

Investigating the Effects of Visual Saliency on Deictic Gesture Production from a Humanoid Robot

Aaron St. Clair, Ross Mead, Maja J Mataric
Interaction Lab, Viterbi School of Engineering
University of Southern California
Los Angeles, CA 90089
{astclair, rossmead, mataric}@usc.edu

ABSTRACT

In many collocated human-robot interaction scenarios, robots are required to accurately and unambiguously indicate an object or point of interest in the environment. Realistic, cluttered environments containing many visually salient targets can present a challenge. In this paper, we describe an experiment and results detailing the effects of visual saliency and pointing modality on human perceptual accuracy of robot deictic gestures (head and arm pointing) and compare these results to the perception of human pointing.

Categories and Subject Descriptors

I.2.9 [Robotics]: Operator interfaces. H.1.2 [Models and Principles]: User/Machine System.

General Terms

Design, Experimentation, Human Factors, Verification.

Keywords

Human-robot interaction, deixis, communicative behavior, social signaling, visual saliency, nonverbal.

1. INTRODUCTION

To carry on sustained interactions with people, autonomous robots must be able to effectively and naturally communicate with humans in many different interaction contexts. Besides natural language, humans employ both coverbal (e.g., beat gestures) and nonverbal modalities—including facial expression, proxemics, eye gaze, head orientation, and arm gestures, among others—to signal their own intentions and to attribute intentions to the actions of others. Prior work has demonstrated that robots can successfully employ these channels [1, 2, 3, 4]. Our aim is to develop a general, empirical understanding of design factors when employing multimodal communication channels with a robot. In this paper, we limit our focus to a study of deictic gestures since (1) their use in human communication has been widely studied [6, 7, 8], (2) they are relatively simple to map to intentional

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

HRI'11, March 6–9, 2011, Lausanne, Switzerland.
Copyright 2010 ACM 1-58113-000-0/00/0010...\$10.00.

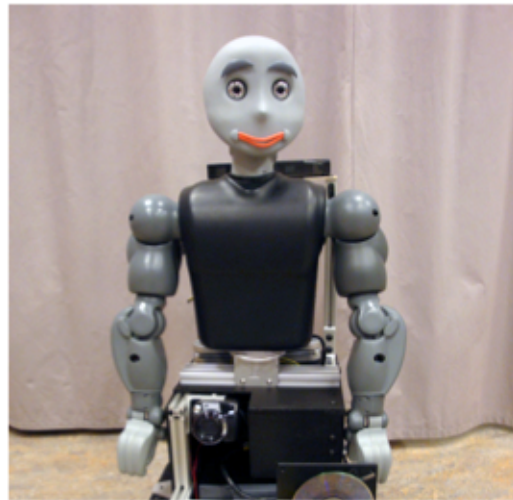


Figure 1. The Bandit humanoid on top of a Pioneer base.

constructs in context [10], and (3) they are generally useful to robots interacting in shared environments since they serve to focus attention and refer to objects. To achieve robust deixis via gesture in a human-robot context, it is necessary to validate the perceived referent. This paper presents results from an experimental study of human perception of a robot's deictic gestures under a set of different environmental visual saliency conditions and pointing modalities using our upper-torso humanoid robot, Bandit.

2. EMBODIED DEICTIC GESTURE

Multi-disciplinary research from neuroscience and psychology has demonstrated that human gesture production is tightly coupled with language processing and production [11, 12]. There is also evidence that gestures are adapted by a speaker to account for the relative position of a listener and can, in some instances, substitute for speech functions. Bangerter [8] and Louwerse & Bangerter [10] demonstrated this substitution effect for deictic gestures by studying performance on a target disambiguation task, where they found that deictic speech combined with deictic gesture offered no additional performance gain compared to one or the other used separately.

These findings have important implications for the field of human-robot interaction. Despite on-going work in speech

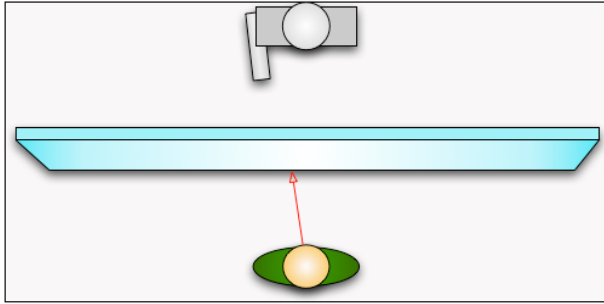


Figure 1. The diagram of the experiment layout. The robot performs a gesture relative to the board and the person indicates an estimate of the gesture target with a laser

recognition and natural language processing, it remains difficult to employ verbal communication as a reliable channel beyond simple commands, especially on mobile robots. Robots interacting with humans in a shared physical environment should be able take advantage of other social channels to both monitor and communicate intent during the course of an interaction without complete reliance on speech or environment knowledge. To make this possible it is necessary to gain an empirical understanding of how to map well-studied human gestures to robots of varying capabilities and embodiments. Specifically, we are interested in identifying variables for proper production of robot gestures to best realize some fixed interpretation by a human audience. In general, this is difficult, for the same reasons that processing natural language is difficult; many gestures are context-dependent and rely on accurately estimating a mental model of the scope of attention and possible intentions for people’s actions given only low-level perceptual input.

Deictic gestures, however, are largely consistent in their mapping to linguistic constructs, such as “that” and “there”, and serve to focus the attention of observers to a specific object or location in the environment, or perhaps to indicate an intended effect involving such an object (e.g., “I will pick up that”). These characteristics, while simplifying their interpretation and production, also make the gestures useful for referring to objects and grounding attention. Intentional analysis and timing are still nontrivial, except in the context of performing a specific pre-determined task.

Both recognition [13, 14, 15, 16, 17] and production [18, 19, 20, 21] of deictic gestures has been studied in human-human, human-computer, and human-robot interaction settings. Our work adds to this field as a step toward obtaining an empirically grounded HRI model of deictic gestural accuracy between humans and robots, with implications for the design of robot embodiments and control systems that perform situated distal pointing.

A study of the literature on human deictic pointing behavior suggests a number of possible variables that could potentially affect the robustness of referent selection when having a robot employ deictic gestures. The physical embodiment of the robot constrains the appearance of the part of the robot used for pointing. Pointing with a blunt or irregularly shaped object has lower resolution than a sharp, pointed object with a clear vector interpretation. It has also been shown that humans take into account the position and orientation of their audience when staging a gesture in order to improve the audience’s interpretation [6]; this consideration of audience perspective is even a core concept in hand-drawn [22] and computer animation [23].

Assuming mobility a robot could relocate or reorient itself relative to its audience or to the referent target to improve viewer’s interpretation accuracy [24]. Other considerations beyond physical variables include parameterizing the timing and appearance of the gesture itself. Humans in most cases attempt to minimize effort when performing deixis and thus do not tend to make maximally expressive gestures unless necessary, such as when an object is far away [25]. Most robots, however, point without accounting for distance to the target by using a more or less constant arm extension or head gaze that is reoriented appropriately [18, 5]. Finally, since the gesture is grounded with respect to a specific referent in the environment, the robot must be able to correctly segment and localize visually salient objects in the environment at a similar granularity to the people with whom it is interacting. Biologically inspired methods to assess and map visually salient features and objects in an environment exist [26, 27], as do models of human visual attention selection [28], but the role of visual saliency during deictic reference by a robot is largely uninvestigated.

3. EXPERIMENTAL DESIGN

Given the large number of possible variables when considering how to best reference, via gesture, a point in the environment with a particular robot embodiment, we conducted an initial pilot study to test distance and angle to target, distance and angle to audience, and modality for a face-to-face interaction between a person and our upper-torso humanoid robot, Bandit (Figures 1). No strong correlations emerged in early testing, so we narrowed the set of conditions and hypotheses.

3.1 Hypotheses

We conducted a factorized IRB-approved experiment over three robot pointing modalities: the head with 2 degrees-of-freedom (DOF), the arm with 7 DOF, and both together (i.e., head+arm) with two saliency conditions: a blank (or non-salient) environment and an environment with several highly and equally visually salient targets. Since our results, particularly the modality condition, may be specific to Bandit, we also conducted a similar, but smaller, test with a human performing the pointing gestures, for comparison.

3.1.1 Modality

The conditions tested include head (Figure 1a), away-from-body, straight arm (Figure 1b), cross-body, bent arm (Figure 1c), and combined head and arm (Figure 1d). We hypothesized that the arm modality would lead to more accurate perception since, when fully extended, it is the most expressive and easily interpreted as a vector from the robot to the screen. Since Bandit’s head does not have moveable eyes, the point of reference is somewhat ambiguous and could lead to pointing error. Our kinematic calculations solved for the midpoint of the eyes to align to the target; this information was not shared with experiment participants. Additionally, we expected to see an effect between away-from-body, straight arm (Figure 1b) points that occur on one side of the screen and the cross-body bent-arm (Figure 1c) gesture used for the other side, with the bent-arm being more difficult to interpret, since it is staged in front of the robot’s body rather than laterally [22, 23, 24]. A similar effect was seen in human pointing in work by Bangerter [9], showing that humans are capable of estimating vectors accurately from body pose, although the exact mechanisms for referential grounding remain unclear. Finally, we hypothesized using both modalities together would reduce error

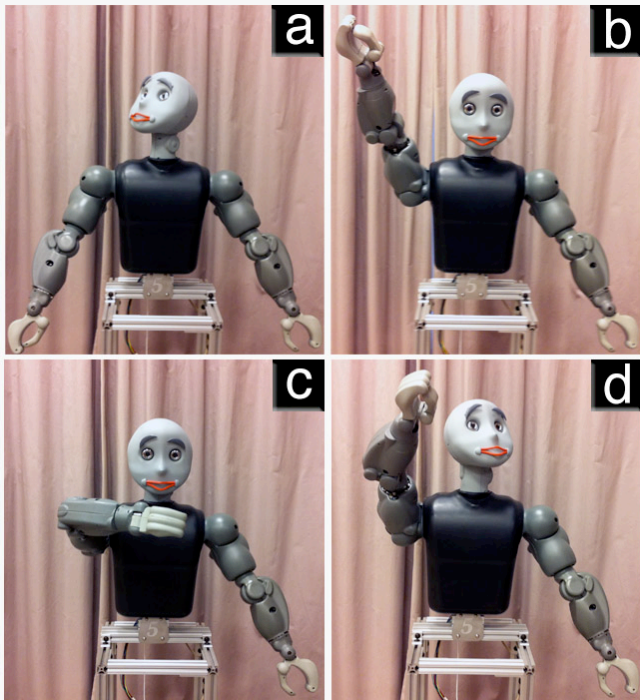


Figure 2. Bandit pointing with its head (a), straight-arm (b), bent-arm (c), and head+arm (d).

relative to a single modality, since participants would have two gestures on which to base the estimate.

3.1.2 Saliency

For the two saliency conditions, we hypothesized that the salient objects would affect people’s interpretations of the points. Specifically, we anticipated that people would “snap to” the salient objects, thus, reducing error for points whose targets were on or near markers, whereas in the non-salient condition there are no points of reference to bias estimates. We did not expect to see any difference in the performance of each pointing modality when comparing the salient and non-salient conditions.

3.2 Implementation

In the experiments, the participant is seated directly facing Bandit from a distance of 6 feet (1.8 meters). The robot and the participant are separated by a transparent, acrylic screen measuring 12 feet by 8 feet (2.4 by 3.6 meters) (see Figure 2). The screen, thus, covers a horizontal field of view from approximately -60 to 60 degrees and a vertical field of view -45 to 60 degrees. The robot performs a series of deictic gestures and the participant is asked to estimate their referent location on the screen. All gestures are static and are held indefinitely until the participant estimates a location, after which the robot returns to a home location (looking straight forward with its hands and its sides) before performing the next gesture. Participants are given a laser pointer to mark their estimated location for each gesture. These locations are recorded using a laser rangefinder, which is placed facing upwards at the base of the screen. For each gesture, an experimenter places a fiducial marker over the indicated location,

which is subsequently localized using the rangefinder data within approximately 1 cm. The entire experiment is controlled via a single Nintendo Wiimote™, with which the experimenter can record marked locations and advance the robot to point to the next referent target.

The face-to-face nature (see Figure 3) of the experiment was chosen intentionally although other work in gesture perception [8] has tested human deictic pointing accuracy when the pointer and the observer are situated more or less side-by-side, observing a scene. In our work, and in most human-robot interaction settings the robot is typically facing the participant, rather than side-by-side, which, in terms of proxemics, is more likely to occur when coordinated motion and interaction are concurrent (e.g. the robot and human walking together while an interaction is taking place [29]). For this reason, our design tests the face-to-face scenario, since it is more applicable to the types of proxemic configurations we have tended to encounter during our prior work.

3.2.1 Modality

The robot gestures to locations by moving its head, arm, or both together. A single arm was used resulting in cross-body, bent-arm gestures for points on one side of the screen and away-from-body straight-arm gestures for points on the other side of the screen. The arm was not modified for pointing and, thus, the end-effector was simply the 1-DOF gripper in a closed position. All the gestures were static, meaning the robot would leave the home position, reach the gesture position, and hold it indefinitely until it would return to the home position for the next gesture. This was intended to minimize any possible timing effects by giving participants as long as they needed to estimate a given gesture. Participants were also asked to only turn the laser pointer after they had visually selected an estimated (perceived) target location to prevent them from using the laser pointer to line up the robot arm and head with the actual target location.

3.2.2 Saliency

The screen itself was presented with two visual saliency conditions: one in which it was completely empty (i.e., non-salient) and one in which it was affixed with eight round markers distributed at random (i.e., salient). In the salient case, the markers were all 6 inches (15 cm) in diameter and were identical in shape

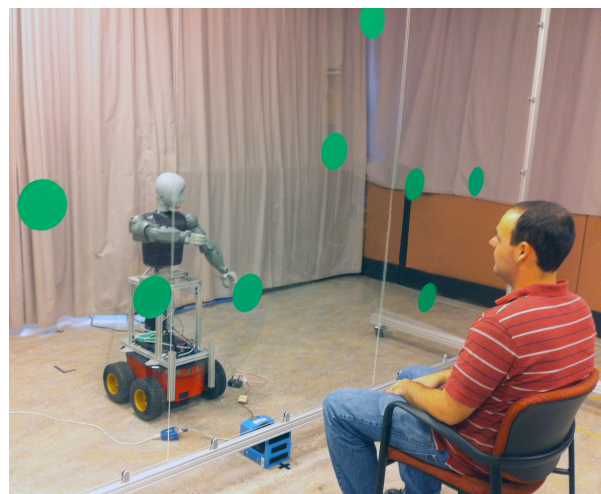


Figure 3. A picture of the experimental setup, with Bandit indicating a cross-body point to the subject in the foreground.

and color. Experiments were conducted in two phases, reflecting the two saliency conditions. In the salient condition, the robot's gestures included 60 points toward randomly selected salient targets and 60 points chosen randomly to be on the screen, but not necessarily at a salient target. In the non-salient condition, 74 of the points within the bounds of the screen were chosen at random and the remaining 36 were chosen to sample a set of 4 calibration locations with each pointing modality. These calibration points were used to assess the consistency and normality of the error in the robot's actual pointing and the participant's perception to determine whether a between-subjects comparison was possible. All three pointing modalities (head-only, arm-only, and head+arm) were used in both the salient and non-salient cases.

3.2.3 Human Pointing

The human-human pointing condition was conducted by replacing Bandit with an experimenter with the intention of anecdotally comparing robot pointing with typical human pointing in the same scenario. Since people point by aligning a chosen end-effector with the referent target using their dominant eye [8], conducting the experiment with a human pointer introduces the confound that they cannot point with the arm-only modality. We also found the head-only modality difficult to measure accurately and, for this reason, only the head+arm modality was tested. These experiments were conducted by replacing the robot with an experimenter who held a Nintendo Wiimote™ in his or her non-pointing hand. Two different vibration patterns signaled whether to point to a target location or a location selected arbitrarily. The experimenter pointed while holding a laser pointer concealed in the palm of their hand and held their pose. Another experimenter then marked both the subject's and experimenter's points, as before.

3.2.4 Surveys

In addition to the pointing task, we also administered a survey asking participants to estimate their average error with respect to modality and location on the screen and to rate each modality in terms of preference on a Likert scale. The surveys also collected background information such as, handedness and level of prior experience with robots.

4. RESULTS

4.1 Participants and Data

A total of 30 runs of the experiment were conducted as described, with 17 (12 female, 7 male) participating in the non-salient condition and 12 (7 male, 5 female) participants in the salient condition. In total, approximately 3600 points were estimated and recorded. The conditions are close to equally weighted with the exception of the non-salient arm-only and head+arm conditions, which was done initially to allow for comparison of the cross-body versus away-from-body arm gestures. The number of points collected for each condition is presented in Table 1. Participants were recruited from on-campus sources and all were undergraduate or graduate students at USC from various majors. The participants were roughly age and sex matched with an average age of 20. The data collected for each run included a log file recording the desired target on the screen (i.e., the desired location the robot should have pointed to), the actual target on the screen for each modality (i.e., the location the robot actually pointed to), and the perceived point as indicated by the participant and recorded by the laser rangefinder. We also captured timing

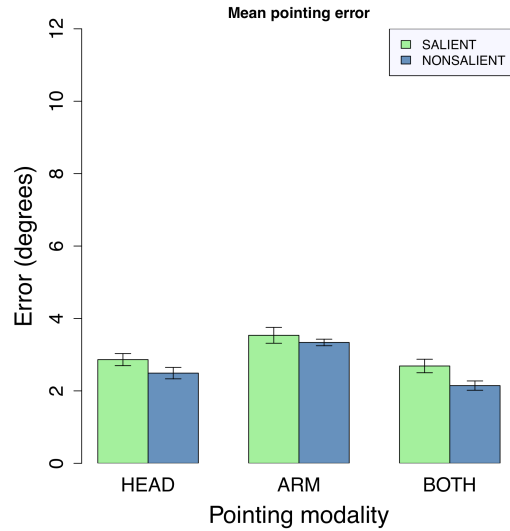


Figure 4. Mean angular pointing error.

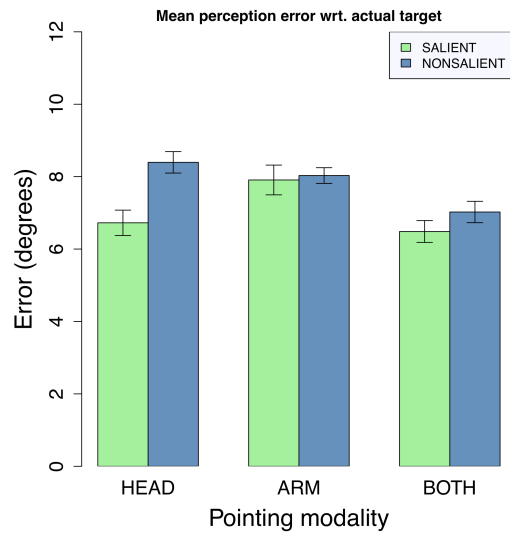


Figure 5. Mean error between perceived and desired targets.

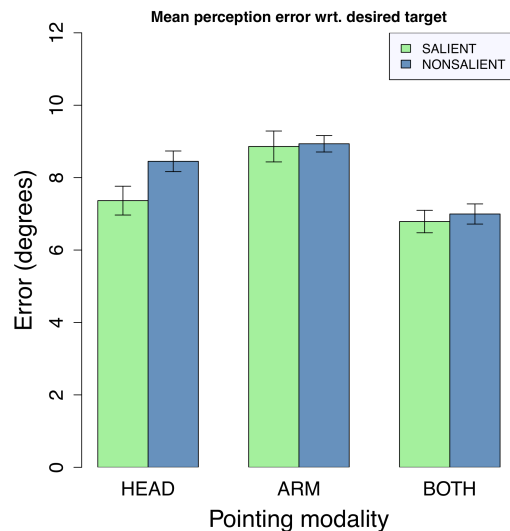


Figure 6. Mean error between perceived and actual targets.

data for each point and video of the sessions taken from a camera mounted behind and above the robot.

Table 1. Data counts. Arm is over-represented in the non-salient condition to compare away-from-body and cross-body.

	Head	Arm	Both
Non-salient	565	890	298
Salient	474	481	489

4.2 Perceived Error Analysis

We conducted a two-way analysis of variance (ANOVA) with modality and saliency as the independent factors. Both were found to have significant effects on the angular error between perceived and desired target points as well as on perceived and actual target points (see Table 2). Additionally, the interaction effects between the modality and saliency factors were found not to be significant. Mean angular error computed from the perspective of the person and confidence intervals are shown in the graphs in Figures 4-6. We used angular error as a metric to effectively normalize for different distances to target. For comparison purposes, human perceptual error when estimating human pointing gestures (arm or eye gaze) has been measured to be approximately 2-3 degrees for people up to 2.7 meters apart [8, 30]. We conducted post-hoc analysis using the Tukey’s honestly significant differences (HSD), which revealed that mean error tends to be about 1.5 degrees higher for arm points with $p < 0.01$, and that using both modalities tends to outperform the arm modality in most cases with the means differing by 1.8 degrees in the salient case and 2.0 degrees in the non-salient case with significance of $p < 0.01$. The arm alone, however, was equally bad in both saliency conditions.

Table 2. F and p values for 2-way ANOVA with perceived-desired and perceived-actual angular error.

	Perceived-Desired		Perceived-Actual	
	F	p	F	p
Saliency	F=4.53,	p<0.03*	F=11.4	p<0.01*
Modality	F=9.6	p<0.01*	F=5.0	p<0.01*

To compare cross-body bent-arm versus away-from-body straight-arm gesture, we looked at arm points in the non-salient case. Partitioning them into two sets, depending on the side of the screen they were on, resulted in 450 cross-body points and 440 straight-arm points. Conducting a one-way ANOVA with arm gesture type as the dependent variable, we find there is a significant difference between the straight-arm case ($M=5.6$ degrees, $SD=3.7$) and the bent-arm case ($M=10.4$ degrees, $SD=7.9$) with $p < 0.001$ (see Figure 7). We also obtained similar results when conducting a full 3-way ANOVA with the other two conditions, although there were interactions between some of the factors in addition to significant main effects, likely due to the high variance in arm accuracy.

To assess whether the accuracy of the other modalities varied with angle to target, we first fit a linear regression model with angular error as the dependent variable and the desired target as the independent variable. The resulting model did not perform well upon cross validation suggesting that the error was nonlinear

in nature. To cope with the nonlinearity, we then binned points by

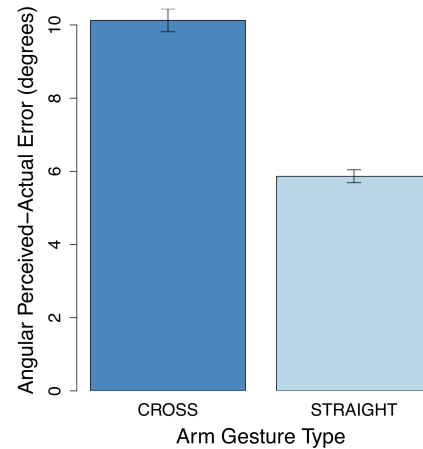


Figure 7. Mean angular error for straight versus bent arm.

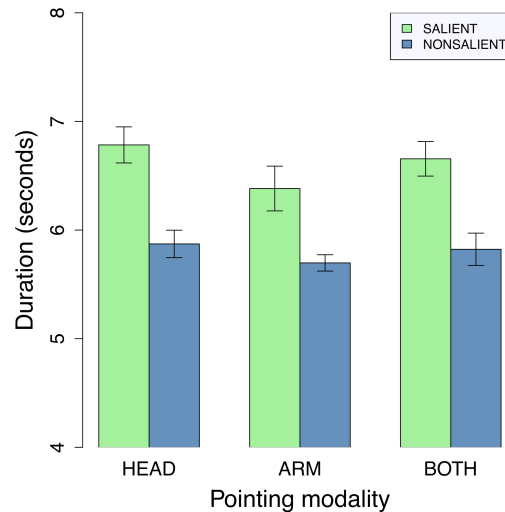


Figure 8. Mean perception time for all conditions.

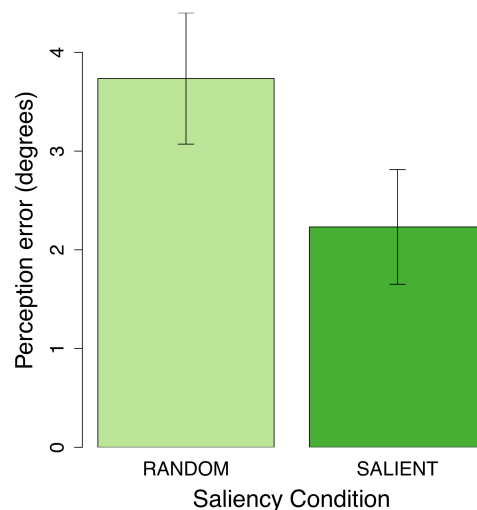


Figure 9. Mean error by saliency condition, human pointer.

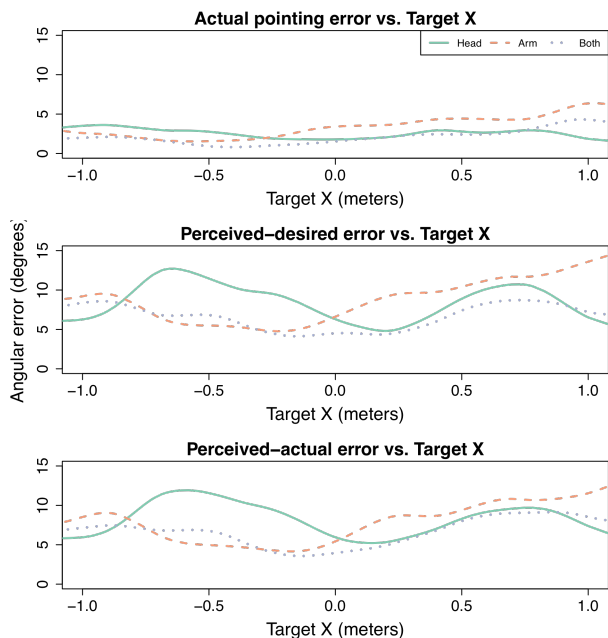


Figure 10. Mean angular error with respect to horizontal target position as seen from the participant’s perspective.

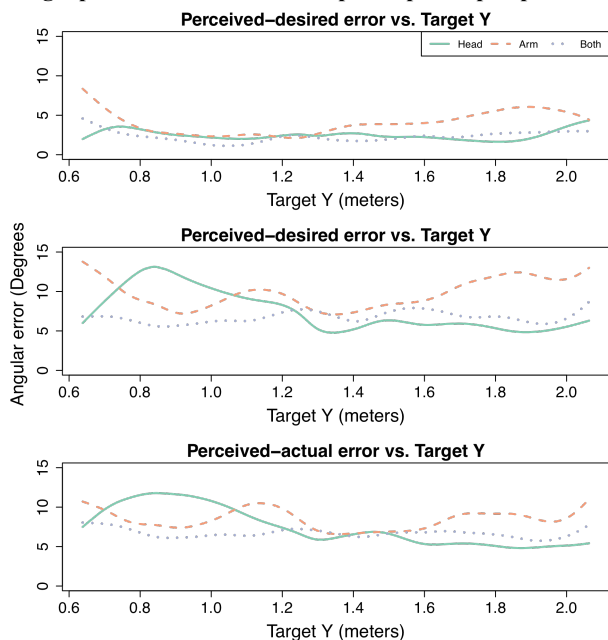


Figure 11. Mean angular error with respect to vertical target position as seen from the participant’s perspective

angle to target in 9 uniform intervals covering the extent of the screen. We then performed an n -way ANOVA with target x and y coordinates as a conditional factor, and found that the head-only and head+arm conditions were significantly better ($p < 0.001$) in the center of the screen with error increasing about halfway to the edge before leveling out. These effects were largely symmetric for the head-only and head+arm conditions. The arm-only condition, as described above, was asymmetric and was significantly worse in almost all cases except for the middle of the left side of the

screen corresponding to away-from-body straight-arm points. Finally, all modalities tended to result in more erroneous estimates in the lower extreme of the screen. Figure 9 depicts a smoothed cubic fit to the entire dataset for each modality with respect to target x -coordinates and Figure 10 with respect to the target y -coordinates. Each is pictured from the participant’s perspective, meaning cross arm gestures correspond to the right side of the horizontal graph.

The average time taken to estimate each point was nearly a second faster in the non-salient case ($M=6.6$, $SD=2.6$) than in the salient case ($M=5.8$, $SD=2.3$). This effect was found to be significant with $p < 0.001$, while there was no significant difference for the modality conditions (see Figure 8).

4.3 Human Pointing

For the human pointing phase of the experiment, we collected a total of 70 data points from 2 participants. While this is not a considerable amount of data and should include more participants before drawing conclusions, we did find a significant ($p < 0.09$) effect when comparing the two saliency conditions. For the salient condition ($M=2.23$, $SD=2.46$), the error was small enough that both the pointer and the observer were able to “hit” all the salient targets, while the non-salient condition ($M=3.74$, $SD=2.74$) resulted in a 60% increase in error (see Figure 9). When plotting the error versus intended referent x -coordinate, we see that points directed at the center of the screen appear to result in lower perception error than points directed between the center and the periphery. Overall, the error in estimating human-produced points appears to have a similar profile to that of the robot-produced points; however, more investigation is necessary.

4.4 Survey Responses

In the responses to the survey, participants in the non-salient condition estimated that their points were within an average of 28 centimeters (11 inches); this is very close to the mean error of 27 centimeters we found in practice. There was no significant difference between participants’ estimated error when comparing across the two conditions. Pointing with the head-only and with the head+arm were preferred by the majority of the participants, with only 4 (or 16%) stating a preference for the arm modality. When asked if there was a noticeable difference in straight-arm points versus the bent-arm, 65% said there was with the remainder not seeing a difference. Ten out of 12 (or 83%) of the participants in the salient condition said that the markers would have an effect on their estimate of the referent target.

5. DISCUSSION

5.1 Visual Saliency

The mean error, as computed (using the perceived point and the desired target points), tells us how close to a desired target (either a randomly chosen one in the non-salient case or one of the markers in the salient case) the robot was actually able to indicate. The performance of the head-only in the salient condition is improved by approximately 1 degree, the arm-only is not appreciably different, and head+arm exhibits only a modest gain. This suggests that the snap-to-target effect that we expected to see when salient objects were introduced is modest at best, resulting in a best-case improvement of approximately 1 degree. This is also seen when we consider the mean perceived-actual error and find that nearly every condition is slightly more accurate. This suggests that participants estimate the referent to be closer to the actual point the robot is physically indicating than the nearest

salient object. This could be potentially useful because it allows us to consider pointing without having to assess scene saliency beforehand. That is, if we have not specified the referential target of a point *a priori*, through some other means, such as verbal communication or previous activity, people tend to evaluate the point in an ad hoc manner by taking their best guess. When disambiguating referents, if there are unknown salient objects in the environment, we can anticipate their effects on the perception of a given gesture to be small enough in most cases that a precise point to our actual target should suffice to communicate the referent.

5.2 Modality

As we hypothesized, the modalities did result in different pointing accuracy profiles. When considering modalities, pointing with the head+arm does appear to perform appreciably better than either the arm-only or the head-only, in most cases. One possible explanation of this is that it more closely emulates typical human pointing, in which people tend to align arm gestures with their dominant eye [8], or that multiple modalities provide more diverse cues that indicate the referential target resulting in better priming of the viewer to interpret the gesture. The poor performance of the arm in the salient condition was also somewhat unexpected. This might be due to its higher actual error compared to the head, which we discovered to be related to a small “dead-band” created in the Bandit firmware that can become compounded over each joint in the arm, resulting in poor overall performance. Another source of the error could be use of the cross-body arm gesture, which, while equally weighted, resulted in nearly twice the perceptual error as compared to the away-from-body arm. This might be a result of the reduced length of the arm, which forces people to estimate the vector based on only the forearm versus the entire arm as in the away-from-body case. Another explanation is that the gesture is staged against the body, that is, with minimal silhouette and is, thus, more difficult for people to see [Disney, Mead]. In either case, roughly one-third of the participants did not notice a difference in the arm gestures while their performance was, in fact, affected. This illustrates the impact that gesture and embodiment design can have on interpretation, and underscores the need to validate gestural meaning with people.

When considering the horizontal and vertical target position analyses, we see that people are best at estimating points directly between the participant and the robot. Performance then drops off when the target is located laterally, above, or below. This effect could be due to a field-of-view restriction, preventing the viewer from seeing both the robot’s gesture and the target at the same time in high acuity, foveal vision. Estimating these points then requires the viewer to saccade their head between the two points of interest. We believe the slight improvement at the far periphery for some of the modalities is due to the fact that we informed participants that the points would be on the screen, thus, creating a bound for points near the screen edges.

5.3 Human Pointing

The results of our smaller-scale investigation of human pointing did find that the salient condition resulted in approximately 1.5 degrees less error than in the non-salient condition, which is consistent with our finding using the robot pointer. Also, the 2-degree perceptual accuracy that we found when testing a human pointer seems to agree with prior study of humans conducted in relevant literature. It is also worthy of note that, although the deictic pointing performance of the robot is several times worse than what we saw in the human experiment or

would expect from literature, we can use the estimate of our resolving power (i.e., the minimum angle between referents that we could hope to convey) to inform controller design and ensure that the robot repositions itself or gets close enough to prevent these effects. The salient condition also resulted in a 16% increase (or approximately 1 second) in time needed to estimate the gesture. This is intuitive, as the participants were presented with more stimuli in the form of the salient objects and, thus, take some extra time to ground the point, possibly checking to see if it is coincident with any objects first. This information could be useful in developing methods for effective timing control.

6. FUTURE WORK

One obvious next step is to conduct a similar experiment with a robot of a different embodiment to evaluate whether the same general conclusions hold true or if they are tightly coupled to the specific appearance of Bandit. Since the project was developed using Willow Garage’s Robot Operating System (ROS), substituting different robots into the experiment will involve minimal changes to the codebase. We are also currently developing a deictic gesture (pointing) software package for the PR2 and other robots that have an URDF visualization specification, and plan on using it to run the experiment with the PR2. Formal studies of other relevant variables mentioned in the introduction (such as angle to target, timing, and embodiment, among others), as well as comparable studies with human pointing, is also necessary to develop a better understanding of how robot deictic gestures compare and contrast to their human counterparts.

We also plan on analyzing head movements over time as participants perform the task. There is a noticeable saccade effect, in which people saccade back and forth between an estimated point and the robot before finally glancing at the experimenter to mark the point. This process appears to be similar in nature to the regressive eye movements used to measure gesture clarity in [8]. By tracking head movement [31], we plan to check for significance in the number of saccades as an indicator of perceptual difficulty as well as a classification system capable of establishing when and where joint attention has been established.

We are seeking ways to automatically measure estimated points, thereby allowing us to remove the potentially-distracting experimenter from the room.

Finally, we plan on using the data presented to construct a parameterized error model to allow Bandit to perform effective deixis to objects in a mapped environment. Specifically, we plan on integrating this work for performing valid deixis in the context of collaborative tasks.

7. CONCLUSION

In this paper, we presented the results of a study of human perception of robot gestures intended to test whether visual saliency and embodied pointing modality have an effect on the performance of human referent resolution. Our results suggest that environmental saliency, when employing only deictic gesture to indicate a target, results in only a modest bias effect. We also demonstrated that by pointing with two combined and synchronized modalities, such as head+arm, we are able to consistently outperform one or the other individually. Additionally, we found that the physical instantiation of the gesture (i.e. how it is presented to the observer) can have drastic effects on perceptual accuracy, as noted when comparing bent-arm and straight-arm performance.

8. ACKNOWLEDGMENTS

The authors would like to thank Hieu Minh Nguyen and Karie Lau for their help with data collection. This work was supported in part by National Science Foundation (NSF) grants CNS-0709296, IIS-0803565, and IIS-0713697 and the ONR MURI program (N00014-09-1-1031 and N00014-08-1-0693). The second author was supported by an NSF Graduate Research Fellowship

9. REFERENCES

- [1] B. Scassellati, "Investigating models of social development using a humanoid robot," vol. 4, pp. 2704 – 2709 vol.4, jul. 2003.
- [2] A. Brooks and C. Breazeal, "Working with robots and objects: Revisiting deictic reference for achieving spatial common ground," in *Proceedings of the 1st ACM SIGCHI/SIGART conference on Human-robot interaction*, p. 304, ACM, 2006.
- [3] C. Breazeal, C. Kidd, A. Thormaz, G. Hoffman, and M. Merlin, "Effects of nonverbal communication on efficiency and robustness in human-robot teamwork, *iee/rsj int*," in *IEEE/RSJ Int. Conf. on Intelligent Robots and Systems (IROS2005)*, pp. 383–389, 2005.
- [4] C. Sidner, C. Kidd, C. Lee, and N. Lesh, "Where to look: a study of human-robot engagement," in *Proceedings of the 9th International Conference on Intelligent User Interfaces*, p. 84, ACM, 2004.
- [5] B. Mutlu, T. Shiwa, T. Kanda, H. Ishiguro, and N. Hagita, "Footing in human-robot conversations: how robots might shape participant roles using gaze cues," in *Proceedings of the 4th ACM/IEEE International Conference on Human Robot Interaction, (HRI'09)*, 2009.
- [6] A. Ozyurek, "Do speakers design their co-speech gestures for their addressees? the effects of addressee location on representational gestures," *Journal of Memory and Language*, vol. 46, no. 4, pp. 688–704, 2002.
- [7] N. Nishitani, M. Schurmann, K. Amunts, and R. Hari, "Broca's region: From action to language," *Physiology*, vol. 20, no. 1, p. 60, 2005.
- [8] A. Bangerter, "Accuracy in detecting referents of pointing gestures unaccompanied by language," *Gesture*, vol. 6, no. 1, pp. 85–102, 2006.
- [9] A. Bangerter, "Using pointing and describing to achieve joint focus of attention in dialogue," *Psychological Science*, vol. 15, no. 6, p. 415, 2004.
- [10] M. Louwerse and A. Bangerter, "Focusing attention with deictic gestures and linguistic expressions," in *Proceedings of the 27th Annual Meeting of the Cognitive Science Society*, 2005.
- [11] R. Mayberry and J. Jaques, *Gesture production during stuttered speech: Insights into the nature of gesture-speech integration*, ch. 10, pp. 199–214. Cambridge University Press, 2000.
- [12] S. Kelly, A. Ozurek, and E. Maris, "Two sides of the same coin: Speech and gesture mutually interact to enhance comprehension," *Psychological Science*, vol. 21, no. 2, pp. 260–267, 2009.
- [13] R. Cipolla and N. Hollinghurst, "Human-robot interface by pointing with uncalibrated stereo vision," *Image and Vision Computing*, vol. 14, no. 3, pp. 171–178, 1996.
- [14] D. Kortenkamp, E. Huber, and R. Bonasso, "Recognizing and interpreting gestures on a mobile robot," in *Proceedings of the National Conference on Artificial Intelligence*, pp. 915–921, 1996.
- [15] K. Nickel and R. Stiefelhagen, "Visual recognition of pointing gestures for human-robot interaction," *Image and Vision Computing*, vol. 25, no. 12, pp. 1875–1884, 2007.
- [16] P. Pook and D. Ballard, "Deictic human/robot interaction," *Robotics and Autonomous Systems*, vol. 18, no. 1-2, pp. 259–269, 1996.
- [17] N. Wong and C. Gutwin, "Where are you pointing? : the accuracy of deictic pointing in cves," in *Proceedings of the 28th international conference on Human factors in computing systems*, pp. 1029–1038, ACM, 2010.
- [18] M. Marjanovic, B. Scassellati, and M. Williamson, "Self-taught visually guided pointing for a humanoid robot," in *From Animals to Animats 4: Proceedings of the Fourth International on Conference Simulation of Adaptive Behavior*, pp. 35–44.
- [19] J. Trafton, N. Cassimatis, M. Bugajska, D. Brock, F. Mintz, and A. Schultz, "Enabling effective human–robot interaction using perspective-taking in robots," *IEEE Transactions on Systems, Man, and Cybernetics*, vol. 35, no. 4, pp. 460–470, 2005.
- [20] O. Sugiyama, T. Kanda, M. Imai, H. Ishiguro, N. Hagita, and Y. Anzai, "Humanlike conversation with gestures and verbal cues based on a three-layer attention-drawing model," *Connection science*, vol. 18, no. 4, pp. 379–402, 2006.
- [21] Y. Hato, S. Satake, T. Kanda, M. Imai, and N. Hagita, "Pointing to space: modeling of deictic interaction referring to regions," in *Proceeding of the 5th ACM/IEEE International Conference on Human-robot Interaction*, pp. 301–308, ACM, 2010.
- [22] F. Thomas and O. Johnston, *The Illusion of Life: Disney Animation*. Hyperion, 1981.
- [23] J. Lasseter, "Principles of traditional animation applied to 3D computer animation," in *ACM Computer Graphics*, vol. 21, no. 4, pp. 35–44, July 1987.
- [24] R. Mead and M.J. Mataric, "Automated caricature of robot expressions in socially assistive human-robot interaction," in *The 5th ACM/IEEE International Conference on Human-Robot Interaction (HRI2010) Workshop on What Do Collaborations with the Arts Have to Say about HRI?*, Osaka, Japan, March 2010.
- [25] A. Bangerter, "Using pointing and describing to achieve joint focus of attention in dialogue," *Psychological Science*, vol. 15, no. 6, p. 415, 2004.
- [26] L. Itti, C. Koch, and E. Niebur, "A model of saliency-based visual attention for rapid scene analysis," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 20, no. 11, pp. 1254–1259, 1998.
- [27] D. Walther, L. Itti, M. Riesenhuber, T. Poggio, and C. Koch, "Attentional selection for object recognition: a gentle way," in *Biologically Motivated Computer Vision*, pp. 251–267, Springer, 2010.
- [28] R. Desimone and J. Duncan, "Neural mechanisms of selective visual attention," *Annual review of neuroscience*, vol. 18, no. 1, pp. 193–222, 1995.
- [29] G. Butterworth and S. Itakura, "How the eyes, head and hand serve definite reference," *British Journal of Developmental Psychology*, vol. 18, no. 1, pp. 25–50, 2000.
- [30] R. Mead, "Space: a social frontier," poster presented at the *Workshop on Predictive Models of Human Communication Dynamics*, Los Angeles, California, August 2010.
- [31] L. Morency, J. Whitehill, and J. Movellan, "Generalized adaptive view-based appearance model: Integrated framework for monocular head pose estimation," in *8th IEEE International Conference on Automatic Face & Gesture Recognition, 2008. FG'08*, pp. 1–8, 2008.