

# FACIAL POSE ESTIMATION USING ACTIVE APPEARANCE MODELS AND A GENERIC FACE MODEL

Thorsten Gernoth, Katerina Alonso Martínez, André Gooßen and Rolf-Rainer Grigat  
Vision Systems E-2, Hamburg University of Technology, Harburger Schloßstr. 20, 21079 Hamburg, Germany

**Keywords:** Pose estimation, Active appearance model, Infrared imaging, Face recognition.

**Abstract:** The complexity in face recognition emerges from the variability of the appearance of human faces. While the identity is preserved, the appearance of a face may change due to factors such as illumination, facial pose or facial expression. Reliable biometric identification relies on an appropriate response to these factors. In this paper we address the estimation of the facial pose as a first step to deal with pose changes. We present a method for pose estimation from two-dimensional images captured under active infrared illumination using a statistical model of facial appearance. An active appearance model is fitted to the target image to find facial features. We formulate the fitting algorithm using a smooth warp function, namely thin plate splines. The presented algorithm requires only a coarse and generic three-dimensional model of the face to estimate the pose from the detected features locations. The desired field of application requires the algorithm to work with many different faces, including faces of subjects not seen during the training stage. A special focus is therefore on the evaluation of the generalization performance of the algorithm which is one weakness of the classic active appearance model algorithm.

## 1 INTRODUCTION

In the modern society there is a high demand to automatically and reliably determine or verify the identity of a person. For example, to control entry to restricted access areas. Using biometric data to identify a target person has some well known conceptual advantages, such as the identification procedure is immutable bound to the person which should be identified. Using facial images as a biometric characteristic has gained much attention and commercially available face recognition systems exist (Zhao et al., 2003, Phillips et al., 2007). However unconstrained environments with variable ambient illumination and changes of head pose are still challenging for many face recognition systems.

The appearance of a face can vary drastically if the intensity or the direction of the light source changes. This problem can be overcome by employing active imaging techniques to capture face images under invariant illumination conditions. In this work we use active near-infrared (NIR) illumination (Gernoth and Grigat, 2010). Possible surrounding light in the visible spectrum is filtered out.

Another benefit of active near-infrared illumination is the *bright pupil effect* which can be employed to assist eye detection. Pupils appear as unnaturally

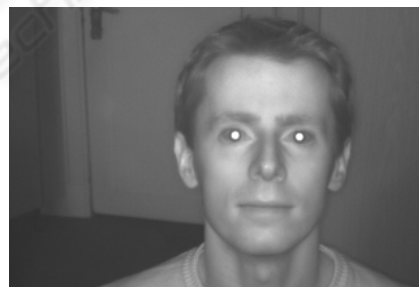


Figure 1: The bright pupil effect perceivable under active near-infrared illumination.

bright spots when an active near-infrared radiation source is mounted close to the camera axis (Figure 1). We use image processing to detect these bright spots in the images and thus can reliably detect the eyes (Zhao and Grigat, 2006).

Challenging for face recognition systems are also changes of head pose. Appearance-based face recognition systems use the texture of faces in the form of two-dimensional frontal images to identify a target person. But faces are three-dimensional objects and due to head pose changes, their appearance in images can change significantly.

There are three different main strategies to over-

come this problem in appearance-based approaches for face recognition. The first is to use features which are invariant to these deformations, e. g. invariant to changes of the facial pose relative to the camera. Another strategy is to use or synthetically generate a large and representative training set. A third approach is to separate the factors which code the identity of a person from other sources of variation, such as pose changes. This is addressed in this paper. The posture of the head in-front of the camera is estimated from monocular images. The additional pose information may be utilized to register the facial images very precise and thereby make it possible to perform face recognition using a pose normalized representation of faces.

A survey of head pose estimation in computer vision was recently published by Murphy-Chutorian and Trivedi, 2009. We use active appearance models (AAM) to detect facial features in the images. Subsequently, the head pose is determined from a subset of the localized facial features using an analytical algorithm (DeMenthon and Davis, 1995, Martins and Batista, 2008). The algorithm can estimate the pose from a single image using four or more non-coplanar facial features positions and their known relative geometry. Using three-dimensional model points from the generic Candide-3 face model (Ahlberg, 2001, Dornaika and Ahlberg, 2006) and their image correspondences estimated using the active appearance model, the posture of the head in-front of the camera can be estimated.

Active appearance models are a common approach to build parametric statistical models of facial appearance (Cootes et al., 2001, Stegmann et al., 2000). The desired field of application requires the algorithm to work with many different faces, including faces not seen during the training stage (Gross et al., 2005). We use simultaneous optimization of pose and texture parameters and formulate the fitting algorithm using a smooth warping function (Bookstein, 1989). The thin plate spline warping function is parametrized efficiently to achieve some computational advantages. A special focus is on the evaluation of the generalization performance of the model fitting algorithm.

In Section 2 we introduce statistical models of facial appearance. Section 3 describes the smooth warping function. The pose estimation algorithm is explained in Section 4. With experimental results and discussion in Section 5, we conclude in Section 6.

## 2 STATISTICAL MODELS OF FACIAL APPEARANCE

We parametrize a dense representation of facial appearance using separate linear models for shape and texture (Matthews and Baker, 2004). The shape and texture parameters of the models are statistically learned from a training set.

### 2.1 Facial Model

Shape information is represented by an ordered set of  $l$  landmarks  $\mathbf{x}_i, i = 1 \dots l$ . These landmarks describe the planar facial shape of an individual in a digital image. The landmarks are generally placed on the boundary of prominent face components (Figure 2a). The two-dimensional landmark coordinates are arranged in a shape matrix (Matthews et al., 2007)

$$\mathbf{s} = (\mathbf{x}_1 \quad \mathbf{x}_2 \quad \dots \quad \mathbf{x}_l)^\top, \quad \mathbf{s} \in \mathbb{R}^{l \times 2}. \quad (1)$$

Active appearance models express an instance  $\mathbf{s}_p$  of a particular shape as mean shape  $\mathbf{s}_0$  and a linear combination of  $n$  eigenshapes  $\mathbf{s}_i$ , i.e.

$$\mathbf{s}_p = \mathbf{s}_0 + \sum_{i=1}^n p_i \mathbf{s}_i. \quad (2)$$

The coefficients  $p_i$  constitute the shape parameter vector  $\mathbf{p} = (p_1, \dots, p_n)^\top$ . The mean shape  $\mathbf{s}_0$  and shape variations  $\mathbf{s}_i$  are statistically learned using a training set of annotated images (Figure 2a). Since reliable pupil positions are available (Zhao and Grigat, 2006), the training images can be aligned with respect to the pupils in a common coordinate system  $\mathcal{I} \subset \mathbb{R}^2$  using a rigid transformation. The images are rotated, scaled and translated using a two-dimensional similarity transform such that all the pupils fall in the same position (Figure 2b). The mean shape  $\mathbf{s}_0$  and basis of shape variations  $\mathbf{s}_i$  are obtained by applying principal component analysis (PCA) on the shapes of the aligned training images (Cootes et al., 2001).

The texture part of the appearance is also modeled using an affine linear model of variation. Texture is defined as the intensities of a face at a discrete set  $\mathcal{A}_0$  of positions  $\mathbf{x}$  in a shape-normalized space  $\mathcal{A} \subset \mathbb{R}^2$ . The texture of a face is vectorized by raster-scanning it into a vector. Similar to the shape,  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_m)^\top$  denotes a vector of texture parameters describing a texture instance

$$\mathbf{a}_\lambda = \mathbf{a}_0 + \sum_{i=1}^m \lambda_i \mathbf{a}_i. \quad (3)$$

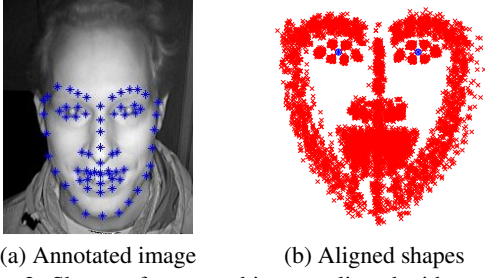


Figure 2: Shapes of annotated images aligned with respect to pupil positions.

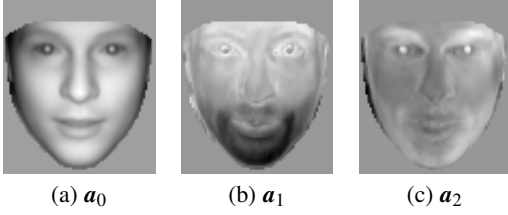


Figure 3: The mean texture  $a_0$  and the first two basis of texture variations  $a_i$ .

The texture at position  $x \in \mathcal{A}_0$  of  $a_\lambda$  is a function of the domain  $\mathcal{A}$ , with

$$a_\lambda : \mathcal{A} \rightarrow \mathbb{R}; \quad x \mapsto a_\lambda(x). \quad (4)$$

To create a texture model, all the aligned training images are warped into the shape-normalized space. The shape-normalized space is given by the mean shape  $s_0$  of the shape model. A smooth warping function that maps one image to another by relating two sets of landmarks is used as described in Section 3.  $\mathcal{A}_0$  contains positions that lie inside the mean shape  $s_0$ . PCA is applied on the training textures to obtain the mean texture  $a_0$  and basis of texture variations  $a_i$ .

Photometric variations of the texture  $a_\lambda$  are modeled separately by a global texture transformation  $T_u(a_\lambda(x)) = (u_i + 1)a_\lambda(x) + u_2$  (Baker et al., 2003). The intensities of the texture vector  $a_\lambda$  are scaled by a global gain factor  $(u_i + 1)$  and biased by  $u_2$ .

To simplify the notation, the parameters describing shape, texture and photometric variations are combined into the single parameter vector

$$\mathbf{q} = (\mathbf{p}^\top \quad \mathbf{u}^\top \quad \lambda^\top)^\top. \quad (5)$$

## 2.2 Model Fitting

The parameters of the generative model described in Section 2.1 need to be estimated to fit the model to a target image. The target image can be aligned to the common coordinate system with respect to the pupils

(Section 2.1). The target image is regarded as a continuous function of the domain  $\mathcal{I}$ :

$$I : \mathcal{I} \rightarrow \mathbb{R}; \quad x' \mapsto I(x'). \quad (6)$$

Fitting the model to an image is generally done by minimizing some error measure between the modeled texture and the target image. The error at the position  $x \in \mathcal{A}_0$  between the generated texture and the target image is

$$e(\mathbf{x}, \mathbf{q}) = a_\lambda(\mathbf{x}) - T_u(I(W(\mathbf{x}, \mathbf{p}))). \quad (7)$$

$W(\mathbf{x}, \mathbf{p})$  is a non-linear warping function that maps positions  $\mathbf{x} \in \mathcal{A}$  of the model to positions  $\mathbf{x}' \in \mathcal{I}$  of the target image. The warping function is parametrized by the shape parameters  $\mathbf{p}$  as described in Section 3.

Typically the sum-of-squared error of all positions  $\mathbf{x}$  is minimized to find the parameters  $\mathbf{q}$ , such that

$$\operatorname{argmin}_{\mathbf{q}} \frac{1}{2} \sum_{\mathbf{x} \in \mathcal{A}_0} [e(\mathbf{x}, \mathbf{q})]^2. \quad (8)$$

This is a non-linear optimization problem. General optimization techniques can be used to find a solution. Commonly used is an iterative Gauß-Newton style algorithm (Matthews and Baker, 2004). We assume a current estimate of  $\mathbf{q}$  and solve for incremental updates  $\Delta \mathbf{q}$  in each step. The update can be combined with the previous estimate in several ways (Matthews and Baker, 2004). The simplest update is the linear additive increment  $\mathbf{q} \leftarrow \mathbf{q} + \Delta \mathbf{q}$ . The following expression is minimized with respect to  $\Delta \mathbf{q}$ :

$$\operatorname{argmin}_{\Delta \mathbf{q}} \frac{1}{2} \sum_{\mathbf{x} \in \mathcal{A}_0} [e(\mathbf{x}, \mathbf{q} + \Delta \mathbf{q})]^2. \quad (9)$$

Performing a first-order Taylor expansion of the residual  $e(\mathbf{x}, \mathbf{q} + \Delta \mathbf{q})$  around  $\mathbf{q}$  yields:

$$e(\mathbf{x}, \mathbf{q} + \Delta \mathbf{q}) \approx e(\mathbf{x}, \mathbf{q}) + \left( \frac{\partial e(\mathbf{x}, \mathbf{q})}{\partial \mathbf{q}} \right)^\top \Delta \mathbf{q} \quad (10)$$

$$\approx e(\mathbf{x}, \mathbf{q}) + \left( \frac{\partial e(\mathbf{x}, \mathbf{q})}{\partial \mathbf{p}} \right)^\top \Delta \mathbf{p} + \quad (11)$$

$$\left( \frac{\partial e(\mathbf{x}, \mathbf{q})}{\partial \mathbf{u}} \right)^\top \Delta \mathbf{u} + \left( \frac{\partial e(\mathbf{x}, \mathbf{q})}{\partial \lambda} \right)^\top \Delta \lambda.$$

The Gauß-Newton algorithm uses the update

$$\Delta \mathbf{q} = -\mathbf{H}^{-1} \sum_{\mathbf{x} \in \mathcal{A}_0} e(\mathbf{x}, \mathbf{q}) \frac{\partial e(\mathbf{x}, \mathbf{q})}{\partial \mathbf{q}} \quad (12)$$

with

$$\mathbf{H} = \sum_{\mathbf{x} \in \mathcal{A}_0} \frac{\partial e(\mathbf{x}, \mathbf{q})}{\partial \mathbf{q}} \left( \frac{\partial e(\mathbf{x}, \mathbf{q})}{\partial \mathbf{q}} \right)^\top. \quad (13)$$

Solving for  $\Delta \mathbf{q}$  and using a Gauß-Newton step to optimize Eq. (9) involves computing an approximation of the Hessian matrix (Eq. (13)) and its inverse in each iteration. Assuming a constant Hessian matrix results in significant computational savings (Amberg et al., 2009, Cootes et al., 2001).

Matthews and Baker, 2004 showed that the linear additive increment is not the only parameter update strategy. They introduced compositional update strategies which permit overall cheaper algorithms for person-specific active appearance models. Their *simultaneous inverse compositional algorithm* for person-independent active appearance models (Gross et al., 2005) is however not as computationally efficient.

### 3 THIN PLATE SPLINE WARP

A warp function maps positions of one image to positions of another image by relating two sets of landmarks. The most common warp functions are the piecewise affine warp (Glasbey and Mardia, 1998) and the thin plate spline (TPS) warp (Bookstein, 1989). The affine warp function has the advantage of being simple and linear in a local region. But although it gives a continuous deformation, it is not smooth. Thin plate spline warping as an alternative produces a smoothly warped image. However, it is more expensive to calculate and non-linear due to the interpolating function used. In this paper we focus on the thin plate spline warp.

In the case of our active appearance model, the warp maps the positions  $\mathbf{x}$  from the shape-normalized space  $\mathcal{A}$  to positions  $\mathbf{x}' \in \mathcal{I}$  of the target image. The transformation is such that the landmarks  $\mathbf{x}_i, i = 1 \dots l$  are mapped to corresponding landmarks  $\mathbf{x}'_i, i = 1 \dots l$  of a shape instance  $s_p$  in the target image. Since the landmark positions in the target image depend on the shape parameters  $\mathbf{p}$ , we parametrize the warp function by the shape parameter vector  $\mathbf{p}$ .

The thin plate spline warp function  $\mathbf{W} : \mathcal{A} \rightarrow \mathcal{I}$  is vector valued and defined as (Bookstein, 1989, Cootes and Taylor, 2004):

$$\mathbf{W}(\mathbf{x}, \mathbf{p}) = \left( \sum_{i=1}^l w_i U(\|\mathbf{x} - \mathbf{x}_i\|) \right) + \mathbf{c} + \mathbf{C}\mathbf{x} \quad (14)$$

$$= \underbrace{\mathbf{W}(\mathbf{p})}_{2 \times (l+3)} \cdot \underbrace{\mathbf{k}(\mathbf{x})}_{(l+3) \times 1}, \quad (15)$$

with

$$\mathbf{W}(\mathbf{p}) = [\mathbf{w}_1 \quad \dots \quad \mathbf{w}_l \quad \mathbf{c} \quad \mathbf{C}], \quad (16)$$

$$\mathbf{k}(\mathbf{x}) = [U(r_1(\mathbf{x})) \quad \dots \quad U(r_l(\mathbf{x})) \quad 1 \quad \mathbf{x}^\top]^\top, \quad (17)$$

where  $r_i(\mathbf{x}) = \|\mathbf{x} - \mathbf{x}_i\|$  is the Euclidean distance between a position  $\mathbf{x}$  and landmark  $\mathbf{x}_i$  of the mean shape  $s_0$ .  $U(r)$  is the TPS interpolating function (e.g.  $U(r) = r^2 \log r^2$  with  $U(0) = 0$ ) that makes the warp function non-linear.

$\mathbf{W}(\mathbf{p})$  contains the warp weights. The weights  $\mathbf{c}, \mathbf{C}$  represent the affine part of the mapping in Eq. (14). The warp weights are defined by the sets of source and destination landmarks to satisfy the constraints  $\mathbf{W}(\mathbf{x}_i, \mathbf{p}) = \mathbf{x}'_i \forall i \in \{1 \dots l\}$  and to minimize the bending energy. Combining all constraints yields in a linear system (Bookstein, 1989):

$$\begin{bmatrix} \mathbf{K} & \mathbf{P} \\ \mathbf{P}^\top & \mathbf{O} \end{bmatrix} \mathbf{W}(\mathbf{p})^\top = \mathbf{L} \mathbf{W}(\mathbf{p})^\top = \begin{bmatrix} \mathbf{s}_p \\ \mathbf{o} \end{bmatrix} \quad (18)$$

where  $\mathbf{K}$  is a  $l \times l$  matrix and  $K_{ij} = U(\|\mathbf{x}_i - \mathbf{x}_j\|)$ , the  $i$ 'th row of the  $l \times 3$  matrix  $\mathbf{P}$  is  $(1 \quad \mathbf{x}_i^\top)$ ,  $\mathbf{O}$  is a  $3 \times 3$  matrix of zeros and  $\mathbf{o}$  is a  $3 \times 2$  matrix of zeros. If  $\mathbf{L}$  is non-singular the warp weights are given by

$$\mathbf{W}(\mathbf{p}) = (\mathbf{L}^{-1} \mathbf{B} \mathbf{s}_p)^\top \quad (19)$$

with

$$\mathbf{B} = \underbrace{\begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \dots & 1 \\ 0 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 \end{bmatrix}}_{(l+3) \times l}. \quad (20)$$

During model fitting (Section 2.1) we estimate iteratively the parameters  $\mathbf{p}$  of the shape  $s_p$  that corresponds to the shape of the target image. Since the residual in Eq. (7) is defined within the shape-normalized space  $\mathcal{A}$  given by  $s_0$ , the matrix  $\mathbf{L}$  depends only on the landmarks  $\mathbf{x}_i$  of the mean shape  $s_0$ . The mean shape does not change for a give training set.  $\mathbf{L}$  and its inverse can therefore be precomputed.

Parametrization of the warp function with respect to  $\mathbf{p}$  yields another computational advantage. The texture is defined at a discrete set of positions  $\mathcal{A}_0$  (Section 2.1). The residual in Eq. (7) only need to be evaluated at these positions. The positions depend on  $s_0$  and do not change for a give training set.  $\mathbf{k}(\mathbf{x})$  can also be precomputed for all  $\mathbf{x} \in \mathcal{A}_0$ .

Using Gauß-Newton to minimize Eq. (9) requires the Jacobian matrix  $\frac{\partial \mathbf{W}(\mathbf{x}, \mathbf{p})}{\partial \mathbf{p}}$  of the warp function with respect to the shape parameters (Eq. (11)):

$$\left(\frac{\partial e(\mathbf{x}, \mathbf{q})}{\partial \mathbf{p}}\right)^\top = -\left(\frac{\partial T_u(I(\mathbf{W}(\mathbf{x}, \mathbf{p})))}{\partial \mathbf{p}}\right)^\top \quad (21)$$

$$= -\left(\frac{\partial T_u(I(\mathbf{W}(\mathbf{x}, \mathbf{p})))}{\partial \mathbf{W}(\mathbf{x}, \mathbf{p})}\right)^\top \frac{\mathbf{W}(\mathbf{x}, \mathbf{p})}{\partial \mathbf{p}} \quad (22)$$

Using the parametrization of the warp function with respect to  $\mathbf{p}$  gives the components of the Jacobian matrix as follows:

$$\frac{\partial \mathbf{W}(\mathbf{x}, \mathbf{p})}{\partial p_i} = s_i (\mathbf{L}^{-1} \mathbf{B})^\top \mathbf{k}(\mathbf{x}). \quad (23)$$

The Jacobian does not depend on the value of the evaluation point  $\mathbf{p}$  and can be precomputed for all  $\mathbf{x} \in \mathcal{A}_0$ .

## 4 POSE ESTIMATION

The pose of a face in-front of the camera is defined as its position and orientation relative to a three-dimensional camera coordinate system. We denote by  $\hat{\mathbf{x}}_i$  a three-dimensional feature point of a face model in a coordinate system attached to the model. The perspective projection of a feature point onto the image plane of the camera is  $\mathbf{x}_i$ . The pose of a face with respect to a camera can be defined with a translation vector  $\mathbf{t} \in \mathbb{R}^3$  and a rotation matrix  $\mathbf{R} \in \text{SO}(3)$ :

- The translation vector is the vector from the origin of the coordinate system attached to the camera to the origin of the face model:  $\mathbf{t} = (t_x \ t_y \ t_z)^\top$ .
- The rotation matrix is the matrix whose rows are the unit vectors of the camera coordinate system expressed in the coordinate system of the face model:  $\mathbf{R} = [\mathbf{i} \ \mathbf{j} \ \mathbf{k}]^\top$ .

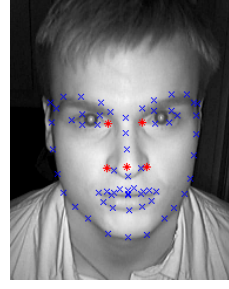
Under perspective projection with a camera having focal length  $f$ ,  $\mathbf{x}_i$  is related to the corresponding feature point of the model  $\hat{\mathbf{x}}_i$  as follows:

$$\mathbf{x}_i = \frac{f}{\mathbf{k}^\top \hat{\mathbf{x}}_i + t_z} \begin{pmatrix} \mathbf{i}^\top \\ \mathbf{j}^\top \end{pmatrix} \hat{\mathbf{x}}_i + \begin{pmatrix} t_x \\ t_y \end{pmatrix}. \quad (24)$$

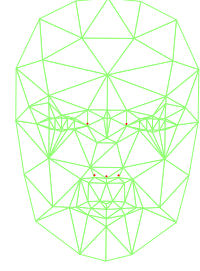
The equation may be written as in DeMenthon and Davis, 1995,

$$\mathbf{x}_i(1 + \epsilon_i) - \begin{bmatrix} t'_x \\ t'_y \end{bmatrix} = \begin{bmatrix} \mathbf{i}'^\top \\ \mathbf{j}'^\top \end{bmatrix} \hat{\mathbf{x}}_i, \quad (25)$$

with  $\mathbf{i}' = \frac{f}{t_z} \mathbf{i}$ ,  $\mathbf{j}' = \frac{f}{t_z} \mathbf{j}$ ,  $t'_x = \frac{f}{t_z} t_x$ ,  $t'_y = \frac{f}{t_z} t_y$  and  $\epsilon_i = \frac{\mathbf{k}^\top \hat{\mathbf{x}}_i}{t_z}$ . By setting  $\epsilon_i = 0$ ,  $\mathbf{x}_i$  equals the scaled orthographic projection of the face model point  $\hat{\mathbf{x}}_i$ . Scaled orthographic projection is similar to perspective projection if the depth of the face is small compared to its distance to the camera. For the case of fixed  $\epsilon_i$  and assuming known projection of the model origin onto the



(a) AAM Landmarks



(b) Candide Model

Figure 4: (a) The selected landmarks for pose estimation and (b) the corresponding Candide-3 model.

image plane, the unknown  $\mathbf{i}$ ,  $\mathbf{j}$  and  $t_z$  can be computed from Eq. (25) from at least 4 non-coplanar feature points (DeMenthon and Davis, 1995, Martins and Batista, 2008). The third row of the rotation matrix  $\mathbf{R}$  can be obtained from the cross product  $\mathbf{k} = \mathbf{i} \times \mathbf{j}$ . Using an iterative scheme and estimating  $\epsilon_i$  from  $\mathbf{k}$  and  $t_z$  of the previous iteration, an approximation of the pose can be computed.

### 4.1 Generic Face Model

A three-dimensional model of the face shape is required to estimate the pose. The Candide-3 model (Ahlberg, 2001) is used. This general model is used in its neutral state but adapted in scale for facial pose estimation of all subjects (Figure 4b).

### 4.2 Feature Point Selection

As stated above, the algorithm requires at least 4 non-coplanar feature points to estimate the facial pose. These points are chosen from the estimated landmark positioning provided by the AAM. A correspondence was established between the Candide-3 model points and the landmarks of the active appearance model.

Since only one generic face model is used, the landmarks which do not vary much between faces of different subjects are chosen as feature points for pose estimation. The variability of each landmark was studied in order to pick the most stable ones. As it can be seen in Figure 4a, two landmarks from the eye contour and landmarks on the nose are chosen. The nose tip is used as the origin of the coordinate system of the face model.

## 5 EXPERIMENTS

We use the TUNIR database (Zhao et al., 2007) for all experiments. The database consists of recordings

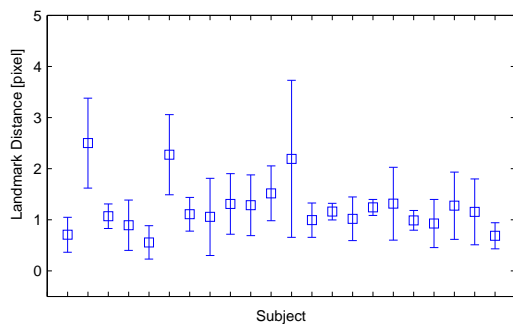


Figure 5: Mean and standard deviation of the distance between estimated and hand-labeled landmarks used for pose estimation for each subject (Person-specific AAM).

of 74 people in a typical access control scenario under active near-infrared illumination. The subjects move in-front of the camera and were asked to speak to recreate a realistic scenario.

### 5.1 Model Fitting

To evaluate the performance of the AAM fitting algorithm with smooth warp function, 5 images of 22 subjects from the TUNIR database were labeled in a semi-automatic way with 67 landmarks. Before fitting, the shapes were prepositioned according to the pupil positions, as described in Section 2.1. A hierarchical approach with two levels was chosen. To evaluate the fitting ability of the algorithm and for comparison, an experiment with known subject but novel image was conducted. Person-specific active appearance models were trained for each subject and evaluated in a leave-one-out cross-validation manner. 90% of shape variance and 95% of appearance variance of each training set were retained in the model. This corresponds to the optimal settings for the experiment with generic active appearance models described below. In Figure 5 the mean and standard deviation of the distance between the estimated and the hand-labeled landmarks are shown for each subject. Only the landmarks which are used for pose estimation contribute to the evaluation of the landmark distance.

Of more interest for the desired field of application is the performance with subjects not seen during training. In a second experiment generic active appearance models were trained from all but one subject. All images of the remaining identity were used to evaluate the performance of the fitting algorithm. Of interest is again the distance between the estimated and the hand-labeled landmarks used for pose estimation. This is shown in Figure 6.

As expected, the mean distance between estimated

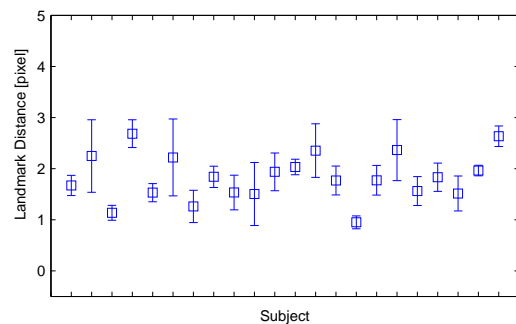


Figure 6: Mean and standard deviation of the distance between estimated and hand-labeled landmarks used for pose estimation for each subject (Generic AAM).

and hand-labeled landmarks is lower for person specific active appearance models. Nevertheless, the total mean distance is also just 1.8 pixels for generic active appearance models. In many cases the estimated landmarks fit well the face components but do not match the hand-labeled landmarks. This is because the hand-labeled positioning is not necessary the optimum one. Data refitting (Gross et al., 2005) could improve the performance. After a visual inspection of the estimated landmarks, an error up to 3 pixels was considered as good performance for the application.

### 5.2 Pose Estimation

To test the accuracy of the pose estimation algorithm quantitatively, we used the three-dimensional Candide-3 model. Of interest was the quality of the pose estimation from landmarks perturbed in the range of what can be expected from the AAM fitting algorithm. The Candide-3 model was situated in 729 different positions. We obtained simulated landmarks by projecting the three-dimensional model points respectively to an image plane corresponding to the application scenario. We uniformly perturbed these landmarks from 1 pixels to 10 pixels and estimated the pose using the algorithm described in Section 4.

In Figure 7, the mean distance of the true model points to the model points of the Candide-3 model with estimated pose for different ranges of landmark perturbation is shown. Figure 7 shows that for a perturbation between 3 and 4 pixels, the mean model point distance is just around 7 pixels.

A typical result of pose estimation is shown in Figure 8. Since the posture of the subjects in-front of the camera is not known for the test images of the TUNIR database, the performance could only be evaluated qualitatively for this database.

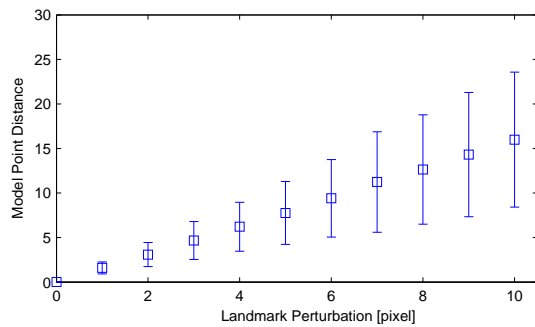


Figure 7: Mean and standard deviation of the distance between true model points and model points with estimated pose for different ranges of landmark perturbation.

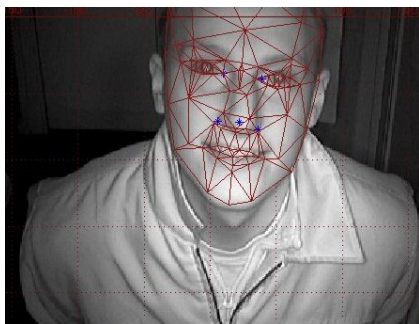


Figure 8: The estimated pose of a subject of the TUNIR database.

## 6 CONCLUSIONS

We presented an approach for facial pose estimation from two-dimensional images using active appearance models. Only a generic three-dimensional face model is required for pose estimation. We formulated the active appearance model fitting algorithm in an efficient manner with a smooth warp function. Our experiments show that the fitting accuracy of the algorithm is sufficient to estimate the pose from the detected landmark positions. Estimated poses from test images of the TUNIR database emphasize this result qualitatively.

## ACKNOWLEDGEMENTS

This work is part of the project *KabTec – Modulares integriertes Sicherheitssystem* funded by the German Federal Ministry of Economics and Technology.

## REFERENCES

- Ahlberg, J. (2001). Candide-3 – an updated parameterized face. Technical Report LiTH-ISY-R-2326, Dept. of Electrical Engineering, Linköping University, Linköping, Sweden.
- Amberg, B., Blake, A., and Vetter, T. (2009). On compositional image alignment, with an application to active appearance models. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2009)*, pages 1714–1721.
- Baker, S., Gross, R., and Matthews, I. (2003). Lucas-Kanade 20 years on: A unifying framework: Part 3. Technical Report CMU-RI-TR-03-35, Robotics Institute, Carnegie Mellon University, Pittsburgh, PA.
- Bookstein, F. (1989). Principal warps: thin-plate splines and the decomposition of deformations. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 11(6):567–585.
- Cootes, T., Edwards, G., and Taylor, C. (2001). Active appearance models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 23(6):681–685.
- Cootes, T. and Taylor, C. (2004). Statistical models of appearance for computer vision. Technical report, University of Manchester, Manchester, U.K.
- DeMenthon, D. F. and Davis, L. S. (1995). Model-based object pose in 25 lines of code. *International Journal of Computer Vision*, 15(1-2):123–141.
- Dornaika, F. and Ahlberg, J. (2006). Fitting 3d face models for tracking and active appearance model training. *Image and Vision Computing*, 24(9):1010 – 1024.
- Gernoth, T. and Grigat, R.-R. (2010). Camera characterization for face recognition under active near-infrared illumination. In *Proc. SPIE*, volume 7529, page 75290Z.
- Glasbey, C. A. and Mardia, K. V. (1998). A review of image-warping methods. *Journal of Applied Statistics*, 25(2):155–171.
- Gross, R., Matthews, I., and Baker, S. (2005). Generic vs. person specific active appearance models. *Image and Vision Computing*, 23(1):1080–1093.
- Martins, P. and Batista, J. (2008). Monocular head pose estimation. In *Proc. 5th international conference on Image Analysis and Recognition (ICIAR '08)*, pages 357–368.
- Matthews, I. and Baker, S. (2004). Active appearance models revisited. *International Journal of Computer Vision*, 60(2):135–164.
- Matthews, I., Xiao, J., and Baker, S. (2007). 2d vs. 3d deformable face models: Representational power, construction, and real-time fitting. *International Journal of Computer Vision*, 75(1):93–113.
- Murphy-Chutorian, E. and Trivedi, M. (2009). Head pose estimation in computer vision: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(4):607–626.

- Phillips, P. J., Scruggs, W. T., O'Toole, A. J., Flynn, P. J., Bowyer, K. W., Schott, C. L., and Sharpe, M. (2007). FRVT 2006 and ICE 2006 Large-Scale Results. Technical Report NISTIR 7408, National Institute of Standards and Technology, Gaithersburg, MD.
- Stegmann, M. B., Fisker, R., Ersbøll, B. K., Thodberg, H. H., and Hyldstrup, L. (2000). Active appearance models: Theory and cases. In *Proc. 9th Danish Conf. Pattern Recognition and Image Analysis*, pages 49–57.
- Zhao, S. and Grigat, R.-R. (2006). Robust eye detection under active infrared illumination. In *Proc. 18th International Conference on Pattern Recognition (ICPR 2006)*, pages 481–484.
- Zhao, S., Kricke, R., and Grigat, R.-R. (2007). Tunir: A multi-modal database for person authentication under near infrared illumination. In *Proc. 6th WSEAS International Conference on Signal Processing, Robotics and Automation (ISPRA 2007)*, Corfu, Greece.
- Zhao, W., Chellappa, R., Phillips, P. J., and Rosenfeld, A. (2003). Face recognition: A literature survey. *ACM Computing Surveys*, 35(4):399–458.



SciTeP  
Science and Technology Publications