and asked to provide descriptions. We chose to provide the workers with no sample descriptions, in order to maximize the variation in the descriptions. The purpose of this was to evaluate our model with data that could accurately represent the ways in which speakers of each language might talk about objects when presenting them to a robot. Despite the lack of priming, many workers in all three languages chose to describe most objects with simple descriptions like “This is a red cube.” However, we also saw noticeable variation in all three languages, where workers chose to provide extra information about objects, such as describing cucumbers as healthy or being good in a salad.

The dataset contains 6,045 descriptions in English, 5,735 descriptions in Hindi, and 5,104 in Spanish. Originally, over 6,000 descriptions were collected for each of Hindi and Spanish, but in both cases we excluded descriptions from workers who explicitly did not follow the directions, such as by responding in the wrong language or consistently responding with text unrelated to the images they were being asked to describe; we did not exclude for other reasons. In

the analysis section, we account for the smaller number of descriptions collected in Spanish and Hindi by randomly subsetting all datasets in such a way that each instance is trained on an equal amount of descriptions in each language. The results are averaged from several of such subsets.

### B. Semantic Processing of Descriptions

Once collected, descriptions were put through a series of preprocessing steps to extract relevant tokens. We defined a relevant token as a word the robot might want to learn to recognize. For example, in the description “This object is a large yellow banana!”, a preprocessed version might be “object large yellow banana” (or “large yellow banana”). We conducted initial experiments to determine methods that extracted the most appropriate tokens in the three languages.<sup>1</sup>

1) *Stemming and Lemmatizing Tokens:* The overall design of this grounded language learning system hinges upon gaining an understanding about the meaning of a new object or attribute descriptor by examining the objects it was used to describe. Given limited training data, the system should be able to recognize when the same word is being used across different examples. In English, nouns have only a singular or plural form, and adjectives are rarely conjugated (with the exception of comparative terms like “larger”). In contrast, in both Spanish and Hindi nouns can be gendered and adjectives are often conjugated to match the gender of the nouns. For example, in Spanish, “The red table” is “La mesa roja,” where “roj” is the stem of the adjective “red,” and “roja” is the feminine singular form. A robot may therefore unnecessarily learn both a masculine and feminine word describing the same concept, subdividing the training data and weakening the resulting classification.

Past work applied a *lemmatizer* to the English data to remove conjugation from words. Lemmatizers take words that have been conjugated in some way and reduce them to

<sup>1</sup>Following standard language approaches, we lower-cased all words (which was only a valid step for Spanish and English) and removed punctuation, including language-specific punctuation like the Spanish upside-down question mark and the Hindi full-stop.

their unconjugated forms—for example, reducing “running” to “run.” When examining this step in the context of Hindi and Spanish, due to a lack of readily available non-English lemmatizers, we replaced lemmatization with the simpler but related step of *stemming*. Stemming is very similar to lemmatization in that both tools take conjugated words and remove the conjugation. The difference is that stemmers simply chop off conjugations instead of attempting to find the original unstemmed form. This simplification makes them potentially less effective (e.g., a stemmer might reduce “running” to “runn” instead of “run”), but also much easier to implement. For this paper, we use the NLTK Snowball Stemmer [26] for Spanish and English text while For Hindi we used the simple stemmer described in [27].

2) *Stop Word Identification and Removal*: An additional challenge in using unconstrained natural language is identifying what words the system should attempt to learn physical meanings for and what words it should ignore. In natural language processing, it is common practice to remove “stop words” from language data as a preprocessing step. Stop words are defined as words that are necessary to form grammatical sentences, but do not contribute to the overall meaning of the sentence. We considered two kinds of stop words, general and domain-specific. There are publicly available lists of general stop words, such as “the” or “and.” For English and Spanish, we used NLTK’s stop word lists in each language. Since NLTK did not have a list for Hindi, we instead utilized the stop word list compiled in [29].

In our data, terms such as “object” or “item,” were broadly used as generic terms to describe most elements in the dataset. We identified these domain-specific stop words by their Inverse Document Frequency (IDF), an effective method for information retrieval in which each word was scored according to the log of the number of instances in the dataset divided by the total number of instances where that word appeared in some description. Words with low IDF scores appeared in many different descriptions and were removed. We discuss this further in later sections.

### C. Model Learning

This section discusses the grounded language learning system’s components, including feature extraction and identification of “ground truth” positive and negative examples.

1) *Tokenization and Positive/Negative Example Identification*: After the descriptions were preprocessed, the next step was to extract relevant tokens and identify positive and negative examples of these tokens. As the only human guidance we have is the collection of unprompted and unguided descriptions (as opposed to annotations over pre-defined labels), the “ground truth” for what instances were and were not representative of various concepts had to be extracted from the descriptions.

To identify relevant tokens, we concatenated all descriptions for each instance (approximately 80 descriptions per instance) and counted how many times each unique word was used. If a word was used more than five times for a particular instance, the instance was deemed a “positive

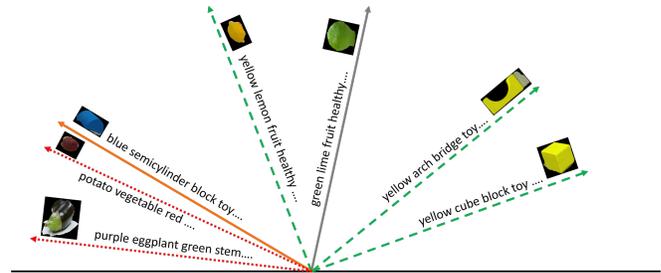


Fig. 4. A visualization of negative example selection for the token “yellow.” Note that the negative examples selected (shown with red dotted arrows) were the ones that were the farthest away from the most positive instances (shown with green dashed arrows).

example” of that word. Words that did not have at least three positive examples were excluded from consideration, as it was deemed that the robot had not seen this concept enough times to learn it. By extracting the tokens in this way, the system ignored any context in which the token was used outside of the image it was paired with.

Once positive examples were found for each token, our system also needed negative examples. Negative examples are rarely given in descriptions: users are much more likely to describe an object with positive properties, rather than enumerate negative properties. To address this, we used the approach of [25]. We took the combined descriptions of each object instance and represented these descriptions in vector space using the Distributed Memory Model of Paragraph Vectors (PV-DM) [22], [17]. We used the PV-DM model to measure dissimilarity between descriptions; in this model, semantically similar documents are similar in vector space. We used the cosine similarity statistic measure to find the dissimilar vectors and used those objects as negative examples in our language learning model. The intuition behind this approach is that instances that are negative examples of other instances are likely to have been described using very different language.

One consideration when choosing negative instances was that a token may have positive instances that are dissimilar to each other (see Fig. 4). We minimize this risk by choosing negative examples that are the farthest away in space from the most positive instances. This was decided using a weighted vote. For a term, all objects that had ever been described using that term were discarded. Remaining objects were sorted by the cosine similarity of their descriptions in vector space, and the last 2/3rds of the list were retained. Each positive instance then returned this set of candidates, weighted by similarity. The final scores for each negative instance candidate came from the sum of its weighted vote. The candidates were sorted by this score, and the top 25% of the candidates were selected as negative examples. This value was chosen experimentally. Tokens with no identified negative examples were removed from consideration.

2) *Extracting Image Features*: For the classification task, two kinds of features were extracted from the RGB-D images. For generating color features, the RGB images of

the objects were used. The RGB features of all pixels in the images were clustered using  $k$ -means based on RGB values. Then histograms were created based on the density of all clusters. The ‘R G B’ values of the centroid with the highest density were selected as the color features. We used HMP-extracted kernel descriptors [16], [4] to take the location and depth data for each image and extract shape features.

3) *Training and Testing Classifiers*: For each remaining token, three binary classifiers were trained using logistic regression on positive instances (objects described using that token) and negative instances selected as above. The underlying idea behind training all three types of classifiers per token was that a new word might describe a color, shape, or object and a robot with no previous knowledge would not know which category the token should belong to. Color classifiers were trained using RGB features, shape classifiers were trained using HMP features, and object classifiers were trained using a combination of the above. During evaluation, each classifier attempted to classify held-out instances.

#### IV. EXPERIMENTAL RESULTS

In this work, we expand the language learning system previously tested with English data to additional languages. To assess the quality of the language groundings, we applied the learned token classifiers to an object selection task. One instance of every object category and its descriptions were held out during training for testing. For each token used in the descriptions of these instances, positive and negative examples were found from the testing instances, as discussed above. The token classifiers learned in the training phase were then scored by how accurately they could identify which instances were positive and which were negative examples of the token when presented with images of those instances. The choices of which instances to test on and which image of each instance to present were randomized, and the scores found for each token were the average from choosing ten times. The entire evaluation setup was run twenty times for each time the model was trained on a different train-test split. The model was trained nine times for each result reported.

In addition, as each token had a separate color, shape, and object classifier, three F1-scores were reported for each token. The final F1-scores presented here average the F1-score across all test tokens for each category. We note that we did not do any manual selection of which words would normally belong in each category, thereby depressing the overall scores. Tokens that appeared in testing but not training data received an F1-Score of 0.

Next we discuss the performance of the system, and how it was impacted by both the origins of the training data and the language processing steps applied to it.

##### A. Performance Across Data Sources and Processing Techniques

1) *The Utility of Translated Data*: For a robotic system to learn language data in new languages, it would be convenient if one could train such a system on translations of the

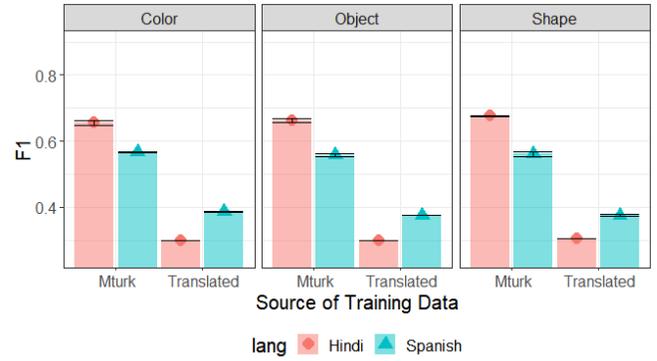


Fig. 5. System performance as compared between training the model on a translated version of the English dataset and training it on using descriptions collected in the target language using crowd-sourcing. Note that the translated dataset was not sufficient for learning many of the tokens used by native speakers when describing the images.

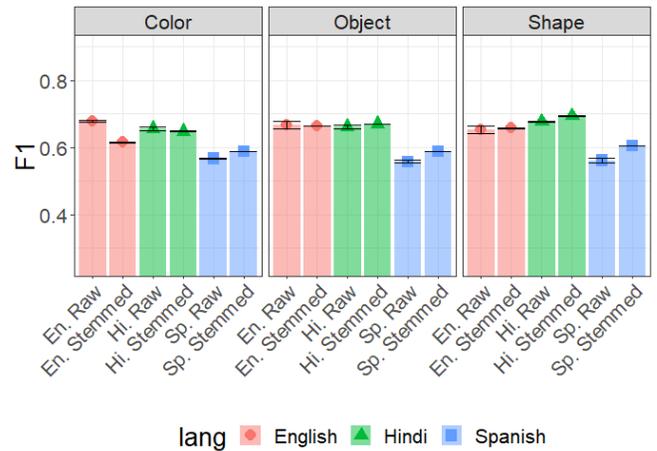


Fig. 6. Overall scores of all three languages. This compares the performance over the un-stemmed “raw” datasets, and the stemmed versions. Note that to allow for fair comparison, the Hindi, English, and Spanish datasets have been subset to have equal amounts of descriptions per instance. The error bars provide the variance in the scores.

language data already collected. We explored this possibility in Fig. 5, which shows the performance of the system when it was trained using either a translated version of the English dataset (using Google Translate’s API [35]), or the crowd-sourced descriptions. In both cases, the trained models were tested on the tokens and positive/negative instances identified from the crowd-sourced descriptions, as these most accurately represent how a native speaker might describe the objects to a robot. The figure shows that the classifiers trained using the translated corpora were not sufficient when faced with the native language data. A major contribution to the lower performance was the large number of tokens that were used in the Mechanical Turk descriptions, and not in the translated corpora. Intuitively, a direct translation of a corpus in English was unable to accurately represent the wide variety of ways that a native speaker in the target language might describe the same objects.

2) *Results Across Languages:* Fig. 6 shows the average F1-scores of the model across the three languages when trained and tested on the Mechanical Turk descriptions. The model performed comparably across the three languages. Average scores were somewhat low due to the large number of previously unseen words occurring in the test data; this was especially the case for Spanish, with an average of fifty previously unseen tokens per testing run (where English and Hindi each averaged approximately thirty). This is primarily a product of our relatively small dataset; larger initial training data collected from Mechanical Turk would likely improve performance by ensuring that more total descriptions occurred in the test split.

3) *Language Processing Considerations:* When expanding the grounded language system to Hindi and Spanish, some consideration had to be made of the language processing techniques chosen in order to ensure that relevant tokens were identified and correctly conflated. Stemming had a relatively low impact on the scores. Qualitative analysis of the results showed that stemming did enable the system to correctly conflate different gendered forms of adjectives in Hindi and Spanish. Fig. 7 shows the effect of varying the IDF score threshold for removal, showing that the optimal threshold varied by language. For Spanish, several important color words were used often enough that they appeared in the bottom 2% of tokens by IDF score, while for Hindi, the bottom 3% of terms were safe to remove. These are comparatively small differences, but suggest that although it can be beneficial to remove unnecessary words, care must be taken when defining what tokens to remove.

### B. Design Criteria for Language-Agnostic Learning

In the course of this work, we have presented a comparative analysis of the performance of a particular grounded language learning method when applied to two novel languages. Beyond this detailed analysis, an additional contribution of this work is suggested design criteria that groups designing grounded language learning systems may wish to consider.

First, consistent with expectations, we find that the more sophisticated a semantic processing step is, the less likely it is to work in a new language without significant modifications. For example, the simpler approach of stemming is more accessible than lemmatization for the novel languages considered. A robotics audience may then wish to explore learning methods that do not rely heavily on natural language preprocessing, or to focus on methods that are themselves relatively language-agnostic.

Second, training data beyond that provided by end-users should be provided by fluent speakers of the new language whenever possible. Fig. 6 shows that for our test case, using training data drawn from a simple machine translation approach was inadequate to support classifier performance when tested with fluent speakers, who tend to take wide advantage of the rich variety of terms and idioms available.

Finally, in general, the words-as-classifiers approach is well suited to transitioning across languages. The only modifications required were in identifying meaningful words in

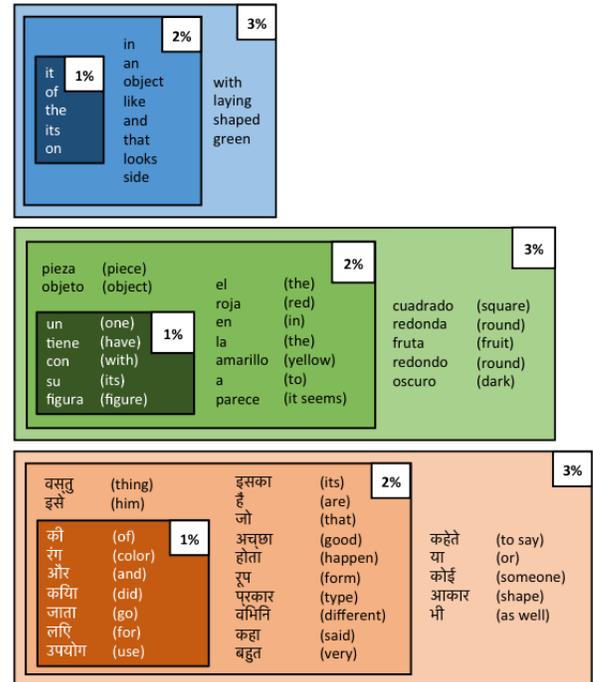


Fig. 7. Stop words as selected by IDF, showing the 1–3% of tokens occurring in the most documents, in each of English (top), Spanish (middle), and Hindi (bottom). Color words incorrectly occur in the top 2% for Spanish and 3% for English; in Hindi no color or shape words appeared in the top 3%. This suggests that the percentage of words removed may need to be tuned by language.

the data. This reduces possible complications that would be introduced by more sophisticated NLP techniques. For example, disregarding the relative placement of words in the descriptions meant that the model was unaffected by the fact that Hindi has much looser rules for word ordering than English or Spanish.

## V. DISCUSSION AND FUTURE WORK

Grounding natural language in perception is an essential task in human robot interaction. In this paper, we have taken an existing grounded language learning system and demonstrated that it can be easily extended to handle data in new languages. We found that when designing a system with unconstrained language data and noisy perceptual data, it is important to minimize and simplify the natural language processing preprocessing steps. We collected two new language corpora in Spanish and Hindi, and demonstrated that translated data was not sufficient for training. We will make our novel corpus available upon publication. In the future, we will implement the modified system or a similar learning system on a mobile robot in collaborative setting, where it must learn from and then interact with people using one or more novel languages.

## REFERENCES

- [1] P. Anderson, Q. Wu, D. Teney, J. Bruce, M. Johnson, N. Sünderhauf, I. Reid, S. Gould, and A. van den Hengel, “Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real

- Caroline Kery, Nisha Pillai, Cynthia Matuszek, Francis Ferraro. "Building Language-Agnostic Grounded Language Learning Systems." In *Proceedings of the 28<sup>th</sup> Int'l. Conference on Robot and Human Interactive Communication (Ro-Man)*, 2019. Submitted draft.
- environments," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, vol. 2, 2018.
- [2] Y. Artzi and L. Zettlemoyer, "Weakly supervised learning of semantic parsers for mapping instructions to actions." *Transactions of the Association for Computational Linguistics (TACL)*, vol. 1, pp. 49–62, 2013.
- [3] V. Balakrishnan and E. Lloyd-Yemoh, "Stemming and lemmatization: A comparison of retrieval performances," *Lecture Notes on Software Engineering*, vol. 2, no. 3, 2014.
- [4] L. Bo, K. Lai, X. Ren, and D. Fox, "Object recognition with hierarchical kernel descriptors," in *Computer Vision and Pattern Recognition*, 2011.
- [5] J. Y. Chai, Q. Gao, L. She, S. Yang, S. Saba-Sadiya, and G. Xu, "Language to action: Towards interactive task learning with physical agents." in *IJCAI*, 2018, pp. 2–9.
- [6] D. L. Chen, J. Kim, and R. J. Mooney, "Training a multilingual sportscaster: Using perceptual context to learn language," *Journal of Artificial Intelligence Research*, vol. 37, pp. 397–435, 01 2010.
- [7] D. Elliott, S. Frank, K. Sima'an, and L. Specia, "Multi30k: Multilingual english-german image descriptions," in *Proceedings of the 5th Workshop on Vision and Language*, 2016, pp. 70–74.
- [8] N. FitzGerald, Y. Artzi, and L. S. Zettlemoyer, "Learning distributions over logical forms for referring expression generation." in *EMNLP*, 2013, pp. 1914–1925.
- [9] S. Frank, D. Elliott, and L. Specia, "Assessing multilingual multimodal image description: Studies of native speaker preferences and translator choices," *Natural Language Engineering*, vol. 24, no. 3, pp. 393–413, 2018.
- [10] S. Gella, R. Sennrich, F. Keller, and M. Lapata, "Image pivoting for learning multilingual multimodal representations," in *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. Copenhagen, Denmark: Association for Computational Linguistics, Sept. 2017, pp. 2839–2845. [Online]. Available: <https://www.aclweb.org/anthology/D17-1303>
- [11] S. Harnad, "The symbol grounding problem," *Physica D: Nonlinear Phenomena*, vol. 42, no. 1-3, pp. 335–346, 1990.
- [12] T.-H. K. Huang, F. Ferraro, N. Mostafazadeh, M. Misra, A. Agrawal, J. Devlin, R. Girshick, X. He, D. Batra, C. L. Zitnick, D. Parikh, L. Vanderwende, M. Galley, and M. Mitchell, "Visual storytelling," in *NAACL*, 2016, equal contribution: TH, FF.
- [13] Á. Kádár, D. Elliott, M.-A. Côté, G. Chrupała, and A. Alishahi, "Lessons learned in multilingual grounded language learning," in *Proceedings of the 22nd Conference on Computational Natural Language Learning*. Brussels, Belgium: Association for Computational Linguistics, Oct. 2018, pp. 402–412. [Online]. Available: <https://www.aclweb.org/anthology/K18-1039>
- [14] R. A. Knepper, S. Tellex, A. Li, N. Roy, and D. Rus, "Recovering from failure by asking for help," *Autonomous Robots*, vol. 39, no. 3, pp. 347–362, Oct 2015.
- [15] T. Korenius, J. Laurikkala, K. Järvelin, and M. Juhola, "Stemming and lemmatization in the clustering of finnish text documents," in *Proceedings of the Thirteenth ACM International Conference on Information and Knowledge Management*, ser. CIKM '04. New York, NY, USA: ACM, 2004, pp. 625–633. [Online]. Available: <http://doi.acm.org/10.1145/1031171.1031285>
- [16] K. Lai, L. Bo, X. Ren, and D. Fox, "Rgb-d object recognition: Features, algorithms, and a large scale benchmark," in *Consumer Depth Cameras for Computer Vision: Research Topics and Applications*, 2013, pp. 167–192.
- [17] Q. Le and T. Mikolov, "Distributed representations of sentences and documents," in *International conference on machine learning*, 2014, pp. 1188–1196.
- [18] S. Levine, P. Pastor, A. Krizhevsky, J. Ibarz, and D. Quillen, "Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection," *The International Journal of Robotics Research*, vol. 37, no. 4-5, pp. 421–436, 2018.
- [19] C. Matuszek, L. Bo, L. Zettlemoyer, and D. Fox, "Learning from unscripted deictic gesture and language for human-robot interactions," in *Proceedings of the 28<sup>th</sup> AAAI Conference on Artificial Intelligence (AAAI)*, 2014.
- [20] C. Matuszek, N. FitzGerald, E. Herbst, D. Fox, and L. Zettlemoyer, "Interactive learning and its role in pervasive robotics," in *ICRA Workshop on The Future of HRI*, St. Paul, MN, 2012.
- [21] C. Matuszek, N. FitzGerald, L. Zettlemoyer, L. Bo, and D. Fox, "A joint model of language and perception for grounded attribute learning," in *Proceedings of the 2012 International Conference on Machine Learning*, Edinburgh, Scotland, June 2012.
- [22] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *Advances in neural information processing systems*, 2013, pp. 3111–3119.
- [23] R. Mooney, "Learning to connect language and perception," in *Proceedings of the 23rd AAAI Conference on Artificial Intelligence (AAAI)*, Chicago, IL, 2008, pp. 1598–1601.
- [24] R. Paul, J. Arkin, N. Roy, and T. M Howard, "Efficient grounding of abstract spatial concepts for natural language interaction with robot manipulators," in *Proceedings of Robotics: Science and Systems (R:SS) 2016*. Robotics: Science and Systems (RSS), 2016.
- [25] N. Pillai and C. Matuszek, "Unsupervised selection of negative examples for grounded language learning," in *Proceedings of the 32nd National Conference on Artificial Intelligence (AAAI)*, New Orleans, USA, 2018.
- [26] M. F. Porter, "Snowball: A language for stemming algorithms," *Retrieved March*, vol. 1, 01 2001.
- [27] A. Ramanathan and D. D. Rao, "A lightweight stemmer for hindi," in *The Proceedings of EACL*, 2003.
- [28] D. Schlagen, S. Zarrieß, and C. Kennington, "Resolving references to objects in photographs using the words-as-classifiers model," in *Proc. of the 54<sup>th</sup> Annual Meeting of the Association for Computational Linguistics (ACL 2016)*, 2016.
- [29] S. Siddiqi and A. Sharan, "Construction of a generic stopwords list for hindi language without corpus statistics," *International Journal of Advanced Computer Research*, vol. 8, no. 34, pp. 35–40, 2018.
- [30] S. Tellex, P. Thaker, J. Joseph, and N. Roy, "Learning perceptually grounded word meanings from unaligned parallel data," *Machine Learning*, vol. 94, no. 2, pp. 151–167, 2014.
- [31] J. Thomason, S. Zhang, R. J. Mooney, and P. Stone, "Learning to interpret natural language commands through human-robot dialog," in *Twenty-Fourth International Joint Conference on Artificial Intelligence*, 2015.
- [32] J. D. Thomason *et al.*, "Continually improving grounded natural language understanding through human-robot dialog," Ph.D. dissertation, University of Texas at Austin, 2018.
- [33] E. van Miltenburg, Á. Kádár, R. Koolen, and E. Krahmer, "DIDEC: The Dutch image description and eye-tracking corpus," in *Proceedings of the 27th International Conference on Computational Linguistics*. Santa Fe, New Mexico, USA: Association for Computational Linguistics, Aug. 2018, pp. 3658–3669. [Online]. Available: <https://www.aclweb.org/anthology/C18-1310>
- [34] G. Wilcock and K. Jokinen, "Multilingual wiktalk: Wikipedia-based talking robots that switch languages," in *Proceedings of the 16th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, 2015, pp. 162–164.
- [35] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey, *et al.*, "Google's neural machine translation system: Bridging the gap between human and machine translation," in *CoRR*, 2016.
- [36] L. Yu, H. Tan, M. Bansal, and T. L. Berg, "A joint speaker-listener-reinforcer model for referring expressions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 7282–7290.