Gaze Tracking for Near to Eye Displays

Timo Pylvänäinen timo.pylvanainen@nokia.com Toni Järvenpää toni.jarvenpaa@nokia.com Ville Nummela ville.nummela@nokia.com

Lixin Fan lixin.fan@nokia.com Nokia Research Center P.0. Box 1000, FI-33721 Tampere

Abstract

A novel gaze tracking principle for near-to-eye displays, based on collimated reference lights, is introduced. The paper presents pupil and reference light reflection detection algorithms based on maximizing the posterior probability. The method incorporates prior knowledge learned during tracking. A simple calibration procedure for the system is described. The system is perfectly suited for virtual and augmented reality applications due to its light weight and computational efficiency and because it can be used with a see-through display.

1. Introduction

Near-to-Eye Displays (NED), such as the one used in our prototype depicted in Figure 1, can bring a high-resolution display to a compact mobile device. Typically the magnified virtual image from the micro-display of a NED is perceived larger than the device itself. NEDs have appeared in science fiction as sunglasses type of wearable displays. In reality, commercial products have barely been mobile, due to limitations on size and viewing ergonomics.

Diffractive optical elements on planar waveguides have been proposed as miniaturized, good quality and ergonomically acceptable NED optics. The transparent nature of the waveguides also enable a see-through mode for mobile applications [4, 5]. The immersive quality of the NED calls for a private and fluent input method to interpret user's intentions and to facilitate smooth interaction with the device. Many alternatives, such as mobile keyboard/keypad/mouse, voice controls, hand/finger motion, gesture recognition and head tracking have been proposed. The subconscious nature of gaze, however, opens up new possibilities for intelligent systems that can work with the user without explicit interaction.

Eye tracking technology has been available for many



Figure 1. The eye tracking camera is visible in (a). The brighter areas on the right eye display are the reference light outputs. The small NTSC camera is visible in (c), glued behind the right eye display. A design for the enclosure of the system is shown in (d).

years using a variety of methods (*e.g.* [1, 8, 6]), but these approaches are not readily applicable for NED due to special ergonomics requirements. This paper presents a near-to-eye gaze tracking approach which reconstructs (up to scale) the exact camera-eye-geometry. Notably, the reconstruction is invariant to camera location relative to the eye, and thus the user can freely adjust the NED glasses without the need for recalibration. In a calibration stage, two parameters related to eye anatomy i.e. *pupil depth* and *optical axis offset* (two angles) are determined for different users. This calibration procedure is required only once for each user. The presented approach has achieved robust real time gaze tracking on a laptop PC and is lightweight enough for even commercial mobile devices.

The paper is organized as follows. Section 2 investigates

existing video-based eye tracking approaches. Section 3 describes the basic principle behind the gaze tracking. Sections 5 and 6 describe the computer vision algorithms used to detect relevant image features from the eye-camera. Section 4 briefly describes a calibration method for the system. Section 7 presents a test results on the accuracy of the system and finally Section 8 discusses future development and applications of the system.

2. Related Work

Many video-based eyetracking systems [1, 8, 6] use an off-axis infrared (IR) source to illuminate the eye such that the pupil is the darkest region in comparison to sclera, iris, and eye lids etc. The pupil center is detected by fitting an ellipse to the pupil contour, which corresponds to high intensity gradient pixels in the captured eye images. The first-surface specular reflection of the IR source is detected and used to compute the vector difference between the pupil center and the reflection point. In the calibration stage, a homographic mapping between locations in the world objects and the vector difference is determined from at least four correspondences. The derived homographic mapping is then used to determine users gaze point from the vector difference in consecutive frames [6].

The homographic mapping approach relies on two assumptions: (1) the pupil centers and reflection points across different frames are all on the same plane, i.e. a planar eye surface assumed; (2) the reflection points are fixed across different frames. The planar eye surface approximation is sufficiently accurate, as in [6], if the distance from camera to eye is large in relative to cornea radius. For near-to-eye display, however, the average camera-eye distance is small ($\approx 2 - 3cm$) and this approximation is expected to be too crude. Most importantly, for NED glasses the fixed reflection point assumption becomes invalid when the camera location relative to the eye is changed, as users quite often adjust the display for better viewing experience.

3. Near-to-Eye Gaze Tracking Principle

The principle of the near-to-eye gaze tracker is based on the reconstruction of the camera-eye-geometry. A general theory of gaze estimation based on pupil center and corneal reflections is presented in [2]. A special case not discussed in [2] is when the light sources are at infinity. This special case is enabled by the diffractive optics used in our system and has some special properties that – unlike the solution in [2] – lead to a closed form solution for gaze estimation.

Up to scale reconstruction of the camera-eye-geometry is made possible by two reference lights provided by special diffractive optics fitted to the display optics. Diffractive plates transfer the display image for both eyes. In a similar fashion, a single IR-light beam is split and transferred to two output areas. The IR-light exits the diffractive plates as two collimated wave fronts at two distinct angles. We will call the two collimated wave fronts reference lights. The principle is illustrated in Figure 2.



Figure 2. The reference lights provide two specular reflections that can be used to determine the cornea center. The optical axis passes through the cornea and pupil centers. Gaze angle follows the optical axis of the eye at an offset angle.



Figure 3. Given the pupil depth r_p , the pupil center is determined by the intersection of the sphere of radius r_p centered at cornea center and the ray from the camera at the angle at which the pupil center was observed.

ICAT 2008 Dec. 1-3, Yokohama, Japan ISSN: 1345-1278 Assume a left-handed coordinate system, where the camera is at the origin facing along the positive z-axis. Let vectors $\vec{l_1}$ and $\vec{l_2}$ be unit vectors representing the directions of the reference lights as in Figure 3. The reference lights are fixed relative to the camera. The eye-observing camera sees the reference lights as specular reflections on the cornea surface. Let $\vec{s_i}$ be the location of the specular reflection due to $\vec{l_i}$ in homogeneous 2D coordinates. Given the camera calibration matrix K_c , the reflection vector $\vec{r_i}$ is given by

$$\vec{r_i} = K_c^{-1} \vec{s_i} \vec{r_i} = \frac{r_i'}{\|\vec{r_i}\|}.$$
(1)

Notice that $\vec{r_i}$ are in inhomogeneous 3D coordinates.

The normal of the cornea surface at the reflection point, normalized to unit length, is then given by the law of reflection as

$$\vec{n_i} = \frac{\vec{l_i} + \vec{r_i}}{\|\vec{l_i} + \vec{r_i}\|}.$$
(2)

Here, the normal is towards the center of the cornea.

The cycle from the camera center, through the first specular reflection to the center of the cornea and back via the second specular reflection is described by

$$\alpha \vec{r_1} + \vec{n_1} - \vec{n_2} - \beta \vec{r_2} = 0.$$
(3)

Here, a unit length for the cornea radius is assumed. The scalars α and β represent the unknown distances to the cornea surface along the rays $\vec{r_1}$ and $\vec{r_2}$. This can be rearranged as

$$R\begin{bmatrix}\alpha\\\beta\end{bmatrix} = \vec{n_2} - \vec{n_1},\tag{4}$$

where $R = [\vec{r_1} - \vec{r_2}]$. In the ideal case, there is an exact solution and the two points $\alpha \vec{r_1} + \vec{n_1}$ and $\beta \vec{r_2} + \vec{n_2}$ are the same point: the center of the cornea. In the presence of noise, however, Equation 4 does not have a solution in general.

Assuming Gaussian noise in the measurements $\vec{s_i}$, the maximum likelihood solution is given by

$$\min_{\vec{c}} \sum_{i=1}^{2} \|(\vec{s_i} - f(\vec{l_i}, \vec{c})\|,$$
(5)

where the function $f(\vec{l}, \vec{c})$ returns the location of the specular reflection from a unit sphere centered at \vec{c} given the light direction \vec{l} . Notice that f returns the camera projection of the specular reflection and is therefore also dependent on the camera calibration K_c . It could also take into account any distortion models of the camera.

The optimization problem 5 can be solved with standard non-linear optimization methods, such as the Levenberg-Marquardt algorithm. The function f, however, appears to

be quite tricky to compute efficiently and is generally done by iterative numerical methods [10]. In practice, this approach may be too slow for real time implementations.

A faster and simpler approach is, of course, to take the minimum norm solution given by

$$\begin{bmatrix} \alpha \\ \beta \end{bmatrix} = (R^T R)^{-1} R^T (\vec{n_2} - \vec{n_1})$$
(6)

(in pratice computed using *e.g.* the QR decomposition). An estimate of the cornea center is then obtained from the mid-way point

$$\vec{c} = \frac{1}{2}((\alpha \vec{r_1} + \vec{n_1}) + (\beta \vec{r_2} + \vec{n_2})).$$
(7)

While not optimal, this solution is usually acceptable in practice.

The optical axis of the eye passes through the cornea and pupil centers. Typically, the pupil is near the halfway point from the cornea center to the cornea surface, i.e. $r_p \approx 0.40$ [7]. It has, however, some variation from person to person and needs to be determined by a calibration procedure. Given the pupil depth relative to cornea radius r_p , the pupil center is obtained from the intersection of a sphere and the pupil center ray from the camera, as illustrated in Figure 3. Out of the two possible intersections, the one closer to the camera is the pupil center.

It is also possible to take into account the deflection at the cornea surface. A simple ray-tracing method can be employed. The pupil image center defines the ray \vec{p} . Snell's law is applied to the ray at the intersection of \vec{p} and the unit sphere at \vec{c} . The deflected ray continues from the intersection point and intersects the sphere of radius r_p centered at \vec{c} at some point. This is the point through which the optical axis passes.

The point of sharp vision, or fovea, also falls of the optical axis. The offset from the optical axis varies from person to person and is another parameter that needs to be determined by the calibration procedure.

The tracking principle presented here requires a simple calibration that determines parameters related to the users eye. Notably, the camera location relative to the eye does not affect the calibration. The user can freely adjust the position of the glasses, without the need for recalibration.

Camera rotation is implicitly present in the calibration. The system cannot determine the eyes rotation about the optical axis. The offset angles of the fovea are therefore determined in a specific coordinate system defined by the orientation of the camera relative to the eye.

4. Calibration

As explained in Section 3, three calibration values need to be determined: pupil depth radius and the two offset angles that determine the location of the fovea relative to the optical axis. A simple one-shot procedure can be used for user calibration.

The user is asked to look at certain (possibly randomly) selected spots on the screen and the pupil center and reference point location data is collected when looking at these points.

Let $\vec{k_i}$ be the 3-vector in the image plain representing the location of the i:th calibration point on the screen, i.e. so that $K_c \vec{k_i}$ gives the screen coordinates of the calibration point. Given the ray representing the optical axis $\vec{o_i}(r_p)$ for a given pupil depth r_p , when looking at point $\vec{k_i}$, the final coordinates of the gaze in the image plane would be

$$x = \tan(\operatorname{atan}(o_i^1(r_p)) + \alpha_h)$$

$$y = \tan(\operatorname{atan}(o_i^2(r_p)) + \alpha_v).$$
(8)

When the field-of-view of the display is not too wide, the tangent function is fairly linear within the display area. So for a fixed pupil depth, good approximations for the horizontal and vertical correction angles α_1, α_2 are obtained from

$$\alpha_{1} = \frac{1}{N} \sum_{n=1}^{N} \left(\operatorname{atan}(k_{i}^{1}) - \operatorname{atan}(o_{i}^{1}(r_{p})) \right)$$

$$\alpha_{2} = \frac{1}{N} \sum_{n=1}^{N} \left(\operatorname{atan}(k_{i}^{2}) - \operatorname{atan}(o_{i}^{2}(r_{p})) \right).$$
(9)

The squared error in the image plane is given by

$$e = \sum_{i=1}^{N} \sum_{j=1}^{2} (k_i^j - \tan(\tan(o_i^j(r_p)) + \alpha_j))^2.$$
(10)

The mean correction angle can be used as a starting point for iterative optimization of the residual e. The first and second derivatives of Equation 10 for α_i can be computed analytically, so one may simply use Newton's method.

The minimum residual may then be minimized over the pupil depth r_p , in essence performing nested optimizations. The pupil depth is mathematically constrained to be within 0 and 1. Physiologically, the maximum and minimum pupil depths are closer to 0.5. Since this calibration procedure is run only once, time complexity is not critical and one can simply use brute force search. On a modern PC, this optimization takes less than a second.

Alternatively, one could use the Levenberg-Marquardt algorithm on the three parameters to minimize the residual in Equation 10.

5. Pupil Detection

The pupil detection was inspired by the Starburst algorithm [6]. The main idea of Starburst is to search for the pupil contour along rays radiating from some starting point. Keypoints are generated at rising edges along the ray. RANSAC is then used to fit an ellipse to the detected edge points. Here, both of these steps are improved by taking advantage of prior knowledge, some of which is learned during tracking. The RANSAC step is replaced by hill climbing [9] MAPSAC [11].

Rising edges are detected along a limited number of rays extending from a starting point. This provides a set of keypoints along the contour of the pupil. Due to proximity of the camera to the eye and fixed camera optics, focusing problems often blur the edge of the pupil. Instead of thresholding the derivative along the ray as in the original Starburst, thresholding is done on the cumulative change over a fixed segment along the ray. Improved localization of the keypoint is obtained by finding the local maxima in the derivative. Further, a keypoint is required to be near a dark patch. In summary, a point along a ray is marked as a keypoint if the following conditions are satisfied:

- 1. The point is a local maxima in the first derivative along the ray.
- 2. The cumulative change in the value along the ray around the point exceeds a threshold.
- 3. There is a sufficient distance from the last generated keypoint along the ray.
- 4. A dark value was seen not too far back along the ray.
- 5. There are not too many bright pixels near the current position back along the ray.

A total of 96 rays are cast uniformly from starting point all the way to the image border. The starting point for the ray search is either the pupil center from the previous iteration or the image center. The image center is used if the pupil center from the previous iteration is so close to the image border that it is unlikely.

Figure 4 shows the generated keypoints. There are some invalid keypoints at the edges of the corneal reflections, but they are in the vast minority and MAPSAC can easily determine them as outliers.

Instead of RANSAC, a Bayesian MAPSAC [11] is used in conjunction with the hill climbing strategy [9]. The image of the pupil is expected to exhibit certain characteristics. The camera-eye geometry provides initial probability density functions for eccentricity, semi-axis lengths, and location. These prior expectations are then adapted during tracking based on observations.

An ellipse is fit to each set of five randomly selected points. Typical ellipse parameter estimation methods will lead to solutions such that the coefficient vector norm is fixed to a constant [3]. This leads to a situation where some ellipses have geometrically much wider inlier areas than others.



Figure 4. The pupil contour points generated by the ray casting method. The shaded ellipse is the best fit given by MAPSAC and the shading indicates the MAPSAC inlier threshold.

This problem is addressed by the following normalization procedure. Given the ellipse $ax^2 + 2bxy + cy^2 + dx + ey + f = 0$, the ellipse center (x_0, y_0) can be computed. The intersection points of the x- and y-coordinate axes translated to x_0, y_0 with the ellipse are given by

$$\begin{split} \delta_{x} &= \left(4b^{2}y_{0}^{2} + 4by_{0}d + d^{2} - 4acy_{0}^{2} - 4aey_{0} - 4af\right)^{1/2} \\ x_{p} &= x_{0} - \frac{1}{2a}\left(2ax_{0} + 2by_{0} + d + \delta_{x}\right) \\ \delta_{y} &= \left(4b^{2}x_{0}^{2} + 4bx_{0}e + e^{2} - 4cax_{0}^{2} - 4cdx_{0} - 4cf\right)^{1/2} \\ y_{p} &= y_{0} - \frac{1}{2c}\left(2bx_{0} + 2cy_{0} + e + \delta_{y}\right). \end{split}$$
(11)

Finally, the algebraic distance at points $(x_0, y_p - 1)$ and $(x_p - 1, y_0)$ is computed. The coefficients of the ellipse are normalized such that larger of these two algebraic distances is 1.

After this normalization the algebraic distance has a vague geometric meaning near the edge at either the vertical or horizontal intersection. For inlier threshold of 1, the inlier region now extends one pixel vertically or horizontally at this point.

The fraction of points captured as inliers by the model is expected to follow a linear distribution: $P(f_i | \text{model is correct}) = 2f_i$, where

$$f_i = \frac{\text{number of inliers}}{\text{total number of keypoints}}.$$

Let the probability distribution functions for eccentricity, major axis length, and location be $P(\varepsilon \mid \text{mic}), P(m_a \mid \text{mic})$ and $P(\vec{x_0} \mid \text{mic})$, respectively. Where mic, is an abbreviation for "model is correct". Assuming that the distributions are independent, the Bayes rule yields

$$P(\text{model is correct}|\varepsilon, m_a, \vec{x_0}, f_i) \propto P(\varepsilon, m_a, \vec{x_0}, f_i \mid \text{mic}) = (12) 2f_i P(\varepsilon \mid \text{mic}) P(m_a \mid \text{mic}) P(\vec{x} \mid \text{mic}).$$

The characteristics ε , m_a , $\vec{x_0}$ and f_i can be computed from the coefficients and the posterior probability is given by Equation 12. Each model is scored on the posterior probability and the hill climbing strategy is used to search for the highest scoring model.

The score of the proposed model is linearly proportional to the number of inliers, so the hill climbing strategy for ellipse fitting presented in [9] can be used without modification. The weighting is based solely on the inlier-outlier categorization, but is only updated when a new model is accepted based on MAP-score.

The overhead of the hill climbing strategy is only one random number per iteration and one pass over the keypoints to regenerate the inlier set whenever the proposed model is accepted. The latter happens relatively few times for each model fitting task.

6. Specular Reflection Detection

The specular reflections from the reference lights appear as small bright spots in the camera image. Since similar reflections creep into the system from uncontrolled light sources, there are many possible candidates for reference spots. Figure 4 shows the camera image under low ambient light conditions, where the reference spots are clearly visible.

A simple spot trigger is used in an area near the pupil center. The spot kernel is defined as

(0	0	25	0	0 \	
0	0	0	0	0	
25	0	1	0	25	
0	0	0	0	0	
0	0	25	0	0 /	

Local maxima of the spot response, above a dynamically determined threshold, are treated as possible reference reflection points. The threshold is set so that a certain number of pixels remain above the threshold.

Again, certain properties are expected from the correct reflection pair. Reflections should be fairly high contrast (high response to the kernel), nearly horizontally aligned, near the pupil center and within certain distance range from each other. Bayesian belief based probabilities are given for each of the possible attributes.

The highest scoring pair is then selected as the true reflection points. This kind of Bayesian scoring proves to be highly robust. The Bayesian scoring encapsulates the prior knowledge of the appearance and location of the two point constellation. It can reliably find the two relevant spots from a set of tens of possible individual reflection points. The combined computational cost of the spot detector and scoring is fairly low.

After the proper pair is found, the location is refined to sub-pixel accuracy by means of parabolic interpolation of

	A1	A2	B1	B2	C1	C2
Subject 1	0.73	0.86	1.00	1.24	0.73	0.92
Subject 2	0.74	0.65	0.78	1.24	0.71	0.89
Subject 3	0.59	1.06	1.63	0.94	1.00	0.80
Average	0.69	0.86	1.14	1.14	0.81	0.87

Table 1. Average angular error when looking at the 9 test points. Glasses were removed between tests A,B and C. The test grid was viewed two times in each test.

the pixel intensity values. The function

$$f(x,y) = ax^{2} + by^{2} + cx + dy + c$$
(13)

is fitted to the data in the neighboring pixels to the spot translated to the origin. The function has maximum at

$$\left[\delta_x, \delta_y\right] = \left[\frac{-c}{2a}, \frac{-d}{2b}\right],\tag{14}$$

which gives the offset for the spot location.

7. System Validation

The validation method used in [6] was adapted for our prototype near-to-eye tracker. Since the system is expected to be tolerant to changes in device position on the head, the test subjects were asked to remove the glasses between measurements.

Three subjects with normal vision, all of whom had been using the device before, participated in the test. The calibration algorithm was run for each, using 9 calibration points and the results were saved. Each subject was asked to follow 9 test targets on a 3 by 3 test grid. The test targets were displayed one at a time in a random order.

Each subject viewed the test targets twice in a row. Then the glasses were handed to the next subject. The saved calibration values for the subject in question were loaded and the test process repeated for him. The whole process was repeated three times, so that each subject removed the glasses twice between measurements and viewed the test grid six times in total. Calibration values were estimated only once for each subject.

The results are displayed in Table 1. The average angular error over all tests and subjects was 0.92 degrees. While the best accuracy was obtained right after calibration, no clear deterioration of accuracy was observed due to device position drift. It should be noted, however, that if the glasses are forced to an extreme position beyond where they naturally set on the user's face, angular offset errors up to 3-5 degrees can be observed. The error is significantly less for normal positions of the glasses.

8. Conclusion and Future Work

A special case of corneal reflection based tracking, where the light sources are at infinity, has been described.

It is particularly suited for near-to-eye displays, where the camera must view the eye behind the display and thus the camera-to-eye distance must be quite small. This situation is difficult for existing gaze trackers.

A set of evolutionary steps in pupil and corneal reflection detection were described that significantly increase the robustness of the system. Finally, a user friendly calibration procedure for the system was described. A working prototype based on these principles has been tested and found to perform very well.

During development, the system has been tested by a large number of subjects. For most it works well but for some it fails to track the pupil. The biggest issue is the illumination and the FOV of the camera. The geometry of the human face has so much variance that for some people the pupil tends to fall outside the camera view or is not properly illuminated. These are more likely to have a mechanical, rather than algorithmic, solution.

The experiments in this paper were done on a prototype that has a geometry similar to what is depicted in Figure 2. The display was designed for video viewing and the back side is covered in dark plastic to block out the outside world. A camera looking straight at the eye can be used in a nontransparent system. In this configuration the eye is always observed at approximately the same distance and the focus of the camera is not an issue.

A system where the camera is looking at the eye from an angle as in Figure 3 is in an early prototype phase but already working. In this system, the back side of the display will be covered with a controllable polarizer that can be used to control the amount of light from the outside world, while keeping the system transparent. The requirements on illumination and camera sensitivity are higher in this configuration, as the pupil moves in depth relative to the camera and a larger depth of field is required.

Hopefully, this new prototype will provide the basis for a lightweight, mobile and discreet augmented reality system that can be adaptive to the users gaze.

References

- [1] J. S. Babcock and J. B. Pelz. Building a lightweight eyetracking headgear. In ACM Eye Tracking Research and Applications Symposium, pages 109–114, San Antonio, TX, USA, March 2004. 1, 2
- [2] E. D. Guestrin and M. Eizenman. General theory of remote gaze estimation using the pupil center and corneal reflections. *IEEE Transactions on Biomedical Engineering*, 53(6):1124–1133, June 2006. 2
- [3] R. Hartley and A. Zisserman. Multiple View Geometry in Computer Vision. Cambridge University Press, second edition, 2003. 4
- [4] T. Levola. Diffractive optics for virtual reality displays. *Journal of the Society for Information Display*, 14(5):467–475, 2006.

- [5] T. Levola. Novel diffractive optical components for near to eye displays. In SID International Symposium. Digest of Technical Papers, Vol. XXXVII, Book I, pages 64–67, 2006.
- [6] D. Li, W. David, and D. J. Parkhurst. Starburst: A hybrid algorithm for video-based eye tracking combining featurebased and model-based approaches. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2005. 1, 2, 4, 6
- [7] C. H. Morimoto and M. R. Mimica. Eye gaze tracking techniques for interactive applications. *Computer Vision and Image Understanding*, 98:4–24, April 2005. 3
- [8] J. B. Pelz, R. L. Canosa, D. Kucharczyk, J. Babcock, A. Silver, and D. Konno. Portable eyetracking: A study of natural eye movement. In *Proceedings of the SPIE, Human Vision and Elecronic Imaging*, pages 566–582, San Jose, CA, USA, 2000. 1, 2
- [9] T. Pylvänäinen and L. Fan. Hill climbing method for random sample consensus methods. In *International Symposium on Visual Computing*, 2007. 4, 5
- [10] D. Roger and N. Holzschuch. Accurate specular reflections in real-time. *Computer Graphics Forum (Proceedings of Eurographics 2006)*, 25(3), September 2006. 3
- [11] P. H. Torr and C. Davidson. IMPSAC: Synthesis of importance sampling and random sample consensus. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(3):354–364, March 2003. 4