

# Presentation Agents That Adapt to Users' Visual Interest and Follow Their Preferences

Arjen Hoekstra<sup>1</sup>, Helmut Prendinger<sup>2</sup>, Nikolaus Bee<sup>3</sup>, Dirk Heylen<sup>1</sup>, and Mitsuru Ishizuka<sup>4</sup>

<sup>1</sup> Computer Science, Human Media Interaction  
University of Twente, PO Box 217, 7500 AE Enschede The Netherlands  
[a.h.hoekstra@student.utwente.nl](mailto:a.h.hoekstra@student.utwente.nl), [heylen@cs.utwente.nl](mailto:heylen@cs.utwente.nl)

<sup>2</sup> National Institute of Informatics  
2-1-2 Hitotsubashi, Chiyoda-ku, Tokyo 101-8430, Japan  
[helmut@nii.ac.jp](mailto:helmut@nii.ac.jp)

<sup>3</sup> Institute of Computer Science, University of Augsburg  
Eichleitnerstr. 30, D-86135 Augsburg, Germany  
[nikolaus.bee@gmail.com](mailto:nikolaus.bee@gmail.com)

<sup>4</sup> Graduate School of Information Science and Technology, University of Tokyo  
7-3-1 Hongo, Bunkyo-ku, Tokyo 113-8656, Japan  
[ishizuka@i.u-tokyo.ac.jp](mailto:ishizuka@i.u-tokyo.ac.jp)

**Abstract.** This research proposes an interactive presentation system that employs eye gaze as an intuitive and unobtrusive input modality. Eye movements are an excellent clue to users' attention, visual interest, and preference. By analyzing and interpreting eye behavior in real-time, our system can adapt to the current (visual) interest state of the user, and thus provide a more personalized and 'attentive' experience of the presentation. The system implements a virtual presentation room, where research content is presented by a team of two highly realistic 3D agents in a dynamic and interactive way. A small preliminary study was conducted to investigate users' gaze behavior with a non-interactive version of the system. A demo video based on our system was awarded as the best application of life-like agents at the GALA event in 2006.<sup>5</sup>

## 1 Introduction

The challenge of giving a good presentation is to provide relevant and interesting content in an easily accessible way while keeping the attention of the audience during the entire presentation time. Human presenters often obtain feedback from listeners regarding their level of attention by simply looking at their behavior, specifically whether they are looking at the currently presented material, typically visualized on slides, at the presenter, or somewhere else. If a presenter, e.g. a museum guide, observes that the attention of the spectators is diverted by other objects, he or she will try to adapt the presentation by taking the interest shift of the audience into account.

<sup>5</sup> <http://hmi.ewi.utwente.nl/gala/>



Although speech conveys the richest information in human-computer interaction, it is not the preferred input modality for scenarios such as presentation settings, which, as monologues, typically do not assume verbal expressions of interest from the audience. To determine the user's current focus of attention and interest, we therefore propose a system that is based on human eye movements. As an input modality, eye gaze has the advantage of being an involuntary signal that reflects the user's visual interest [14], and its signal is robust and can be assessed accurately [4].

Our proposed system can be conceived as reviving the 'self-disclosing display' concept introduced in [19], where eye gaze is utilized as an input modality to recognize and respond to a user's interest. Their system would zoom in to areas of user interest and provide explanations via synthesized speech. Our work extends this concept by detecting both user interest and preference between two (visual) alternatives to continue the presentation, and by embodied life-like characters rather than a disembodied voice.

The remainder of this paper is structured as follows. Section 2 discusses related work. Section 3 and 4 describe our methods to assess (visual) interest and (visual) preference, respectively. Section 5 provides details about the application scenario and the gaze-contingent responses of the agents. In Section 6, we report on the main findings of our preliminary study based on a non-interactive version of the agent application. Section 7 concludes the paper.

## 2 Related Work

Life-like characters are virtual animated agents that are intended to provide the illusion of life or 'suspend disbelief' [2], such that users interacting with those agents will apply social interaction protocols and respond to them as they would to other humans, e.g. by listening to their story and attending to them through eye gaze [13]. Life-like characters have been shown to serve multiple purposes successfully; besides presenters, they can act as tutors, actors, personal communication partners, or information experts [12].

Eyes are an intriguing part of the human face, and are sometimes even seen as 'windows to the soul'. For instance, [6] provides a good summary of the major functions of eye gaze, including paying and signaling attention, conversation regulation, and demonstration of intimacy (see also [1, 10]). Early work on the social functions of gaze direction in dyadic communication can be found in [8]. It has recently been generalized to multiple conversational partners [20].

Recent attempts to integrate eye behavior into interactive systems are reported in [17], which discusses the use of eye tracking in various applications - so-called 'visual attentive interfaces' - such as the *Magic Pointing* and *InVision* systems. Whereas *Magic Pointing* is based on the user's conscious eye behavior in order to control a mouse pointer, *InVision* exploits involuntary gaze movements to estimate a user's plan or needs. Similar to *InVision*, our system exploits the non-conscious nature of eye movements in a non-command fashion.



### 3 Interest Estimation

The focus of interest is determined by a modified version of the algorithm described in [14]. These authors implemented an intelligent virtual tourist information environment (*iTourist*), for which they propose a new interest algorithm based on eye gaze. Two interest metrics were developed: (1) the Interest Score (IScore) and (2) the Focus of Interest Score (FIScore). IScore refers to the object ‘arousal’ level, i.e. the likelihood that the user is interested in a (visual) object. When the IScore metric passes a certain threshold, the object is said to become ‘active’. The FIScore calculates the amount of interest in an active object over time.

Since we were mainly interested in whether a user’s attention is currently on a particular object, a simplified version of the IScore metric was sufficient for our purpose. The basic component for IScore is  $p = T_{ISon}/T_{IS}$ , where  $T_{ISon}$  refers to the accumulated gaze duration within a time window of size  $T_{IS}$  (in our application, 1000 milliseconds). In order to account for factors that may enhance or inhibit interest, [14] characterize the IScore as  $p_{is} = p(1 + \alpha(1 - p))$ . Here,  $\alpha$  encodes a set of parameters that increase the accuracy of interest estimation.

The modification factors are modelled as follows [14]:

$$\alpha = \frac{c_f \alpha_f + c_c \alpha_c + c_s \alpha_s + c_a \alpha_a}{c_f + c_c + c_s + c_a}$$

The terms in this formula are defined as:

- $\alpha_f$  is the frequency of the user’s eye gaze ‘entering’ and ‘leaving’ the object ( $0 \leq \alpha_f \leq 1$ ),
- $\alpha_c$  is the categorical relationship with the previous active object ( $\alpha_c = -1|0|1$ ),
- $\alpha_s$  is the average size of all possible interest objects compared to the size of the currently computed object ( $-1 \leq \alpha_s \leq 1$ ),
- $\alpha_a$  encodes whether the object was previously activated ( $\alpha_a = -1|0$ ), and
- $c_0$ ,  $c_f$ ,  $c_c$ ,  $c_s$ , and  $c_a$  represent empirically derived constant values of the corresponding factors. Some of these factors are domain dependent and are thus not applicable in all contexts.

The factors  $\alpha_c$  and  $\alpha_a$  were not (yet) integrated to our system.  $\alpha_c$  concerns (semantic) relations between objects;  $\alpha_a$  can be used to make the system respond in a different way when an object is activated multiple times.

We continue by explaining  $\alpha_f$  and  $\alpha_s$ , the two remaining factors.  $\alpha_f$  is represented as  $\alpha_f = \frac{N_{sw}}{N_f}$ , where  $N_{sw}$  denotes the number of times eye gaze enters and leaves the object and  $N_f$  denotes the maximum possible  $N_{sw}$  in the preset time window. When the user’s gaze switches to some object many times, the value of the modification factor will increase and hence there will be a higher chance on excitation.  $\alpha_s$  is represented by  $\alpha_s = \frac{S_b - S}{S}$ , whereby  $S_b$  represents the average size of all objects,  $S$  denotes the size of the currently computed object, and the smallest object is never more than twice as small as the average object.



This modification is intended to compensate for the differences between the size of the potential interest objects. Due to some noise in the eye movement signal, larger objects could have a higher chance of being ‘hit’ than smaller ones, which should be avoided.

## 4 Preference Estimation

In order to determine the user’s preference in situations involving a two-alternative forced choice (2AFC), i.e. “how the presentation should continue”, we exploited the so-called ‘gaze cascade’ effect. This effect was discovered in a study where users had to choose the more attractive face from two faces [18]. It could be demonstrated that there was a distinct gaze bias towards the chosen stimulus in the last one and a half seconds before the decision was made.

Our system integrates a recently developed real-time component for automatic visual preference detection, the *AutoSelect* system, which is based on the gaze cascade phenomenon [3]. *AutoSelect* was tested in a study where users were instructed to choose their preferred necktie from two presented neckties, i.e. in a 2AFC setting. There was no input modality available except the subjects’ eye gaze. After the decision of *AutoSelect*, subjects were asked to confirm (or reject) the result of the system. Starting from an initial set of thirty-two pairs of neckties, subjects repeatedly indicated their preference, amounting to sixty-two decisions. The system achieved an accuracy of 81%.

Examples of the exploitation of the gaze cascade effect and of the use of the interest algorithm will be given in the next section.

## 5 Responding to User Interest and Preference

Our implemented system involves a team of two presentation agents that introduce the user to research at the National Institute of Informatics (NII), Tokyo (see Fig. 1 and video<sup>6</sup>). The two agents were designed based on the appearance of two famous Japanese actors. In order to support their life-likeness, the agents are highly expressive. They can perform various gestures, such as greeting and counting, or ‘beat’ and deictic gestures. In addition to body gestures, mimics for joy, surprise, and sadness are available. High-quality synthesized speech is combined with proper lip synchronization,<sup>7</sup> and the head of the agents can be adjusted to any (natural) direction, e.g. to the direction of the other agent when giving turn, or to the virtual slide. The agents and environment are controlled by MPML3D [11], a reactive framework that supports anytime interaction, such as real-time interpreted input from the eye tracker.

The agents will adapt their performance based on user eye gaze in two ways:

<sup>6</sup> A demo video can be found at <http://research.nii.ac.jp/~prendinger/GALA2006/>

<sup>7</sup> When listening to a presentation, paying attention to its visualized content is of key importance. However, the audience will also focus on the presenter’s face to increase comprehension via perception of lip movements in addition to speech, especially when listeners are not native speakers of English, as in our case.



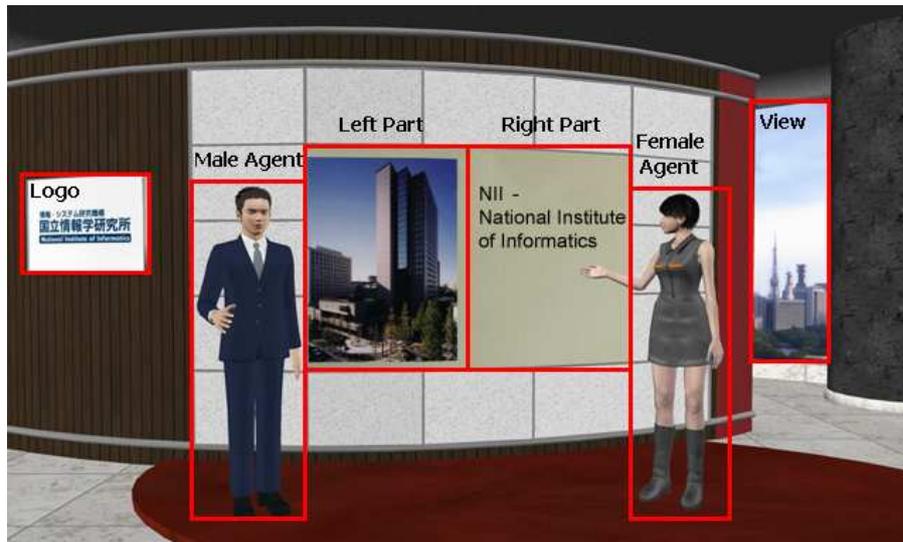


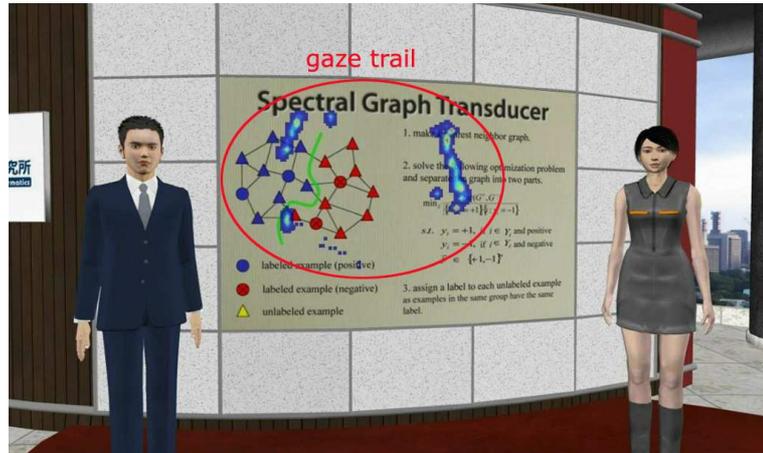
Fig. 1. Interest objects in the virtual environment.

- If the user shows interest in a particular interface object (an ‘interest object’) not currently discussed (e.g. the view), or non-interest in a currently discussed object (e.g. a presentation slide), the agents will interrupt their presentation and react accordingly.
- At decision points in the presentation flow, the user’s preference determines the subsequent topic.

### 5.1 Adapting to User Interest

In the system, the following interest objects are defined (see Fig. 1; from left to right): (a) NII logo; (b) male agent; (c) left part of the slide; (d) right part of the slide; (e) female agent; (f) the view out of the window to the right. For each interest object, the IScore is calculated every second. When the score exceeds the threshold, the object becomes ‘activated’ and the agent(s) will react (if a reaction is defined). Agent responses (or non-responses) are defined for three types of situations:

1. *Continuation of presentation*: If the user attends to the currently explained (part of a) presentation slide (which is desired), the agent will continue with the presentation. Fig. 2 depicts a situation where the user attends to the explanation of the male agent by gazing at the slide content.
2. *Interruption of presentation*: If the user is detected to be interested in an interface object that is not currently discussed, the system chooses between two responses:



**Fig. 2.** User is interested in slide content. The corresponding gaze trail is visualized by ‘heat trails’.

- (a) Suspension: If e.g. the user looks out of the virtual window (at the “View” object) rather than attending to the presentation content explained by the male agent, the female co-presenter agent asks her colleague to suspend the research presentation and continues with a description of the view.
- (b) Redirecting user attention: Here, the presenter agents do not suspend the presentation to comply with the user’s interest. Instead, the co-presenter alerts the user to focus on the presentation content.

The existing implementation of our presentation system handles interruptions in a simple way. If a user’s interest object is not the currently explained object (typically a slide) the presentation will be suspended at first by providing information about that object, and subsequently, the co-presenter agent will try to redirect the user to the presentation content.

## 5.2 Following User Preference

At predefined points during the presentation, the agents ask the user to choose the next presentation topic, while a slide depicting two options is displayed. The gaze cascade phenomenon will occur naturally in this situation. Users alternately look at the left part and the right part of the slide, and eventually exhibit a bias for one part. The decision process occurs within seven seconds. Thereafter, the presentation continues with the selected topic.

## 6 Exploratory Study

A small study was conducted to assess users’ eye behavior when watching a non-interactive version of the research presentation by the agent team, i.e., although



**Fig. 3.** Experimental setup.

eye gaze was recorded, the agents did not adapt to user interest or preference. This approach seemed justified as a first step, given the lack of experience with attentive behavior of human spectators of a presentation performed by two animated agents. Hence, the aim of the study was to assess likely types of gaze behaviors. This information can then be used to refine the functionality of the interactive system, which will be followed by an extensive study.

## 6.1 Method

*Subjects:* The data of four subjects (average 30 years) were analyzed. Subjects received a monetary compensation for participation (1,000 Yen).

*Apparatus and Procedure:* Subjects were seated in front of a 30 inch screen (distance 80 cm) and stereo cameras of the faceLAB eye tracker from Seeing Machines.<sup>8</sup> The cameras and speakers were located below the screen. Two infrared pods were attached at the upper part of the display for illumination of the eyes (see Fig. 3). Then calibration of each subject was performed. Subjects were given no instruction other than watching the presentation.

In the presentation prepared for the study, the agents first introduce themselves, and then explain the researches of three professors of NII. The total length of the presentation is 14:49 min.

*Data Analysis:* The eye tracking software of faceLAB allowed us to extract the coordinates of gaze points on the screen. The recorded data was then processed and analyzed with MATLAB. 'Heat trails' (similar to 'hotspot' maps [15]) were used for visualization, as they present the amount of fixations over time as a

<sup>8</sup> <http://www.seeingmachines.com/>

continuous movie. The heat trails were made transparent with the chroma key (Bluescreen) effect and merged with the captured video of the presentation. The algorithm described in [5] was adjusted to calculate fixations, using a velocity-based algorithm [16]. Animations and (virtual) environment changes were analyzed with the ANVIL annotation tool [9].

## 6.2 Results

The most distinctive result of the study could be found for situations where the agents ask the subject to select the subsequent topic. All of the subjects showed the gaze pattern characteristic of the ‘gaze cascade’ effect in both occurrences of a decision situation. This outcome generalizes the results of [18, 3] to a setting featuring two agents referring to slide content depicting two choices (displayed left and right on the slide). It indicates that the cascade phenomenon can be reliably used to let users select the continuation of the presentation in a non-command fashion.<sup>9</sup> It should be noted, however, that in the non-interactive presentation shown in the study, the subjects’ preference had no effect on the continuation of the presentation.

Deictic arm gestures of embodied agents and agents’ head-turning to slide content are an effective way to redirect the attention of users to the (virtual) slides [13]. We were interested in differences in the effect of deictic gestures depending on whether a new slide is shown, or some textual content of a displayed slide is changed, e.g. a new item is added to a given slide content. In the study, every subject had noticed a new slide within 2 sec (19 new slides presented). On the other hand, changes on slides (18 occurrences) were noticed with some delay, with 97% redirected attention within 3 sec. Although we expected more occasions where an attentive agent would have to alert the user, a 15 min presentation is probably too short to observe a user’s diminishing attention.

The functionality of the interactive system also provides for the possibility that users attend to interface objects not related to the presentation content, such as the NII logo or the view outside the building (see Fig. 1). In the study, however, subjects spent 99% of the total time on the agents or slides. Since the actual view of the subjects was essentially limited to those interface objects (see Fig. 3), there was little room for attending to anything else. Other results regarding cumulative gaze distribution include attention to speaking agent (53%), attention to presented slides (43%), and attention to non-speaking agent (3%).

## 7 Conclusions

The analysis of eye gaze offers a powerful method to adapt a presentation to a user’s interest, to alert the user in case of distraction, and to estimate the user’s preference. Eye gaze as an input modality is particularly beneficial when verbal feedback is either not assumed or difficult to provide. Most importantly, the

<sup>9</sup> Given the small sample size, our results should always be seen as preliminary.



estimation of eye behavior is an unobtrusive method to estimate user interest continuously.

While gaze-contingent interfaces are getting increasingly popular [4], it remains an open question how ‘reactive’ an interface that uses eye gaze as an input should be. The problem of distinguishing between eye movements that are just explorative and those that are meaningful as an input is known as the ‘Midas Touch’ problem: “Everywhere you look, another command is activated; you cannot look anywhere without issuing a command.” [7, p. 156]. Our presentation system avoids the Midas touch problem by (1) strictly confining the screen areas that could yield an agent response (the interest objects), and (2) calculating user interest based on a well-established metric [14].

We have described an interactive presentation system that features two highly realistic and expressive virtual 3D agents capable of responding to a user’s focus and shift of attention and interest in a natural way. In case of interest estimation, the system relies on a previously developed algorithm [14]. User preference estimation is realized by an automated version of the ‘gaze cascade’ effect [3], building on findings from neuroscience [18]. The exploratory study performed with a non-interactive version of the system indicates that this phenomenon occurs naturally when subjects are asked to choose their preferred continuation of the presentation. An open issue is how to handle situations where the system fails to estimate the user’s interest or preference correctly. Currently, the system does not provide a means to ‘undo’ a decision. We leave this problem for future research.

The interactive presentation system was successfully shown at the *NII Open House 2006* event for two days. A video clip based on the system recently won an award as the best application of life-like agents.

## Acknowledgements

The research was supported by the Research Grant (FY1999–FY2003) for the Future Program of the Japan Society for the Promotion of Science (JSPS), by a JSPS Encouragement of Young Scientists Grant (FY2005–FY2007), and an NII Joint Research Grant with the Univ. of Tokyo (FY2006). The first author was supported by the JSPS Encouragement Grant. The third author was supported by an International Internship Grant from NII under a Memorandum of Understanding with the Faculty of Applied Informatics at the Univ. of Augsburg. We would also like to thank Dr. Ulrich Apel (NII) for scripting the dialogues.

## References

1. M. Argyle and M. Cook. *Gaze and Mutual Gaze*. Cambridge University Press, Cambridge, 1976.
2. J. Bates. The role of emotion in believable agents. *Communications of the ACM*, 37(7):122–125, 1994.



3. N. Bee, H. Prendinger, A. Nakasone, E. André, and M. Ishizuka. AutoSelect: What You Want Is What You Get. Real-time processing of visual attention and affect. In *Tutorial and Research Workshop on Perception and Interactive Technologies (PIT-06)*, pages 40–52. Springer LNCS 4021, 2006.
4. A. T. Duchowski. *Eye Tracking Methodology: Theory and Practice*. Springer, London, UK, 2003.
5. R. Engbert and R. Kliegl. Microsaccades uncover the orientation of covert attention. *Vision Research*, 43(9):1035–1045, 2003.
6. D. Heylen. Head gestures, gaze and the principles of conversational structure. *International Journal of Humanoid Robotics*, 3(3):241–267, 2006.
7. R. J. K. Jacob. The use of eye movements in human-computer interaction techniques: What You Look At is What You Get. *ACM Transactions on Information Systems*, 9(3):152–169, 1991.
8. A. Kendon. Some functions of gaze-direction in social interaction. *Acta Psychologica*, 26:22–63, 1967.
9. M. Kipp. *Gesture Generation by Immitation – From Human Behavior to Computer Character Animation*. PhD thesis, Saarland University, 2004. Dissertation.com, Boca Raton, Florida.
10. C. L. Kleinke. Gaze and eye contact: A research review. *Psychological Bulletin*, 100(1):78–100, 1986.
11. M. Nischt, H. Prendinger, E. André, and M. Ishizuka. MPML3D: a reactive framework for the Multimodal Presentation Markup Language. In *Proceedings 6th International Conference on Intelligent Virtual Agents (IVA-06)*, Springer LNAI 4133, pages 218–229, 2006.
12. H. Prendinger and M. Ishizuka, editors. *Life-Like Characters. Tools, Affective Functions, and Applications*. Cognitive Technologies. Springer Verlag, Berlin Heidelberg, 2004.
13. H. Prendinger, C. Ma, and M. Ishizuka. Eye movements as indices for the utility of life-like interface agents: A pilot study. *Interacting with Computers*, 2006. In press.
14. P. Qvarfordt and S. Zhai. Conversing with the user based on eye-gaze patterns. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI-05)*, pages 221–230. ACM Press, 2005.
15. M. Russell. Using eye-tracking data to understand first impressions of a website. *Usability News*, 7(1), 2005.
16. D. D. Salvucci and J. H. Goldberg. Identifying fixations and saccades in eye-tracking protocols. In *Proceedings of the Eye Tracking Research and Applications Symposium*, pages 71–78, New York, 2000. ACM Press.
17. T. Selker. Visual attentive interfaces. *BT Technology Journal*, 22(4):146–150, 2004.
18. S. Shimojo, C. Simion, E. Shimojo, and C. Scheier. Gaze bias both reflects and influences preference. *Nature Neuroscience*, 6(12):1317–1322, 2003.
19. I. Starker and R. A. Bolt. A gaze-responsive self-disclosing display. In *Proceedings CHI-90*, pages 3–9, 1990.
20. R. Vertegaal, R. Slagter, G. van der Veer, and A. Nijholt. Eye gaze patterns in conversations: There is more to conversational agents than meets the eyes. In *Proceedings of CHI-01*, pages 301–308. ACM Press, 2001.

