# GazeProjector: Location-independent gaze interaction on and across multiple displays

**Christian Lander, Sven Gehring, Antonio Krüger**
German Research Center for Artificial Intelligence (DFKI)
{christian.lander, sven.gehring, krueger}@dfki.de

**Sebastian Boring**
University of Copenhagen
sebastian.boring@diku.dk

**Andreas Bulling**
Max-Planck Institute for Informatics
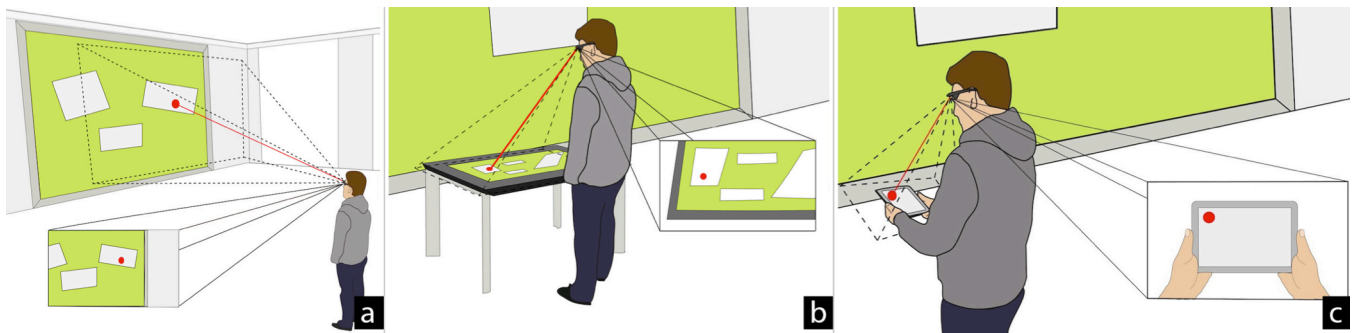bulling@mpi-inf.mpg.de

**Figure 1.** *GazeProjector* **enables seamless gaze-based interaction with multiple displays from arbitrary locations and orientation, such as wall-sized displays (a), horizontal interactive screens (b), and handheld devices (c) – without active recalibration.**

## ABSTRACT

Mobile gaze-based interaction with multiple displays may occur from arbitrary positions and orientations. However, maintaining high gaze estimation accuracy still represents a significant challenge. To address this, we present *GazeProjector*, a system that combines accurate point-of-gaze estimation with natural feature tracking on displays to determine the mobile eye tracker's position relative to a display. The detected eye positions are transformed onto that display allowing for gaze-based interaction. This allows for seamless gaze estimation and interaction on (1) multiple displays of arbitrary sizes, (2) independently of the user's position and orientation to the display. In a user study with 12 participants we compared *GazeProjector* to existing well-established methods such as visual on-screen markers and a state-of-the-art motion capture system. Our results show that our approach is robust to varying head poses, orientations, and distances to the display, while still providing high gaze estimation accuracy across multiple displays without re-calibration. The system represents an important step towards the vision of pervasive gaze-based interfaces.

## Author Keywords

Eye tracking; gaze estimation; calibration; large displays; multi-display environments; natural feature tracking

## ACM Classification Keywords

H.5.m. Information interfaces and presentation (e.g., HCI): Miscellaneous.

## INTRODUCTION

Gaze is a powerful modality for interacting with the rapidly increasing number of public displays around us. Gaze naturally indicates what we visually attend to and what we are interested in [28], and is faster than other pointing devices such as a mouse [21]. Consequently, gaze-based interaction received considerable attention with applications ranging from controlling desktops [31, 13], eye typing [17], to target selection [26], and password entry [7]. Recent advances in mobile eye tracking point the way towards unobtrusive interfaces that will allow us to interact with gaze in everyday settings [6]. A key challenge with such pervasive gaze-based interfaces is how spontaneous and transparent gaze-based interaction can be facilitated on arbitrary displays.

Monocular head-mounted eye trackers are typically equipped with two cameras: a scene camera that captures part of the user's current field of view, and an eye camera that records a close-up video of the user's pupil position and eye movements. Such eye trackers have to be calibrated to a specific user for a specific display before first use. The calibration establishes a mapping between 2D pupil positions and 2D positions in the scene camera's coordinate system. For gaze-based interaction with displays in the environment, these scene camera coordinates have to be mapped further to corresponding 2D gaze positions on the displays. During operation, arbitrary on-screen gaze positions can then be estimated by interpolating between known calibration points.

The problem with head-mounted eye trackers is that this calibration is typically performed for a fixed position and orientation of the user to a particular display. While this is less of an issue for stationary settings and TV-sized displays, mobile everyday-life settings and multiple – potentially large – displays evoke two types of motion: (1) user movements in front of a single display to inspect other parts of the display's content; and (2), head movements to reach targets outside of the ocular motor range [10]. In addition, there might be multiple displays present, evoking further movements. These types of motion considerably reduce gaze estimation accuracy [8]. This problem has been addressed by tracking the eye tracker relative to a particular display by augmenting the environment with visual markers [30, 5] or using vision-based motion capturing systems (e.g., OptiTrack[1]). Although these approaches provide high tracking accuracy, they are impractical for spontaneous gaze interaction, particularly when interacting with multiple displays in mobile everyday-life scenarios. Similarly, marker-free trackers, such as the Kinect, work well for tracking coarse user position but are not accurate enough to track head orientation and fail when people occlude each other.

In this paper we present *GazeProjector*, a system that allows for accurate gaze estimation on arbitrary displays independently of the user's position and orientation (see Figure 1). *GazeProjector* requires a one-time calibration of the head-mounted eye tracker per user, so that the user's pupil position can be mapped to positions in the scene camera's coordinate system. The calibration can be performed on *any* display and not necessarily on the display the user wants to interact with. Afterwards, the system automatically transforms this calibration to other positions and displays in real-time. In contrast to previous approaches, *GazeProjector* does not require a heavyweight motion capturing system or visual markers. Instead, the eye tracker continuously tracks itself relative to different displays using natural feature tracking. *GazeProjector* therefore enables seamless gaze estimation and interaction across different displays, and allows users to freely move around in front of them while maintain a high gaze estimation accuracy.

In a controlled laboratory experiment with 12 participants, we compared our approach to a state-of-the-art OptiTrack motion capturing system as well as a marker-based approach. In the first task, participants looked at on-screen targets from various positions and orientations in front of a large display. In a second task, we compared *GazeProjector* to the marker-based approach on multiple displays (here: a wall-sized display, a tabletop, and a tablet PC). In both tasks, we found that our approach compensated well for head movements (i.e., change of orientation) and user relocation (i.e., change of location).

**RELATED WORK**

---

Our work builds on methods for (1) gaze approximation and estimation on displays, (2) gaze interaction using head-mounted eye trackers, as well as (3) tracking the spatial relationship between users and displays.

**Gaze Approximation and Estimation on Displays**
Several previous works used head orientation as an approximation of where people look. For example, Sippl et al. used a remote camera to detect facial features, such as eyes and nose tip, and estimate head pose on four areas on the display [22]. Nakanishi et al. relied on a stereo face tracking system and the 3D head pose as an approximation of gaze direction [19]. Finally, ViewPointer aimed to detect eye contact between users and devices using a wearable camera and IR tags placed in the environment [24]. While useful for coarse attention measurements, none of these approaches allowed for accurate gaze estimation on the display.

Accurate gaze estimation on displays remains a significant challenge – particularly when remote eye trackers (i.e., eye trackers placed in front of a display) are used. Such trackers only allow a single user to interact with a display at any point in time and any interaction is restricted to the tracking range of typically 50-80 cm in a central area in front of the display, thereby severely limiting users' mobility [23, 26]. Previous work either focused on extending the tracking range of remote trackers [11, 18], or on calibration-free (spontaneous) interaction but was either limited to interaction along a horizontal axis, i.e., without full 2D gaze estimation [32] or required dynamic interfaces [29]. None of these works addressed the problem of users interacting with the display from different positions and orientations.

**Gaze Interaction Using Head-mounted Eye Trackers**
Head-mounted eye trackers are more flexible as they allow the user to move freely in front of the display. Early work on using head-mounted eye trackers for interaction still required calibration to a single, stationary display prior to first use [9]. More recent approaches aimed to estimate gaze dynamically but either required visual markers to be attached to the display [30] or in the environment to detect gaze on a set of pre-defined interaction areas, e.g., to control a TV set or music player [5].

With advances in computer vision, visual markers could be substituted with detecting the display directly in the scene camera's field of view. Mardanbegi et al. described an approach to detect screens based on quadrilaterals found in the scene [16]. Turner et al. extended that work to multiple displays (based on the displays' aspect ratios) by adding a second camera as well as a method for transparently switching between two calibrations [27]. However, in contrast to *GazeProjector*, both approaches required that the display was fully visible in the scene camera's fields of view, which is problematic when users are mobile or close to a very large display. Furthermore, relying on a set of automatically selected key points and features instead of screen borders is more robust to changing light conditions and generalizes better to displays of arbitrary shape and size.

**Tracking Spatial Relationships of Users and Displays**
Tracking the spatial relationship of users (and the users' devices respectively) can be done in two ways. First, external tracking equipment can be used to determine a device's exact position in 3D space (and thus its spatial relationship to a display in the environment). The Proximity Toolkit makes use of such high-precision tracking equipment and provides an interface to acquire spatial relationships [33]. While such a setup results in extremely high accuracy, it is often impractical for outdoor use.

Alternatively the device's camera can be used to identify its spatial relationship to a display. Many approaches exist, such as temporarily showing on-screen visual markers [2] or using dynamic markers following a camera's position [20]. More recently, natural feature tracking was used to determine spatial relationships. Herbert et al. used Scale-Invariant Feature Transform (SIFT) to determine the camera's spatial relationship to a display [12]. Their system tried to identify a screenshot of the display in the device's camera stream. Virtual Projection extended this approach to dynamically updated displays [3]. Touch Projector further allowed for tracking multiple displays provided that display content differs sufficiently [4].

**ENABLING GAZE INTERACTION ON LARGE DISPLAYS**
As mentioned before, estimating a user's gaze on a large display and in multi-display environments using a head-mounted eye tracker faces two challenges: the eye tracker has to be calibrated and used from fixed positions and orientations for all displays. Furthermore, during calibration and use the entire display has to be visible in the eye tracker's scene camera. Ideally, the eye tracker only has to be calibrated once. This can be achieved by (1) calibrating pupil positions to the scene camera coordinate system; and (2) tracking the spatial relationship between the eye tracker and a specific display and (3) mapping 2D gaze positions in scene camera coordinate space to that display.

**Eye Tracker Calibration**
*GazeProjector* uses a one-time calibration to map pupil positions to the scene camera's coordinate system. Because of this, there is no need to perform the calibration on the display one intends to interact with. Instead, the system can be calibrated once on any display in the environment (e.g., a laptop). This independence of the target display has two advantages: The calibration is not dependent on the distance and/or orientation to a display as this is handled by the self-localisation directly; and the calibration does not depend on a single display, thus allowing for seamless gaze estimation across several displays in multi-display environments.

**Tracking the Spatial Relationship to Displays**
To determine the spatial relationship between the eye tracker and a specific display we use the approach described in [3]. Specifically, our system streams the scene camera's video to a server that is aware of all displays in the environment as well as their displayed contents. All displays in the environment repeatedly stream screenshots to the server

to reflect their current content (i.e., in case of quickly updated content, such as videos). The server is thus only aware of the physical dimensions of each display (i.e., size and resolution) as well as their current content, but not their physical location. This is especially important for mobile devices, which frequently change their position and orientation over time.

The server then processes the incoming screenshots as well as incoming frames from the field camera using FAST feature detectors [15] and FREAK feature descriptors [1]. The idea is to use current screenshots as *template* images, which the server tries to find in the *observed* images (here: the field camera's video). If a template matches an observed image, the algorithm calculates the transformation matrix (i.e., a homography), which describes the transformation of points from one image plane (say: a video frame) into another image plane (say: the display's screenshot).

**Gaze Estimation**
As mentioned before, the transformation matrix allows for mapping locations in the scene camera coordinate system to the coordinate system of the target display and vice versa. In this procedure, the display does not have to be visible in full in the scene view. Instead, a sub-region is sufficient given enough features within that region to allow for robust tracking. Similar to Touch Projector (where touch positions are transformed), *GazeProjector* finally uses the transformation matrix to estimate gaze positions on the display.

**Implementation**
Our system consists of three components: (1) a monocular head-mounted eye tracker[2] connected to a laptop [14]; (2) a $10m^2$ back-projected display wall; and (3) a desktop computer driving the display. Laptop and desktop computer are connected via Wifi. The eye tracking software on the laptop is written in Python and is based on PUPIL's open source mobile eye tracking platform[2]. The software running on the desktop computer is written in C# (.NET Framework 4.5). For feature detection, description and matching, we use EmguCV[3] as wrapper for OpenCV[4]. For faster processing, we downscale display screenshots to $384 \times 240$ pixels and camera frames to $320 \times 240$ pixels.

The system allows for distances ranging from 0.5 times the display's diagonal up to six times the display's diagonal. When being further away, the accuracy decreases as the display observed in the camera's field of view decreases in size (thus, removing several features). We believe that a multi-scale approach of screenshots will increase the operational range, yet we decided not to include it in this proof-of-concept implementation. In addition, the tracking compensates for an angular offset of ±60°. While this is sufficient for most interactions, fast eye/head movements will have a slight impact on accuracy. However, we believe that

the increasing processing capabilities of future devices will allow for both faster image processing on larger images
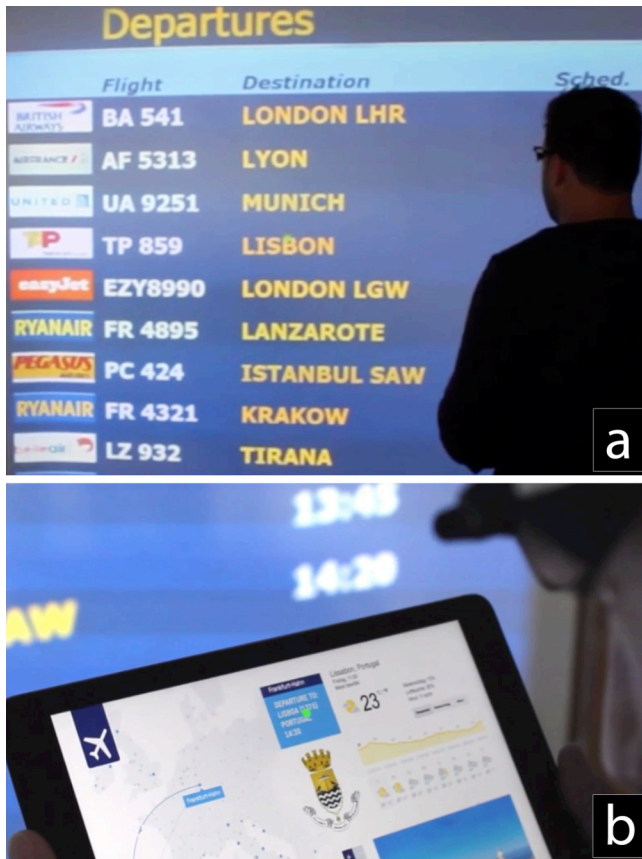


**Figure 2. Our example application: when users select from a time table of departing flights (a), the corresponding information is shown on their mobile device.**

(i.e., less or no scaling required) for higher accuracy.

## Example Application
We built an example application to demonstrate the use of *GazeProjector* in a distributed multi-display environment (e.g., an urban public space). We chose such a setting, as it

underlines how people can make use of continuous gaze estimation with reasonable accuracy while freely moving around in the environment.

Consider a timetable on one of the many large screens at an airport, showing flights departing within the next hours. Users can look at a specific flight and its flight number or destination respectively. Once the system recognizes the point users look at, additional information, such as a picture of the destination or detailed flight data, is transferred to the user's mobile device (e.g., a tablet PC). *GazeProjector* further allows for tracking gaze also on the tablet, allowing for content adaptation. If the user now gazes at the picture, the tablet will show more specific information, such as the weather at the destination. Figure 2 shows the example application, and the video figure explains it in more detail.

## EXPERIMENT I: ASSESSING GAZE ESTIMATION
We conducted a controlled laboratory study to assess *GazeProjector*'s gaze estimation accuracy in comparison to existing but more heavyweight approaches.

### Independent Variables
We had two independent variables in this experiment: *Mode* (i.e., the gaze estimation method used), and *Location* (i.e., where participants stood in front of the display).

*Mode:* We chose three different modes for gaze estimation: *GazeProjector* (*GP*) implemented as described before; *Marker Tracking*[5] (*MT*), which uses a set of on-screen markers for tracking the orientation between the eye tracker and the display provided by the PUPIL framework; and a simple *Head Orientation* (*HO*) approach, which tracks the participant's head using an external OptiTrack system. For each of these modes, we calibrated the eye tracker from two different locations. Both locations, however, were placed centrally in front of the display with one location being close to the display and one being further away. We further calibrated the eye tracker for each participant separately instead of using one calibration (see limitations section for further details).

*Location:* We chose six different locations in front of the display to simulate a more realistic setting. Three of these locations were close to the display and three were further away. The eye tracker was only calibrated for the central *near* and *far* central locations. This is more realistic, as users would not calibrate for every position in a walk-in-and-use scenario. Note that we calibrated the eye tracker for each participant separately instead of using one calibration (see limitations section for further details). Since no visual feedback was given to them and to keep the experiment at a reasonable length, participants had to perform the set of tasks only once. We then computed the gaze estimation accuracy post-hoc for each of the calibrations.

### Task & Procedure
We implemented a gaze pointing task in which participants had to fixate nine different target locations represented as red circles on the display (see Figure 3). A pilot study showed that participants were affected by visualizing their gaze point on the display. Especially if the gaze position was incorrect, people tended to "move" the gaze point to compensate for the error. We therefore opted to not provide any visual feedback to the participants. Participants were instructed to look at each target as quickly and accurately as possible. Each target location was shown for five seconds.

For each *Mode*, participants first calibrated from the *near-center* location and performed the tasks for all other locations. Afterwards, the calibration for the *far-center* location was recorded and gaze positions as well as errors were evaluated post-hoc. Following best practices in gaze estimation experiments, we validated all calibrations by asking

[5] http://www.pupil-labs.com/blog/2013/12/036-release.html

participants to fixate once on each point on the same 9-point pattern. At the end of the study we asked for demographic information.

We collected gaze data from the eye tracker and transformation matrices calculated by *GP* as well as *MT*. Furthermore, we recorded data about the head position and orientation with *OptiTrack*. Data was sampled at 30 Hz (i.e., 150 samples per on-screen target) leading to a total of 1350 samples for each *Mode* and *Location* combination. We discarded samples for which participants' pupil was not detected. We dropped the first two seconds (60 samples) for each target, which was the maximum time required to find the target All together we dropped 276.985 samples. In total we recorded 306.215 samples, 140.532 for *GP*, 165.683 for *MT*, which was the basis sample set used for *HO*.

### Experimental Design
We used a within-subject design with the independent variables *Mode* (*GP-near, GP-far, MT-near, MT-far, HO-near, HO-far*) and *Location* (*front-left, front-center, front-right, back-left, back-center, back-right*).

We counterbalanced the order of *Location* across participants using a Latin Square. Although it is possible to record all position information in parallel, we opted to have *GP* and *MT* separate, as markers would favor *GazeProjector*'s tracking. The *HO* mode was recorded while participants were using the *MT* mode. Half of our participants started with the *MT*, and the other half with *GP*. Thus, each participant performed the task twice per location. For each mode and location, the nine targets (equally distributed in a 3 × 3 grid on-screen) were presented in random order.
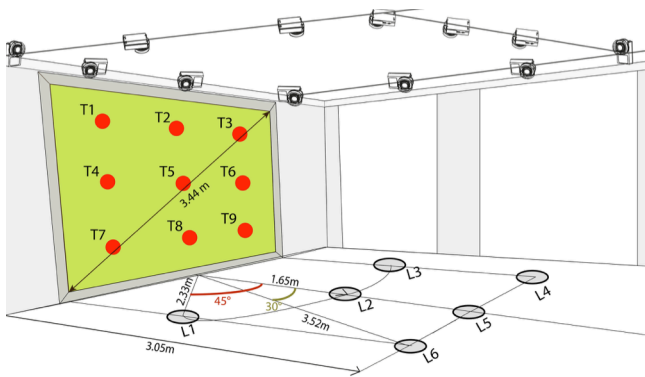


**Figure 3. Experimental setup showing all locations L1-L6 and orientations relative to the display, as well as the nine different positions T1-T9 of the visual targets on the display.**

### Apparatus
Figure 3 shows our experimental setup: we used a large front-projected wall with a size of 2.75 × 2.07 meters (diagonal: 3.44 meters). The six locations were distributed within a nine square meter area in front of the display as follows: three locations at a distance of 1.65 m (*near*), and three locations at a distance of 3.05 m (*far*). The left and right locations for *near* were exactly 2.33 meters away from the

display's centerline (i.e., an angular offset of ±45°); those for *far* were located 3.52 m away from the display's centerline (i.e., an angular offset of ±30°). Naturally, the two center locations for *near* and *far* had an angular offset of 0°. Locations located *far* allow participants to observe the entire screen at once (the display covers 48.52°), while for locations located *near* the display covers 79.60° – thus exceeding the full-scale ocular motor range of ±55° [10]).

### Participants
Twelve participants (three female) between 22 and 32 years (mean = 27.45 years, SD = 3.1 years) were recruited from a local university campus. All participants had normal or corrected to normal vision; none reported any form of visual impairments (e.g., color blindness).

### GAZE ESTIMATION RESULTS
We corrected all reported gaze estimation accuracies by subtracting the mean calibration error. The mean calibration error was 2.04° (SD = 0.69°). To verify that we could do so, we performed a one-way ANOVA with a Bonferroni-corrected post-hoc analysis on calibration accuracies across all *Modes*, and found no significant differences. In all subsequent post hoc analyses, we used Bonferroni-corrected confidence intervals to retain comparisons against $\alpha = 0.05$. Furthermore, we used Greenhouse-Geisser correction in cases where sphericity had been violated.

### Gaze Estimation Error
To assess the gaze estimation error, we calculated the average gaze estimation error in degrees of visual angle. That is, the difference of the visual angle between the predicted on-screen gaze point and the actual fixation targets for all *Modes* and *Locations*. We then performed a 6 × 6 (*Mode × Location*) within subjects ANOVA on gaze estimation errors and found a main effect for *Mode* ($F_{1.989,21.879} = 8.526$, $p < .002$), a main effect for *Location* ($F_{5,55} = 7.363$, $p < .001$), but we did not find an interaction between the two.

We performed post-hoc tests to further understand the main effect of *Mode*. Most importantly, we found significant differences within *MT* and *HO* for the two calibrations near and *far* (all $p < .033$). In both cases, the *near* calibration led to lower estimation errors. *GP*, on the other hand, did not show such an effect, suggesting that the point of calibration does not effect its gaze estimation error significantly, and the difference in means was also lower than for the other two (*GP*: 0.281°; *MT*: 0.931°; *HO*: 0.948°) – yet, also for *GP*, the mean estimation errors were slightly lower for the *near* calibration than for the *far* one.

This is further reflected when comparing across *Modes*: *GP-near* differed significantly from both *MT-far* and *HO-far* (all $p < .01$). However, there was no significant difference between the *Modes* for the *near* calibration. Furthermore, *GP-far* did not differ significantly from any other *Mode* despite having relatively large differences in error.

Overall, *GP-near* showed the lowest error (*M* = 1.80°, SD = 0.20°), followed by *GP-far* (*M* = 2.08°, SD = 0.27°), and

*HO-near* (*M* = 2.09°, *SD* = 0.23°). Also *MT-near* (*M* = 2.23°, *SD* = 0.31°) has an estimated gaze error less than 3 degrees.



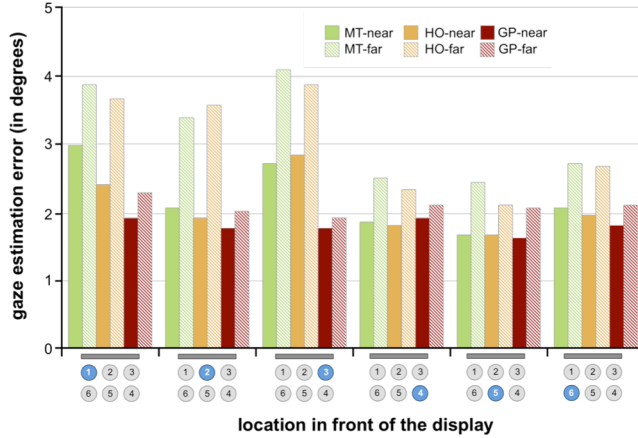**Figure 4. Mean gaze estimation error for every location for MT-near, MT-far, HO-near, HO-far, GP-near and GP-far.**

The other two *Modes* performed slightly worse and had an error of more than 3 degrees, *MT-far* (*M* = 3.16°, *SD* = 0.32°) and *HO-far* (*M* = 3.04°, *SD* = 0.31°). Figure 4 summarizes theses results.

Post-hoc tests on *Location* revealed that the significant main effect stems from whether participants were farther away from the display or not: *front-left* differed significantly from *back-center* and *back-right* (all *p* < .019). *Front-right* also differed significantly from *back-center* (*p* < .011). Overall, *back-center* led to the least estimation errors (*M* = 1.93°, *SD* = 0.23°), followed by *back-right* (*M* = 2.01°, *SD* = 0.17°), and *back-left* (*M* = 2.22°, *SD* = 0.30°). The *front* locations performed worse with *front-center* resulting in the least estimation errors (*M* = 2.45°, *SD* = 0.28°), followed by *front-left* (*M* = 2.86°, *SD* = 0.29°) and *front-right* (*M* = 2.86°, *SD* = 0.30°). Thus, on average, the *back* locations had a lower estimation error of 2.08° (*SD* = 0.23°) compared to the *front* locations with 2.72° (*SD* = 0.29°).



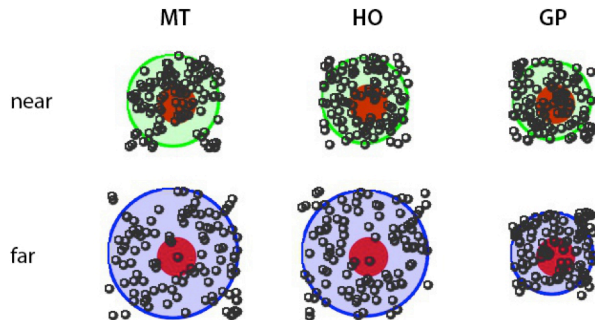**Figure 5. Visualization the mean gaze error (ellipses) for the three modes MT, HO and GP and all calibrations averaged over all targets. Additionally the mean gaze points are visualized by black circles.**

## Differences for On-screen Target Positions
We did not expect high gaze estimation errors for each of the *Modes*. However, we wanted to analyze whether the on-screen targets resulted in different estimation errors and thus analyzed the results separately for each on-screen target. For *MT*, we found no significant main effects on gaze estimation error for *Target*. We found the same for *HO*. Only for *GP* we found significant differences for gaze estimation for *Target*. Our analysis revealed that predominantly the *bottom-left* target T7 differed significantly from few others (T2, T3, T6 and T8) and led to higher estimation errors. We assume that this is due to the scene camera seeing too few features, which in turn increased the error of the transformation matrix. Figure 5 shows gaze estimation errors for the different modes averaged over all targets.

## Eye and Head Movements
We were further interested in whether participants mainly moved their head or their eyes to point at an on-screen target location. As expected [9], we found that the average normalized gaze position in the field camera's video was $x = 0.44$ and $y = 0.47$ ($SD_x = 0.21$; $SD_y = 0.25$). Thus, gaze positions remained near the center of the participants' field of view. We subsequently analyzed the gaze position for every *Location* in front of the display and found no significant differences between them. The largest average difference was 0.03. Table 1 lists these results for each *Location*.

| Location | Mean (x,y) | SD (x,y) | Var (x,y) |
|---|---|---|---|
| *front-left* | 0.43,0.45 | 0.19,0.24 | 0.038,0.060 |
| *front-center* | 0.45,0.47 | 0.20,0.25 | 0.040,0.062 |
| *front-right* | 0.46,0.46 | 0.22,0.27 | 0.052,0.076 |
| *back-left* | 0.46,0.48 | 0.23,0.25 | 0.054,0.065 |
| *back-center* | 0.45,0.48 | 0.20,0.24 | 0.044,0.058 |
| *back-right* | 0.43,0.48 | 0.20,0.23 | 0.043,0.057 |

**Table 1. The table shows the mean, standard deviation and variance for the x,y-coordinates of normalized gaze positions in the participants' field of view.**

The *OptiTrack* data provided detailed information of participants' head orientation (*HO*). We found that the largest head turns covered the entire width of the display (*far*: 51.2°, *near*: 83.66°). On average head motions covered an angle of 31.61° (*SD* = 2.04°). This further confirms our results in that *HO* might be a suitable approximation for gaze estimation with an average error of 2.09° (*SD* = 0.23°) for *HO-near* and 3.04° (*SD* = 0.31°) for *HO-far*.

## EXPERIMENT II: MULTIPLE DISPLAYS
We conducted a second controlled laboratory study to assess *GazeProjector's* gaze estimation accuracy across multiple displays of varying form factors – with only a single calibration performed on one of the displays.

## Independent Variables
We had two independent variables in this experiment: *Mode* (i.e., the gaze estimation method used), and *Screen* (i.e., on which display was the target shown). We did not use fixed positions as we wanted to create a more realistic scenario in

which participants were free to move in the space between the displays.

*Mode:* In this experiment we chose to use only *GazeProjector* (*GP*) and *Marker Tracking* (*MT*), but not head orientation, as we believe it will perform similarly across displays. We again calibrated for two locations (as in the first experiment), but additionally recorded calibrations on a 40" tabletop display (*Surface*), as well as on a 9.7" iPad Air tablet (*iPad*). We chose to do so to investigate the effects on gaze estimation accuracy of calibrating (1) on surfaces not orthogonal to the participant, and (2) on personal devices with a considerably smaller display. Particularly the latter resembles a more realistic scenario in which users calibrate their eye tracker once on their personal device. As in experiment I, calibrations were analyzed post hoc.

*Screen:* In addition to the large display used in the first experiment (*Wall*), we chose to add the other two displays used for calibration as well (here: *Surface*, and *iPad*).

## Task & Procedure

The task used in this experiment was the same as in the first one: participants had to fixate on-screen targets. However, since we had three displays, participants now had to acquire nine targets per display (27 in total) as shown in Figure 6 As mentioned before, participants could freely choose and change their position between the displays. We again opted to not provide any feedback to participants for the same reasons as before. Participants were instructed to look at each target as quickly and accurately as possible. Each target location was shown for ten seconds to give the participants enough time to find the target on the correct display. There was only *one* target on *one* display shown at a time.

The procedure in this experiment was nearly the same as for the first experiment but with an additional calibration for *Surface* and *iPad* after all tasks were completed. On the additional displays we used the same 9-point calibration pattern. At the end of the study we asked for demographic information.

We collected the same gaze data from the eye tracker as well as the transformation matrices from *GP* and *MT* as in the first experiment. Data was sampled at 30 Hz (i.e., 300 samples for each target, 8100 samples for each *Mode*), and samples were discarded if the participants' pupil was not detected. As we expected an increase in search time for the target, we dropped the first five seconds (150 samples) for each target. In total we recorded 259,745 samples: 124,421 for *GP*, for 135,324 for *MT*.

## Experimental Design

We used a within-subject 8 *Mode* (*GP-near*, *GP-far*, *GP-Surface*, *GP-iPad*, *MT-near*, *MT-far*, *MT-Surface*, *MT-iPad*) × 3 *Screens* (*Wall, Surface, iPad*) design. Half of our participants started with *GP*, the other half with *MT* (as in experiment I). The targets were randomized, thus the next target appeared on any of the three *Screens*. The 27 targets

were again placed in $3 \times 3$ grids (i.e., nine per display) on each display. In total, participants acquired 54 targets.

## Apparatus

We used the same front-projected *Wall* as in the first experiment. In addition, we had a 40" Microsoft Surface 2 tabletop display (*Surface*), and a 9.7" iPad Air tablet (*iPad*). Figure 6 shows the setup including the tabletop and the tablet PC. The tabletop display was placed in front of the projection wall in an area where the participant would occlude the beamer projection. Participants held the tablet in hand during the experiment. They could freely choose their location within a nine square meter area.
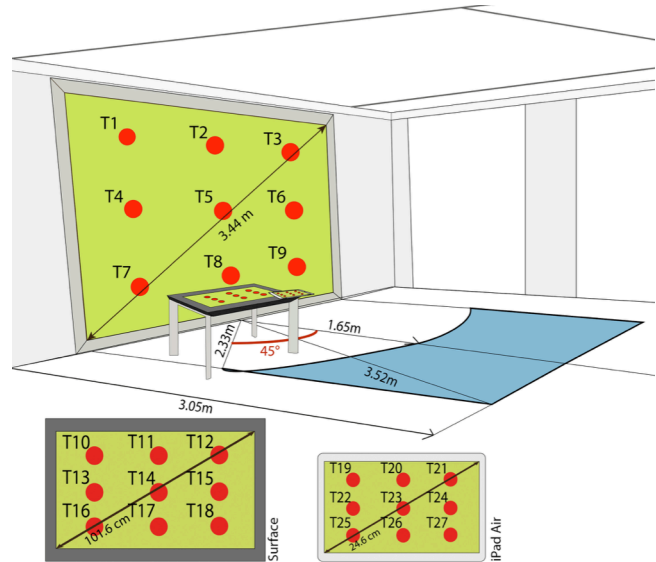


**Figure 6. Experimental setup showing the three used screens and the nine targets on each screen (wall, surface tabletop, iPad). Additionally the area (blue) is marked, where the user was free to choose the location.**

## MULTI-DISPLAY RESULTS

We again corrected gaze estimation accuracy by subtracting the mean calibration error. The mean calibration error was 2.18° (*SD* = 0.69°). We again verified that we could do so by performing an ANOVA with a Bonferroni-corrected post-hoc analysis on calibration accuracies across all *Modes*, and found no significant differences. As in experiment I, we used Bonferroni-corrected confidence intervals in all post hoc analyses and Greenhouse-Geisser correction in cases where sphericity had been violated.

## Gaze Estimation Error

We calculated the average gaze estimation error as in experiment I and subsequently performed a $8 \times 3$ (*Mode* × *Screen*) within subjects ANOVA on them. We found main effects for *Mode* ($F_{7,77} = 21.733$, $p < .001$), and for *Screen* ($F_{2,22} = 82.705$, $p < .001$) as well as an interaction effect between the two ($F_{14,154} = 9.100$, $p < .001$).

Post-hoc pairwise multiple means comparisons revealed that *GP-near* and *GP-far* differed significantly from *MT-*

*far, MT-Surface* and *MT-iPad (*all *p* < .001*)*. Furthermore, *GP-Surface* differed significantly from *MT-Surface* and *MT-iPad* (all *p* < .007). And finally, *GP-iPad* also differed significantly from *MT-far*, *MT-Surface*, and *MT-iPad* (all *p* < .039). It is noteworthy, however, that both *GP* and *MT* did not show any significant differences between their different calibrations, suggesting that the device on which they were calibrated on did not impact accuracy.

Overall, *GP-near* had the lowest estimation error ($M$ = 2.77°, $SD$ = 0.20°), followed by *GP-far* ($M$ = 3.01°, $SD$ = 0.16°), *GP-iPad* ($M$ = 3.24°, $SD$ = 0.17°) and *GP-Surface* ($M$ = 3.31°, $SD$ = 0.16°) across all *Screens*. For all *MT* variations, the estimated gaze errors were larger than 4 degrees. Figure 7 summarizes theses results.

As for the main effect for *Screen*, post-hoc multiple means comparisons revealed that *Wall* was significantly different from the other two *Screens* (all *p* < .001). However, there was no significant difference between *Surface* and *iPad*. Overall, targets on the *Wall* had the least estimation error ($M$ = 2.07°, $SD$ = 0.07°), followed by *Surface* ($M$ = 4.52°, $SD$ = 0.22°) and *iPad* ($M$ = 5.12°, $SD$ = 0.23°).
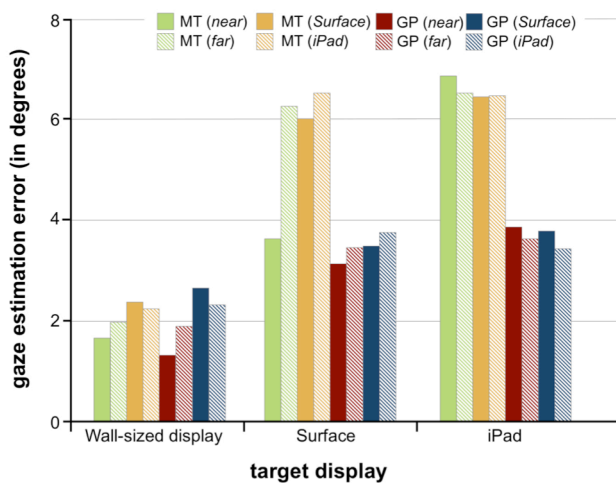


**Figure 7. Mean gaze estimation error for every target screen for MT-near, MT-far, MT-surface, HO-ipad, GP-near, GP-far, GP-surface and GP-ipad.**

Upon inspecting Figure 7, one can see that the source of the *Mode × Screen* interaction is the increased difference between *MT* and *GP* (all calibration modes) between the *Wall* and *Surface*/*iPad*, with the *Wall* resulting in much lower estimation errors than the other two. It is noteworthy, that *MT-near* performs similarly to all *GP* modes on the *Surface*, but its estimation error increases drastically on the *iPad*, although all *GP* modes remain at their level. We subsequently ran separate ANOVAs on *Modes* for each *Screen*, and found several significant effects. On the *Wall*, only *MT-iPad* and *GP-near* differed significantly (*p* < .008) indicating that nearly all modes performed similarly.

On the *Surface*, the differences become more prominent, with *GP-near* and *GP-far* outperforming all *MT* modes

except *MT-near* (all *p* < .016). Furthermore, *GP-Surface* and *GP-iPad* differed significantly from *MT-iPad* (all *p* < .012). We found the most differences on the *iPad*, where all *GP* modes are significantly less error-prone than all *MT* modes (all *p* < 0.03). Here we again did not find any differences within *GP* and *MT* for different calibration modes. Figure 8 visualizes the mean gaze estimation errors for each of the three screens and its nine targets for the two modes *MT* and *GP*.
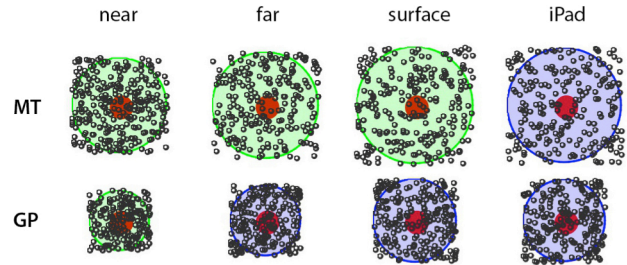


**Figure 8. Visualization of the mean gaze error (ellipses) for all different modes, MT and GP, and all calibrations averaged over all targets over all screens. Additionally the mean gaze estimations are visualized by black circles.**

## DISCUSSION & LIMITATIONS
Our results show that – on a single display – *GazeProjector* achieves an average gaze estimation accuracy of 1.78° compared to 2.64° for *MT*, and 2.65° for *HO*. When used on multiple displays (and only being calibrated on a single screen), *GazeProjector* achieves an average gaze estimation accuracy of 2.47° compared to 3.60° for *MT* over all modes and target screens. Although this accuracy is slightly lower than the 0.5°–1° reported for the PUPIL eye tracking glasses under ideal conditions (i.e., in a stationary desktop setting with a 27" screen and optimal lighting conditions [14]), we achieve this accuracy in a fully unconstrained, pervasive interaction setting.

### Pervasive Settings
The first advantage of *GazeProjector* is its suitability for pervasive gaze interaction settings [6]. Current approaches that allow for gaze interaction on multiple displays using monocular mobile eye trackers require heavyweight external motion capturing systems or visual markers. While motion capture systems allow for high-precision tracking, they are (1) costly and (2) cannot easily be installed in public environments. Markers reduce this, but have another drawback: all displays have to be augmented with them – either with printed ones attached to a display's frame [30, 5], or digital ones shown on the display. However, printed markers quickly clutter the environment, in particular in settings with a large number of displays. While digital markers could only be shown on demand, they still take away display space and "compete" with the main content.

While binocular systems can automatically compensate for vergence error, estimating gaze in display coordinates still requires to track changes in the user's position and orienta-

tion relative to these displays. This severely limits the use of these devices to instrumented environments. *GazeProjector*, however, allows users to interact from arbitrary locations and orientations relative to multiple displays without this need – and, as our experimental results show, *GazeProjector* does so without lowering accuracy. Thus, our approach allows for unconstrained and seamless gaze interaction with multiple displays while on the move.

### Display Visibility
Display tracking using visual markers requires the whole target display to be visible in the eye tracker scene camera's field of view during calibration and interaction. In contrast *GazeProjector* relies on natural feature tracking and provides competitive gaze estimation accuracy even if only a fraction of the target display is visible. Naturally, the larger the visible portion of the display, the lower the tracking error. However, we found that a quarter of the display is usually sufficient, provided that enough features (e.g., high frequencies) are found in that portion. This allows *GazeProjector* to work on much larger displays as well as with more extreme head movements than current eye trackers.

### Multi-Display Interactions
Our results show that *GazeProjector* provides robust gaze estimation accuracy for different displays of different sizes without a need for recalibrating the eye tracker to each of the displays. Instead, the eye tracker only needs to be calibrated once (on *any* display) and gaze estimates are then automatically mapped to the other displays during runtime. Applying our calibration method in the presented experiments, we were still able to achieve an accuracy of 3.24° when the eye tracker had been calibrated on a 9.7" iPad Air screen. This is a significant advancement over state-of-the-art gaze estimation approaches.

### Head Movement and Orientation
We further found that head movements are more prevalent in gaze interaction with large displays compared to smaller displays (e.g., mobile devices). This finding is in line with controlled laboratory studies on human vision: humans employ head movements for gaze shifts with ocular orbital eccentricity exceeding 20° [25]. While head movements pose a significant challenge to current head-mounted eye trackers, *GazeProjector* proved to be robust to head movements, which is an essential feature for using head-mounted eye tracking systems for large screens.

### Limitations
Despite its numerous advantages over state-of-the-art eye-tracking systems, *GazeProjector* also comes with some limitations: first, our current implementation requires continuous snapshots of the target displays to be transferred to a central server. Consequently, all displays need to be registered with such a server a priori. Furthermore, increasing the number of displays also increases the network load for transferring real-time updates of a display's content. However, we believe that future network technologies may overcome this limitation.

Second, *GazeProjector*'s gaze estimation accuracy depends on the quality of the image data. This concerns the scene camera: as these cameras usually come with a wide angle lens to cover as much as possible of a user's field of view, a target display may be rather small within the image. As mentioned before, this may increase errors in the transformation matrix due to insufficient image features. This is, of course, a technical limitation, which can be overcome by using different lenses for scene cameras.

Nevertheless, we believe that *GazeProjector* is a promising alternative that realizes continuous gaze-based interaction in pervasive settings and multi display environments.

## CONCLUSIONS & FUTURE WORK
In this paper, we presented *GazeProjector*, an approach for accurate gaze estimation and seamless interaction with multiple large displays using head-mounted eye trackers. In contrast to existing systems, *GazeProjector* only requires a single calibration performed with an arbitrary display and is robust to the user's location and orientation to the displays as well as head movements. Furthermore, *GazeProjector* works without external tracking equipment, such as motion capturing systems or markers attached to display.

We conducted two experiments in which we compared *GazeProjector* to existing, well-established techniques (which require additional equipment), and found that our approach compares well to these techniques. When being used on multiple displays, the results are even more promising. Overall, our results underline the significant potential to finally bring gaze-based interaction into pervasive settings that involve gaze interaction with multiple displays.

We tested *GazeProjector* in a laboratory environment to gain first insights into its performance compared to existing techniques. However, we want to take our approach one step further. One obvious step is to take it to the real world and evaluate its performance on (1) ultra-large displays, such as media façades, and (2) do so with multiple users simultaneously. This will further add to the eye tracking community as it has virtually been impossible to test eye tracking systems in such large scales.

## ACKNOWLEDGMENTS

## REFERENCES
1. Alahi, A., Ortiz, R. and Vandergheynst, P. FREAK: Fast retina keypoint. *Proc. CVPR* 2012, 510-517

2. Ballagas, R., Borchers, J., Rohs, M. and Sheridan, J. The smart phone: A ubiquitous input device. *IEEE Pervasive Computing* 5(1), 70-77, 2006

3. Baur, D., Boring, S. and Feiner, S. Virtual projection: exploring optical projection as a metaphor for multi-device interaction. *Proc. CHI 2012*, 1693-1702.

4. Boring, S., Baur, D., Butz, A., Gustafson, S. and Baudisch, P. Touch projector: mobile interaction through video. *Proc. CHI 2010*, 2287-2296

5. Breuninger, J., Lange, C., and Bengler, K. Implementing Gaze Control for Peripheral Devices. *Proc. PETMEI 2011*, 3-8

6. Bulling, A., and Gellersen, H. Toward Mobile Eye-based human-computer interaction. *IEEE Pervasive Computing* 9(4):8-12, 2010

7. Bulling, A., Alt, F., and Schmidt, A. Increasing the security of gaze-based cued-recall graphical passwords using saliency masks. *Proc. CHI 2012*, 3011–3020

8. Cerrolaza, J. J., Villanueva, A., Villanueva, M. and Cabeza, R. Error characterization and compensation in eye tracking systems. *Proc. ETRA 2012*, 205-208

9. Eaddy, M., Blasko, G., Babcock, J. and Feiner, S. My own private kiosk: Privacy-preserving public displays. *Proc. ISWC 2004*, 132-135

10. Guitton D. and Volle M. Gaze control in humans: eye-head coordination during orienting movements to targets within and beyond the oculomotor range. *Journal of Neurophysiology* 58:427–459, 1987

11. Hennessey, C. and Fiset, J. Long Range Eye Tracking: Bringing Eye Tracking into the Living Room. *Proc. ETRA 2012*, 249-252

12. Herbert, L., Pears, N., Jackson, D. and Olivier, P. Mobile device and intelligent display interaction via scale-invariant image feature matching. *Proc. PECCS 2011*

13. Jacob, R. J. K. What you look at is what you get: eye movement-based interaction techniques. *Proc. CHI 1990*, 11–18

14. Kassner, M., Patera, W. and Bulling, A. Pupil: An Open Source Platform for Pervasive Eye Tracking and Mobile Gaze-based Interaction. *Adj. Proc. UbiComp 2014*, 1151-1160.

15. Lowe, D.G. Object recognition from local scale-invariant features. *Proc. ICCV 1999*, 1150-1157

16. Mardanbegi, D. and Hansen, D.W. Mobile Gaze-based Screen Interaction in 3D Environments. *Proc. NGCA 2011*, 2:1--2:4

17. Majaranta, P., and Kari-Jouko R. Twenty years of eye typing: systems and design issues. *Proc. ETRA 2002*, 15-22

18. Model, D. and Eizenman, M. A General Framework for the Extension of the Tracking Range of User-Calibration-Free Remote Eye-Gaze Tracking Systems, *Proc. ETRA 2012*, 253-256

19. Nakanishi, Y., Fujii, T., Kiatjima, K., Sato, Y. and Koike, H. Vision-based face tracking system for large displays. *Proc. UbiComp 2002*, 152–159

20. Pears, N., Jackson, D.G. and Olivier, P. Smart phone interaction with registered displays. *IEEE Pervasive Computing* 8(2), 14-21, 2009

21. Sibert, L. E. and Jacob, R. J. K. Evaluation of eye gaze interaction. *Proc. CHI 2000*, 281–288

22. Sippl, A., Holzmann, C., Zachhuber, D. and Ferscha, A. Real-time gaze tracking for public displays. *Proc. AmI 2010*, 167–176

23. San Agustin, J., Hansen, J. P. and Tall, M. Gaze-based interaction with public displays using off-the-shelf components. *Adj. Proc. Ubicomp 2010*, 377–378

24. Smith, J. D., Vertegaal, R. and Sohn, C. ViewPointer: lightweight calibration-free eye tracking for ubiquitous handsfree deixis. *Proc. UIST 2005*, 53-61

25. Stahl, J. S. Amplitude of human head movements associated with horizontal saccades. *Experimental Brain Research* 126.1: 41-54, 1999

26. Stellmach, S. and Dachselt, R. Look & touch: gaze-supported target acquisition. *Proc. CHI 2012*, 2981-2990

27. Turner, J., Bulling, A. and Gellersen, H. Extending the visual field of a head-mounted eye tracker for pervasive eye-based interaction. *Proc. ETRA 2012*, 269-272

28. Vertegaal, R. Attentive user interfaces. *Communications of the ACM* 46(3):30–33, 2003

29. Vidal, M., Bulling, A. and Gellersen, H. Pursuits: Spontaneous Interaction with Displays based on Smooth Pursuit Eye Movement and Moving Targets. *Proc. UbiComp 2013*, 439-448

30. Yu, L., and Eizenman, E. A new methodology for determining point-of-gaze in head-mounted eye tracking systems. *IEEE Transactions on Biomedical Engineering* 51(10):1765-1773, 2004

31. Zhai, S., Morimoto, C. and Ihde, S. Manual and gaze input cascaded (MAGIC) pointing. *Proc. CHI 1999*, 246–253

32. Zhang, Y., Bulling, A. and Gellersen, H. Sideways: A Gaze Interface for Spontaneous Interaction with Situated Displays. *Proc. CHI 2013*, 851-860

33. Marquardt, N., Diaz-Marino, R., Boring, S., and Greenberg, S. The proximity toolkit: prototyping proxemic interactions in ubiquitous computing ecologies, *Proc. UIST 2011*, 315-326.