

Face Alignment Through 2.5D Active Appearance Models

Pedro Martins
pedromartins@isr.uc.pt
Rui Caseiro
ruicaseiro@isr.uc.pt
Jorge Batista
batista@isr.uc.pt

Institute of Systems and Robotics
University of Coimbra
Coimbra, Portugal

This work addresses the fitting of 3D deformable face models from a single view through 2.5D Active Appearance Models (AAM) [4]. The main contribution of this paper is the use of 2.5D AAM that combines a 3D *metric* Point Distribution Model (PDM) and a 2D appearance model whose control points are defined by *full perspective* projections of the PDM. The advantage is that, assuming a calibrated camera, 3D metric shapes can be retrieved from single view images.

1 2.5D Parametric Models

The shape of a non-rigid object can be expressed by a linear combination of a set of n basis shapes (stored in a matrix Φ) i.e. a PDM. A 3D v -point shape is defined by $s = (X_1, \dots, X_v, Y_1, \dots, Y_v, Z_1, \dots, Z_v)^T$ and the 3D PDM, including the full pose variation, is given by

$$s = s_0 + \sum_{i=1}^n p_i \phi_i + \sum_{j=1}^6 q_j \psi_j + \underbrace{\int_0^{t-1} \sum_{j=1}^6 q_j \psi_j dt}_{s_\psi} \quad (1)$$

where \mathbf{p} are the shape parameters, $\mathbf{q} = [w_x, w_y, w_z, t_x, t_y, t_z]^T$ are the pose parameters and s_ψ is the contribution of pose increments over time t . ψ_1, \dots, ψ_6 are a special set of eigenvectors that are expressed w.r.t. the base mesh, s_0 , and derived from first order approximations of the Rodrigues rotation formula (becoming only valid for small changes in pose).

Using a *full perspective* camera, the 3D shape s is projected into the image space as $\mathbf{x}_p = \mathbf{K} [\mathbf{R}_0 \mid \mathbf{t}_0] s$ where \mathbf{K} is the camera matrix, assumed to be known and $\mathbf{R}_0 | \mathbf{t}_0$ is the *Base Pose* - the head reference.

The 2D appearance model is built by texture-warp all the training images into a common reference using a warping function \mathbf{W} . The warp $\mathbf{W}(\mathbf{x}_p, \mathbf{p}, \mathbf{q})$ is a piecewise affine warp and is a function of the shape and pose parameters that defines the 2D texture control points by means of the perspective projection of the mesh s . The appearance model given by $\mathbf{A}(\mathbf{x}_p) = \mathbf{A}_0(\mathbf{x}_p) + \sum_{i=1}^{m+2} \lambda_i \mathbf{A}_i(\mathbf{x}_p)$, $\mathbf{x}_p \in s_{0p}$ where λ is a m dimensional vector of appearance parameters. Two extra eigen images are used to model illumination gain and offset.

2 Model Fitting

Fitting the 2.5D AAM consists in solving

$$\sum_{\mathbf{x}_p \in s_{0p}} \left[\mathbf{A}_0(\mathbf{x}_p) + \sum_{i=1}^{m+2} \lambda_i \mathbf{A}_i(\mathbf{x}_p) - \mathbf{I}(\mathbf{W}(\mathbf{x}_p, \mathbf{p}, \mathbf{q})) \right]^2 \quad (2)$$

simultaneously for \mathbf{p} , \mathbf{q} and λ respectively. $\mathbf{I}(\mathbf{W}(\mathbf{x}_p, \mathbf{p}, \mathbf{q}))$ represents the input image $\mathbf{I}(\mathbf{x}_p)$ warped by $\mathbf{W}(\mathbf{x}_p, \mathbf{p}, \mathbf{q})$. Since the Inverse Compositional approach was proved in [2] to be invalid for the 2.5D AAM, the additive formulation proposed by Lucas-Kanade[3] was adopted. Eq.2 can be solved by the Simultaneous Forwards Additive (SFA) using additive updates to the parameters as

$$\Delta \mathbf{r} = \mathbf{H}_{\text{sfa}}^{-1} \sum_{\mathbf{x}_p \in s_{0p}} \mathbf{SD}(\mathbf{x}_p)_{\text{sfa}}^T \mathbf{E}(\mathbf{x}_p)_{\text{sfa}} \quad (3)$$

where $\mathbf{r} = [\mathbf{p}^T \ \mathbf{q}^T \ \lambda^T]^T$, \mathbf{H}_{sfa} is the Hessian matrix,

$$\mathbf{SD}(\mathbf{x}_p)_{\text{sfa}} = \left[\nabla \mathbf{I} \frac{\partial \mathbf{W}}{\partial \mathbf{p}_1} \dots \nabla \mathbf{I} \frac{\partial \mathbf{W}}{\partial \mathbf{p}_n} \nabla \mathbf{I} \frac{\partial \mathbf{W}}{\partial \mathbf{q}_1} \dots \nabla \mathbf{I} \frac{\partial \mathbf{W}}{\partial \mathbf{q}_6} \mathbf{A}_1(\mathbf{x}_p) \dots \mathbf{A}_{m+2}(\mathbf{x}_p) \right] \quad (4)$$

are the Steepest Descent images that depend on \mathbf{p} and \mathbf{q} (by the Jacobian of the Warp $\frac{\partial \mathbf{W}}{\partial \mathbf{p}}$ and $\frac{\partial \mathbf{W}}{\partial \mathbf{q}}$) and $\mathbf{E}(\mathbf{x}_p)_{\text{sfa}}$ is the Error image.

The Normalization Forwards Additive (NFA) algorithm solves eq.2 by projecting out the appearance images $\mathbf{A}_i(\mathbf{x}_p)$ from the error image, searching only for the shape and pose parameters. The parameters update is given by $\begin{bmatrix} \Delta \mathbf{p} \\ \Delta \mathbf{q} \end{bmatrix} = \mathbf{H}_{\text{nfa}}^{-1} \sum_{\mathbf{x}_p \in s_{0p}} \mathbf{SD}(\mathbf{x}_p)_{\text{nfa}}^T \mathbf{E}(\mathbf{x}_p)_{\text{nfa}}$, where the $\mathbf{SD}(\mathbf{x}_p)_{\text{nfa}}$ are far less dimensional than the equivalent for SFA algorithm.

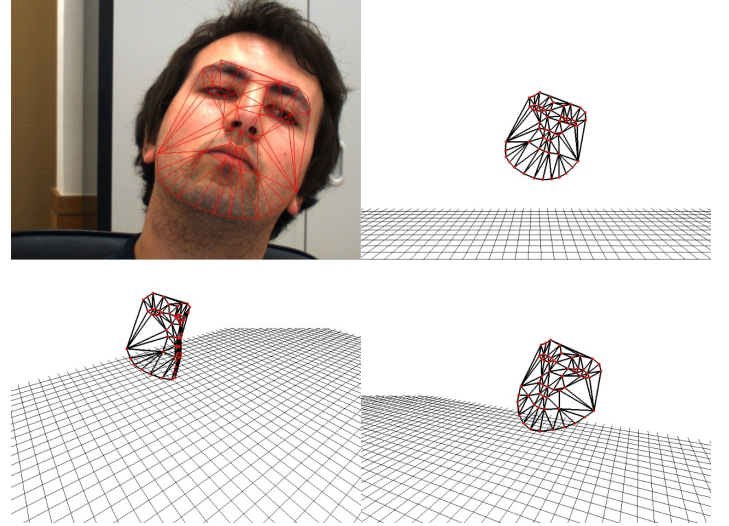


Figure 1: Example of 2.5D AAM fitting.

Some computational load can be reduced by eliminating the need to recompute image gradients at each iteration. Following the idea proposed by [1] we can use the approximation $\nabla \mathbf{I}(\mathbf{W}(\mathbf{x}_p, \mathbf{p}, \mathbf{q})) \approx \nabla \mathbf{A}_0(\mathbf{x}_p) + \sum_{i=1}^{m+2} \lambda_i \nabla \mathbf{A}_i(\mathbf{x}_p)$ which, besides providing extra computation efficiency (the gradients of the template in eq.4 can be precomputed), it has the great advantage of providing better stability to noise sensitivity since it avoids the reevaluation of the gradients in the input image.

Self-occlusion can be modeled as outlier pixels in the appearance model and handled by robust fitting methods. Robust fitting seek to minimize $\sum_{\mathbf{x}_p} \rho(\mathbf{E}(\mathbf{x}_p)_{\text{sfa}}^2, \sigma)$ where $\rho(\cdot)$ is a robust function and σ is the scale parameter. Using Back-face Culling invisible triangles by the camera can be dropped, setting them as outliers and not taking them into consideration in the fitting process.

2.1 The Jacobian of the Warp

The Jacobian for the shape parameters can be decomposed by the chain rule as $\frac{\partial \mathbf{W}(\mathbf{x}_p, \mathbf{p}, \mathbf{q})}{\partial \mathbf{p}} = \sum_{k=1}^v \left[\frac{\partial \mathbf{W}(\mathbf{x}_p, \mathbf{p}, \mathbf{q})}{\partial \mathbf{x}_k} \frac{\partial \mathbf{x}_k}{\partial \mathbf{p}} + \frac{\partial \mathbf{W}(\mathbf{x}_p, \mathbf{p}, \mathbf{q})}{\partial \mathbf{y}_k} \frac{\partial \mathbf{y}_k}{\partial \mathbf{p}} \right]$. The $\frac{\partial \mathbf{W}(\mathbf{x}_p, \mathbf{p}, \mathbf{q})}{\partial \mathbf{x}_k}$ and $\frac{\partial \mathbf{W}(\mathbf{x}_p, \mathbf{p}, \mathbf{q})}{\partial \mathbf{y}_k}$ components are given by $(1 - \alpha - \beta, 0)$ and $(0, 1 - \alpha - \beta)$ respectively, where α and β are barycentric coordinates of the projected base mesh s_{0p} . They depend only on the configuration of the base mesh and thus can be precomputed and efficiently stored as sparse matrices. The same approach is taken to evaluate the Jacobian of the warp for the pose parameters $\frac{\partial \mathbf{W}(\mathbf{x}_p, \mathbf{p}, \mathbf{q})}{\partial \mathbf{q}}$. The remaining terms $\frac{\partial \mathbf{x}_k}{\partial \mathbf{p}}$, $\frac{\partial \mathbf{y}_k}{\partial \mathbf{p}}$ and $\frac{\partial \mathbf{x}_k}{\partial \mathbf{q}}$, $\frac{\partial \mathbf{y}_k}{\partial \mathbf{q}}$ are all *scalars* given by

$$\begin{bmatrix} w_{x_k} \\ w_{y_k} \\ w_z \end{bmatrix} = \mathbf{K} [\mathbf{R}_0 \mid \mathbf{t}_0] \begin{bmatrix} s_0^x + p_i \phi_i^{xk} \\ s_0^y + p_i \phi_i^{yk} \\ s_0^z + p_i \phi_i^{zk} \\ 1 \end{bmatrix} \quad \left| \quad \begin{bmatrix} w_{x_k} \\ w_{y_k} \\ w_z \end{bmatrix} = \mathbf{K} [\mathbf{R}_0 \mid \mathbf{t}_0] \begin{bmatrix} s_0^x + q_i \psi_i^{xk} + s_\psi^x \\ s_0^y + q_i \psi_i^{yk} + s_\psi^y \\ s_0^z + q_i \psi_i^{zk} + s_\psi^z \\ 1 \end{bmatrix} \right. \quad (5)$$

$i = 1, \dots, n$ parameters; $k = 1, \dots, v$ landmarks $\frac{\partial \mathbf{x}_k}{\partial \mathbf{p}} = \frac{\partial}{\partial \mathbf{p}} \left(\frac{w_{x_k}}{w_z} \right)$ and $\frac{\partial \mathbf{y}_k}{\partial \mathbf{p}} = \frac{\partial}{\partial \mathbf{p}} \left(\frac{w_{y_k}}{w_z} \right)$ $\quad \quad \quad i = 1, \dots, 6$ parameters; $k = 1, \dots, v$ landmarks $\frac{\partial \mathbf{x}_k}{\partial \mathbf{q}} = \frac{\partial}{\partial \mathbf{q}} \left(\frac{w_{x_k}}{w_z} \right)$ and $\frac{\partial \mathbf{y}_k}{\partial \mathbf{q}} = \frac{\partial}{\partial \mathbf{q}} \left(\frac{w_{y_k}}{w_z} \right)$.

- [1] G.Hager and P.Belhumeur. Efficient region tracking with parametric models of geometry and illumination. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(10):1025–39, October 1998.
- [2] I.Matthews, J.Xiao, and S.Baker. 2d vs. 3d deformable face models: Representational power, construction, and real-time fitting. *International Journal of Computer Vision*, 75(1):93–113, October 2007.
- [3] S.Baker and I.Matthews. Lucas-kanade 20 years on: A unifying framework. *International Journal of Computer Vision*, 56(1):221–255, March 2004.
- [4] T.F.Cootes, G.J.Edwards, and C.J.Taylor. Active appearance models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 23(6):681–685, June 2001.