

---

# Model-Based Attention Acquisition Strategies for Socially Interactive Robotics

**Eric C. Deng**

University of Southern California  
Los Angeles, CA 90007, USA  
denge@usc.edu

**Maja J. Matarić**

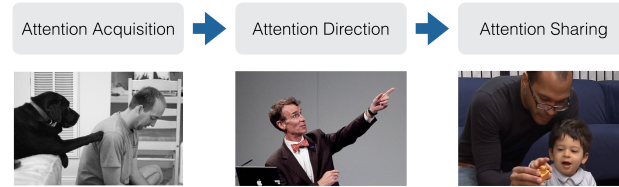
University of Southern California  
Los Angeles, CA 90007, USA  
mataric@usc.edu

## Introduction

Roboticians in the field of human-robot interaction (HRI) are seeking to develop algorithms and methods for interactive robots that teach, help, and play with human users [6]. Our work is exploring facets of *joint attention*, a foundational social skill that is learned in the early stages of development [10] and is of special interest to researchers in fields such as social skill therapy for children with autism [1]. Most human-robot interaction systems use cues such as calling the user by name or making other salient sounds, all of which are known to be effective for *attention acquisition* [7, 12], a necessary step in joint attention.

However, such attention-getting behaviors can be disruptive as well as disturbing, affecting the user's perception of the robot's agency. That is particularly an issue in socially assistive robotics (SAR), the area of HRI that develops methods for assisting users through social rather than physical means, focusing on health and education problem domains [5].

In this work, we focus on gesture-based attention acquisition for joint attention in SAR. We introduce a model-based approach to generating effective, minimal physical attention acquisition behaviors and discuss experimental design of ongoing and future studies as well as recent results. The ultimate goal of using a model-based approach is to *deliver*



**Figure 1:** The stages of establishing joint attention with a social partner

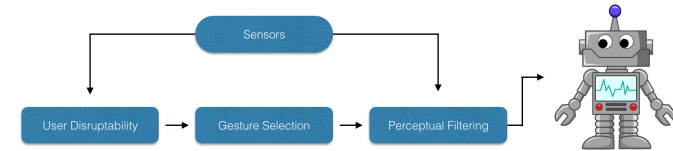
*the minimum amount of gesture that reliably acquires the user's attention so as to not damage the perception and acceptance of the robot.*

The initial focus of this work is in establishing joint attention; however, the work is meant to extend to engaging children in a variety of interactions involving learning and play beyond the screen.

### Attention Acquisition Behavior Generation

We approach the process of actively establishing directed joint attention as a series of three sequential phases: (1) *attention acquisition*, (2) *attention direction*, and (3) *attention sharing* (Figure 1). In this work, we focus on attention acquisition, and specifically on behaviors that aim to acquire the interaction partner's attention.

To generate effective minimal attention acquisition we introduce an approach comprised of four typical HRI modules: (1) perception, (2) user disruptability, (3) gesture selection, and (4) perceptual filtering. *User disruptability* is a learned predictive model of how intrusive an attention acquisition behavior needs to be, based on perception data, in order to successfully acquire the user's attention. *Gesture selection* uses a crowd-sourced general model of relative intensities of gestures in the given gesture space, and a learned



**Figure 2:** Attention Acquisition Behavior Generation Pipeline

relationship between that model and user disruptability. Gesture selection produces a set of appropriate minimal gestures. In the perceptual filtering phase, this set is filtered based on the user's current state relative to the components of the selected gestures.

Figure 2 shows the behavior generation pipeline and the following sections explain each phase of the pipeline.

#### Learning User Disruptability

We represent user disruptability with multivariate Gaussian process regression [11] as a Multivariate General Gaussian Process (MGGP) model [2] of the perceptual inputs, such as the user's on-screen behaviors, head pose, body pose, and gaze. Because disruptability is likely to vary greatly among users, and human annotation of abstract concepts such as disruptability has been shown to be unreliable, we plan to relearn a disruptability model for each user.

We predict that disruptability will be strongly correlated with task difficulty, so by finding peaks in the covariance function of the MGGP, we can autonomously select the difficulty of the next task and in that way personalize the interaction. By manipulating the difficulty of the tasks and collecting data on the user's response to different intensities of attention acquisition gestures, we can sample areas of high variance to minimize variance across the models we are learning.



**Figure 3:** Kiwi the SpriteBot [4]



**Figure 4:** "Code Copying" Activity in Alien Codes from Expedition Deployment built with Phaser led by Eric Deng

### *Learning Gesture Intensities*

There are many parameters of nonverbal gestures that can be adapted by tuning the robot's trajectories [9, 3] and changing properties of the nonverbal sounds that it makes like frequency and amplitude. We represent this multi-modal gesture space as a Riemannian manifold with each dimension representing a parameter of the robot's gesture [9]. This representation lends itself to manifold-related techniques for analyzing and exploring the gesture space and quantitatively comparing gestures based on their relative locations on the manifold.

Prior work has shown that people are consistent in perceiving relative intensities but very inconsistent in perceiving absolute intensities of complex gestures [14]. Thus, we can learn a crowd-sourced, general model of relative robot behavior intensities, resulting in a manifold that represents the perceived disruptability of all gestures in our gesture space.

### *Gesture Selection*

From the components outlined in the previous two sections, we have both a classifier for user disruptability and a relative scale of intensities for nonverbal attention acquisition behaviors. Given those, we can where the two models align for each individual user by holding the relative disruptability of gestures model constant. Using this new relational mapping from one model to the other, together with real-time perceptual data, we can generate a set of gestures of the appropriate disruptiveness for the interaction scenario.

### *Perceptual Filtering*

From the set of produced gestures, the last step of the gesture generation pipeline is to filter these intensity-appropriate gestures to also be "modally-appropriate". Modally-appropriate behaviors are those that are likely to be most clear and salient to the user at the given point in time [13]. Our perception module can return information such as head and

body poses and well as ambient sound and using those data, we can estimate which gestures in our set would be most salient to the user. For instance, if the robot perceives that it is out of the user's field of view, it would use a non-visual action or if there is too much ambient sound, the robot can select a behavior primarily consisting of robot motion.

## **Experimental Validation**

To develop and evaluate the described approach, we will evaluate it with different user populations. We will begin with a convenience population of college-aged adults. Concurrently, we are designing user studies involving children with Autism Spectrum Disorders (ASD) as part of an ongoing research program in our lab funded by the National Science Foundation (NSF). The studies will focus on teaching four target social skills: (1) joint attention, (2) social turn-taking, (3) ordering and sequencing, and (4) perspective taking [15].

We will be using a robot (Figure 3) and tablet setup (Figure 4), and engaging the children in symbolic play, a strategy often used in occupational therapy for children with ASD [8]. The attention acquisition behavior pipeline will be used to initiate joint attention and selectively draw children's attention away from the on-screen activities when appropriate. Because ASD children often have atypical attentional behaviors, the described work can have an impact in a potentially therapeutic context.

## **References**

- [1] Dare A Baldwin, C Moore, and PJ Dunham. 1995. Understanding the link between joint attention and language. *Joint attention: Its origins and role in development* (1995), 131–158.

- [2] Antoni B Chan. 2013. Multivariate generalized Gaussian process models. *arXiv preprint arXiv:1311.0360* (2013).
- [3] Eric Deng and Maja J. Mataric. 2017. Mime-Inspired Behaviors in Minimal Social Robots. In *ACM CHI Workshop on What Actors can Teach Robots*.
- [4] Yifeng Fu Elaine Short, Dale Short and Maja J. Matarić. 2017. *SPRITE: Stewart Platform Robot for Interactive Tabletop Engagement*. Tech report. Department of Computer Science, University of Southern California. <http://robotics.usc.edu/publications/967/>
- [5] David Feil-Seifer and Maja J Matarić. 2005. Defining socially assistive robotics. In *Rehabilitation Robotics, 2005. ICORR 2005. 9th International Conference on*. IEEE, 465–468.
- [6] Terrence Fong, Illah Nourbakhsh, and Kerstin Dautenhahn. 2003. A survey of socially interactive robots. *Robotics and autonomous systems* 42, 3 (2003), 143–166.
- [7] Chien-Ming Huang and Andrea L Thomaz. 2011. Effects of responding to, initiating and ensuring joint attention in human-robot interaction. In *RO-MAN, 2011 IEEE*. IEEE, 65–71.
- [8] Connie Kasari, Stephanny Freeman, and Tanya Paparella. 2006. Joint attention and symbolic play in young children with autism: A randomized controlled intervention study. *Journal of Child Psychology and Psychiatry* 47, 6 (2006), 611–620.
- [9] Jean-Claude Latombe. 2012. *Robot motion planning*. Vol. 124. Springer Science & Business Media.
- [10] Chris Moore and Phil Dunham. 2014. *Joint attention: Its origins and role in development*. Psychology Press.
- [11] Joaquin Quiñonero-Candela and Carl Edward Rasmussen. 2005. A unifying view of sparse approximate Gaussian process regression. *Journal of Machine Learning Research* 6, Dec (2005), 1939–1959.
- [12] Andrea Lockerd Thomaz, Matt Berlin, and Cynthia Breazeal. 2005. An embodied computational model of social referencing. In *Robot and Human Interactive Communication, 2005. ROMAN 2005. IEEE International Workshop on*. IEEE, 591–598.
- [13] Dionysis E Tsoukalas, John Mourjopoulos, and George Kokkinakis. 1997. Perceptual filters for audio signal enhancement. *Journal of the audio Engineering Society* 45, 1/2 (1997), 22–36.
- [14] Yair Weiss, Eero P Simoncelli, and Edward H Adelson. 2002. Motion illusions as optimal percepts. *Nature neuroscience* 5, 6 (2002), 598–604.
- [15] Susan Williams White, Kathleen Keonig, and Lawrence Scahill. 2007. Social skills development in children with autism spectrum disorders: A review of the intervention research. *Journal of autism and developmental disorders* 37, 10 (2007), 1858–1868.