

Effects of 2D and 3D Displays on Turn-taking Behavior in Multiparty Human-Computer Dialog

Samer Al Moubayed Gabriel Skantze

KTH Speech Music and Hearing

Stockholm, Sweden

sameram@kth.se, gabriel@speech.kth.se

Abstract

The perception of gaze from an animated agent on a 2D display has been shown to suffer from the Mona Lisa effect, which means that exclusive mutual gaze cannot be established if there is more than one observer. In this study, we investigate this effect when it comes to turn-taking control in a multi-party human-computer dialog setting, where a 2D display is compared to a 3D projection. The results show that the 2D setting results in longer response times and lower turn-taking accuracy.

1 Introduction

The function of gaze for interaction purposes has been investigated in several studies. Gaze direction and dynamics have been found to serve several different functions, including turn-taking control, deictic reference, and attitudes (Kendon, 1967). Recently, there has been an increasing interest in virtual agents that may engage in multi-party, situated dialogue (e.g., Bohus & Horvitz, 2010). In such settings, gaze may be an essential means to address a person in a crowd, or pointing to a specific object out of many.

It is known that perception of 3D objects that are displayed on 2D surfaces is guided by, what is commonly referred to as, the Mona Lisa effect (Todorovic, 2006). This means that the orientation of the 3D object in relation to the observer will be perceived as constant, no matter where the observer is standing in the room. This effect has important implications for the design of interactive systems, such as embodied conversation agents,

that are able to engage in situated interaction, as in pointing to objects in the environment of the interaction partner, or looking at one exclusive observer in a crowd.

In a previous study (Al Moubayed et al., in press), we have measured how subjects *perceive* gaze direction using an animated agent in 2D and 3D conditions (see Figure 1). The purpose of this study is to investigate how gaze may affect the turn-taking *behavior* of the subjects in a multi-party human-computer dialog, depending on the use of 2D or 3D displays.

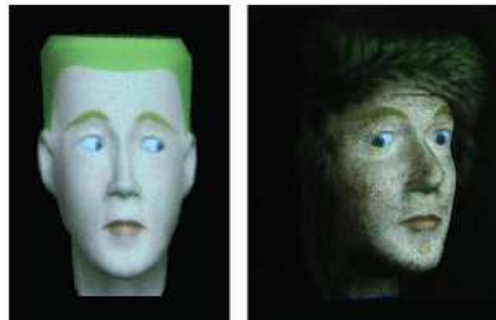


Figure 1: A snapshot of the animated agent projected on a 2D white board (left) and on a 3D head model (right).

2 Method

Two sets of five subjects were asked to take part in the experiment. In each session, the five subjects were seated at fixed positions at an equal distance from each other and from an animated agent. The agent addressed the subjects by directing its gaze in their direction. Two versions of the agent were used, one projected on a 3D head model and one projected on a flat surface (using the same 3D computer model), as shown in Figure 1. The

conversational behavior of the animated agent was controlled using a Wizard-of-Oz setup. For each new question posed by the agent, the gaze was randomly shifted to a new subject. The subjects were given the task of watching a first-person video from a camera navigating around the city of Stockholm, after which the animated agent asked them to describe the route they had just seen. After each video was finished, the animated agent started to ask the subjects about directions on how to reach the landmark the video ended with, starting from the point of view the video started with. The dialogs were in Swedish. An example dialog is shown in Table 1. Each set of subjects did four dialogs in each 2D/3D condition (i.e., a total of eight videos). The ten subjects were not informed about the purpose of the experiment beforehand.

Table 1: The beginning of one of the recorded interactions translated from Swedish.

Agent	Hi, can you help me to find my way to the Finnish church?
Subject 1	Yes
Subject 2	Yes
Agent	Should I continue in this direction?
Subject 3	Yes, a little bit at least
Agent	How far should I go then?
Subject 4	I think it was about 20 meters, maybe a bit longer
Agent	Is there some landmark I should look for?
Subject 2	You should look for the Great Church, when you come to Västerlånggatan you will see the Great Church on your left, then should pass it
Agent	Do you agree?
Subject 5	Yes

3 Results

To measure the efficiency of the gaze control, a confusion matrix was calculated between the intended gaze target and the actual turn-taker. The accuracy for targeting the intended subject in the 2D condition was 53% and 84% for the 3D condition. The mean response time was also calculated for each condition, i.e. the time between the gaze shift of the question and the time takes for one of the subjects to answer. A two sample ANOVA analysis was applied, with the response time as a dependent variable, and the condition as an independent variable. The results show a significant main effect [$F(1)=15.821$, $p<0.001$],

with a mean response-time of 1.85 seconds for the 2D condition, and of 1.38 seconds for the 3D condition. No significant correlation with time was found (Pearson Correlation = -0.094), which means that there is no learning effect on how to perceive the gaze of the agent for either condition.

4 Conclusions

The results show that the use of gaze for turn-taking control on 2D displays is limited due to the Mona Lisa effect. The accuracy of 50% is probably too low in settings where many users are involved. By using a 3D projection, this problem can be avoided to a large extent. However, the accuracy for the 2D condition was higher than what was reported in a previous perception experiment in a similar setting (Al Moubayed et al., in press). A likely explanation for this is that the subjects in this task may to some extent compensate for the Mona Lisa effect – even if they don’t “feel” like the agent is looking at them, they may learn to associate the agent’s gaze with the intended target subject. This comes at a cost, however, which is indicated by the longer mean response time. The longer response time might be due to the greater cognitive effort required making this inference, but also to the general uncertainty among the subjects about who is supposed to answer.

Acknowledgements

This work has been carried out at the Centre for Speech Technology at KTH, and is supported by the European Commission project IURO (Interactive Urban Robot), grant agreement no. 248314, as well as the SAVIR project (Situating Audio-Visual Interaction with Robots) funded by the Swedish Government (strategic research areas).

5 References

- Bohus, D. & Horvitz, E. (2010). Facilitating multiparty dialog with gaze, gesture, and speech. In *Proceedings of ICMI-MLMI*, Beijing, China.
- Al Moubayed, S., Edlund, J., & Beskow, J. (in press). Taming Mona Lisa: communicating gaze faithfully in 2D and 3D facial projections. *ACM Transactions on Interactive Intelligent Systems*.
- Kendon, A. (1967). Some functions of gaze direction in social interaction. *Acta Psychologica*, 26, 22-63.
- Todorovic, D. (2006). Geometrical basis of perception of gaze direction. *Vision Research*, 45(21), 3549-3562.