

Interpreting Multimodal Referring Expressions in Real Time

Miles Eldon¹ and Stefanie Tellex¹

Abstract—Robots that collaborate with humans must be able to identify objects used for shared tasks, for example tools such as a knife for assistance at cooking, or parts such as a screw on a factory floor. Existing work has addressed this problem in single modalities, such as natural language or gesture, but a gap remains in creating real-time multimodal systems that simultaneously fuse information from language and gesture in a principled mathematical framework. We define a multimodal Bayes’ filter for interpreting referring expressions to object using language and gesture in real time. We collected a new RGB-D and audio dataset of people referring to objects in a tabletop setting and demonstrate that our approach successfully integrates information from language and gesture in real time to quickly and accurately identify objects continuously.

I. INTRODUCTION

In order for humans and robots to collaborate in complex tasks, robots must be able to understand people’s references to objects in the external world. To refer to objects, people use a combination of language, gesture, and body language such as eye gaze and looking. These signals are provided continuously to the robot, and a person’s reference can quickly change based on new information about the domain. For example, Figure 1 shows a robot handing a tool to a human collaborator for a manufacturing task; in order to infer the correct tool to deliver, the robot must interpret a person’s language and gesture over time.

Most existing approaches for interpreting language and gesture rely on unimodal models that do not integrate the two information sources, even though people fluidly use language and gesture together. Approaches that fuse information from language and gesture [Matuszek et al., 2014] do not take into account that information appears to the system over a period of time.

In contrast, we propose a Bayes’ filtering approach for interpreting multimodal information from language and gesture [Thrun et al., 2008]. Our framework relies on a factored observation probability that fuses information from language, hand gestures, and head gestures in real time to continuously estimate the object a person is referring to in the real world. We demonstrate our model in simulation, as well as providing quantitative results on a real-world RGB-D corpus of people referring to objects in the environment. These results demonstrate that our approach quickly and accurately fuses multimodal information in real time to continuously estimate the object a person is referring to.

II. RELATED WORK

Our approach is related to Holladay et al. [2014] but focuses on interpreting a person’s gestures rather than enabling

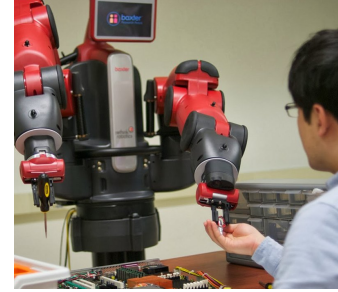


Fig. 1. Robots that collaborate with people need to understand their references to objects in the environment. For example, if a person asks for a tool using language and gesture, the robot needs to interpret the person’s reference in order to pick up the correct tool.

a robot to generate pointing gestures. A large body of work focuses on language understanding for robots [MacMahon et al., 2006, Dzifcak et al., 2009, Kollar et al., 2010, Matuszek et al., 2012]. Guadarrama et al. [2014] presents a framework for interpreting open-domain references to objects but focuses on interpreting language rather than language combined with gesture.

Many existing approaches for interpreting gesture rely on fixed vocabularies of gesture, such as “stop” or “follow” [Waldherr et al., 2000, Marge et al., 2011] without a principled way for fusing information from language and gesture. Our work unifies language and gesture interpretation into a single mathematical framework, and focuses on parameterized gestures such as pointing.

Matuszek et al. [2014] presented a multimodal framework for interpreting unscripted references to tabletop objects using language and gesture. Our approach similarly focuses on tabletop objects but uses language, gesture, and head pose, and integrates these disparate data sources continuously over time using a Bayes’ filtering framework. This approach enables the robot to continuously process new information and produce an estimate that converges over time to the correct object as new information is observed from the person.

III. TECHNICAL APPROACH

Our aim is to estimate a distribution over the object that a person is referring to given language and gesture inputs. We frame the problem as a Bayes’ filter [Thrun et al., 2008], where the hidden state, \mathcal{X} , is the set of m objects in the scene that the person is currently referencing. The robot observes the person’s actions and speech, \mathcal{Z} , and at each time step estimates a distribution over \mathcal{X} :

$$p(x_t | z_0 \dots z_{0:t}) \quad (1)$$

¹Computer Science Department, Brown University

To estimate this distribution, we take a Bayes' filtering approach and alternate performing a time update and a measurement update. The time update updates the belief that the user is referring to a specific subset of objects given previous information:

$$p(x_t|z_{0:t-1}) = \int p(x_t|x_{t-1}) \times p(x_{t-1}|z_{0:t-1}) dx_{t-1} \quad (2)$$

The measurement update combines the previous belief with the newest observation to update each belief state:

$$p(x_t|z_{0:t}) = \frac{p(z_t|x_t) \times p(x_t|z_{0:t-1})}{p(z_t|z_{0:t-1})} \quad (3)$$

$$\propto p(z_t|x_t) \times p(x_t|z_{0:t-1}) \quad (4)$$

A. Observation Model

We assume access to an observation model of the form:

$$p(z_t|x_t) \quad (5)$$

Observations consist of a tuple consisting of a person's actions, $\langle l, r, h, s \rangle$ where:

- l represents the observed origin (l_o) and vector (l_v) for the left arm.
- r represents the observed origin (r_o) and vector (r_v) for the right arm.
- h represents the observed origin (h_o) and vector (h_v) for head.
- s represents the observed speech from the user, consisting of a list of words.

Formally, we have:

$$p(z_t|x_t) = p(l, r, h, s|x_t) \quad (6)$$

We factor assuming that each modality is independent of the others given the state (the true object that the person is referencing):

$$= p(l|x_t) \times p(r|x_t) \times p(h|x_t) \times p(s|x_t) \quad (7)$$

The following sections describe how we model each type of input from the person.

Gesture. We model pointing gestures as a vector through three dimensional space. We calculate the angle between the gesture vector and the vector from the gesture origin to the mean of each cluster, and then use the PDF of a Gaussian with trained variance to determine the weight that should be assigned to that object. Let $\Phi(\langle origin \rangle, \langle point \rangle, \langle point \rangle)$ give the angle between the two points with the given origin.

$$p(l|x_t) \propto \mathcal{N}(\mu_l = 0, \sigma_l, \Phi(l_o, l_v, x_t)) \quad (8)$$

$$p(r|x_t) \propto \mathcal{N}(\mu_r = 0, \sigma_r, \Phi(r_o, r_v, x_t)) \quad (9)$$

Head Pose. Head pose is modeled in the same manner as arm gestures.

$$p(h|x_t) \propto \mathcal{N}(\mu_h = 0, \sigma_h, \Phi(h_o, h_v, x_t)) \quad (10)$$

Speech. We model speech with a simple bag of words model. We take the words in a given speech input and count how

many words in this text match descriptors (denoted x_d) of specific objects.

$$p(s|x_t) = \prod_{w \in s} p(w|x_t) \quad (11)$$

IV. EVALUATION

V. CONCLUSION

REFERENCES

- J. Dzifcak, M. Scheutz, C. Baral, and P. Schermerhorn. What to do and how to do it: translating natural language directives into temporal and dynamic logic representation for goal management and action execution. In *Proceedings of the 2009 IEEE International Conference on Robotics and Automation*, pages 3768–3773, Piscataway, NJ, USA, 2009. IEEE Press. ISBN 978-1-4244-2788-8.
- S. Guadarrama, E. Rodner, K. Saenko, N. Zhang, R. Farrell, J. Donahue, and T. Darrell. Open-vocabulary object retrieval. RSS, 2014.
- R. M. Holladay, A. D. Dragan, and S. S. Srinivasa. Legible robot pointing. RO-MAN, 2014.
- T. Kollar, S. Tellex, D. Roy, and N. Roy. Toward understanding natural language directions. In *Proceedings of HRI-2010*, 2010.
- M. MacMahon, B. Stankiewicz, and B. Kuipers. Walk the talk: connecting language, knowledge, and action in route instructions. In *Proceedings of the 21st national conference on artificial intelligence*, volume 2 of AAAI '06, pages 1475–1482. AAAI Press, 2006. ISBN 978-1-57735-281-5.
- M. Marge, A. Powers, J. Brookshire, T. Jay, O. C. Jenkins, and C. Geyer. Comparing heads-up, hands-free operation of ground robots to teleoperation. *Robotics: Science and Systems VII*, 2011.
- C. Matuszek, E. Herbst, L. Zettlemoyer, and D. Fox. Learning to parse natural language commands to a robot control system. In *Proc. of the 13th Intl Symposium on Experimental Robotics (ISER)*, 2012.
- C. Matuszek, L. Bo, L. Zettlemoyer, and D. Fox. Learning from unscripted deictic gesture and language for human-robot interactions. 2014.
- S. Thrun, W. Burgard, and D. Fox. *Probabilistic robotics*. MIT Press, 2008.
- S. Waldherr, R. Romero, and S. Thrun. A gesture based interface for human-robot interaction. *Autonomous Robots*, 9(2):151–173, 2000.