

A 3D Morphable Eye Region Model for Gaze Estimation

Erroll Wood¹ Tadas Baltrušaitis² Louis-Philippe Morency² Peter Robinson¹ Andreas Bulling³

¹University of Cambridge, United Kingdom, {erroll.wood, peter.robinson}@cl.cam.ac.uk

²Carnegie Mellon University, United States, {tbaltrus, morency}@cs.cmu.edu

³Max Planck Institute for Informatics, Germany, bulling@mpi-inf.mpg.de

Abstract

Morphable models are a powerful tool, but have so far failed to model the eye accurately. We present a multi-part model that includes a new 3DMM of the facial eye region, and an anatomy-based eyeball. By fitting this model to an image, we can estimate 3D gaze robustly.

Model parameters Φ

We parameterize our multi-part model with

$$\Phi = \{\beta, \tau, \theta, \iota, \kappa\}$$

Shape β and texture τ are controlled with linear models. Pose θ is defined by model transforms and procedural animation. We use ambient + directional illumination ι . We assume knowledge of camera parameters κ .

Fitting our model

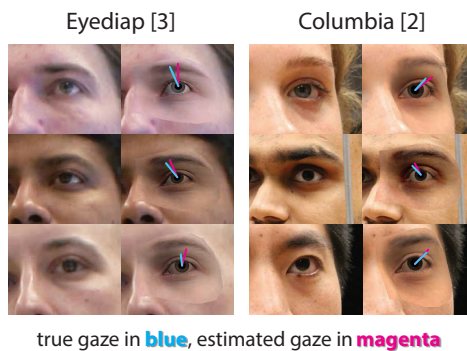
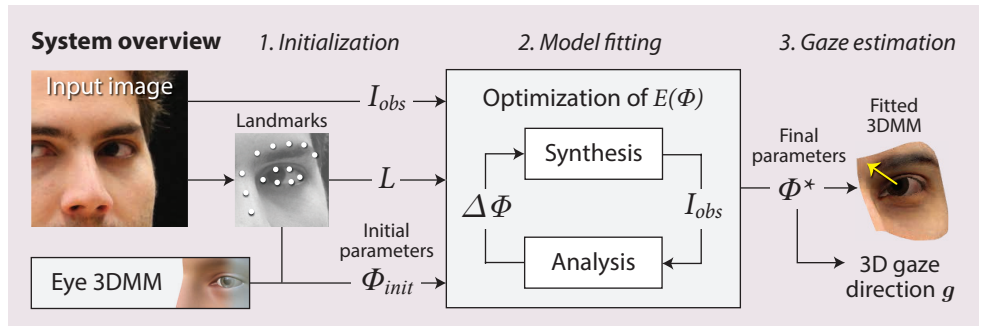
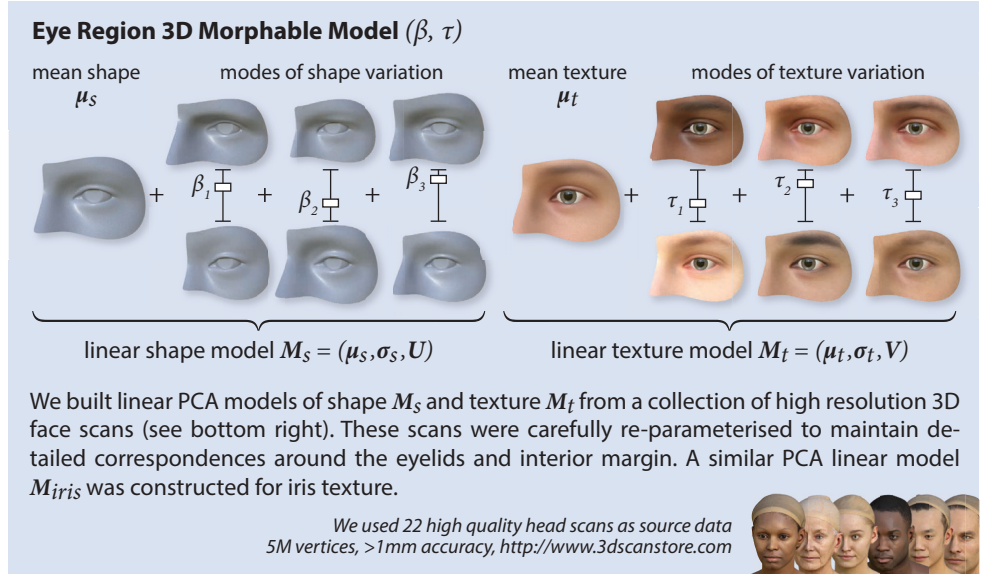
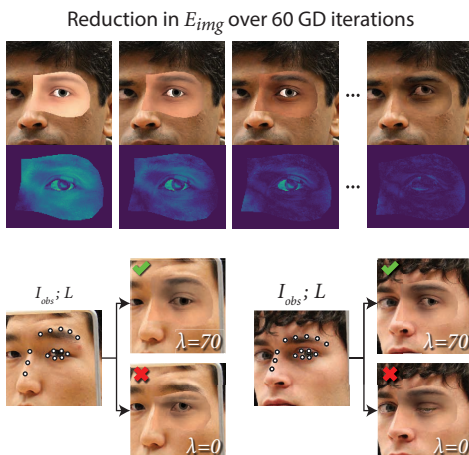
We fit our model with *analysis-by-synthesis*: given an observed image I_{obs} , we produce a synthesized image I_{syn} that best matches it. We search for optimal parameters Φ^* as follows:

$$\Phi^* = \underset{\Phi}{\operatorname{argmin}} E(\Phi)$$

Our energy includes E_{img} that measures the dense pixel-wise similarity between I_{obs} and I_{syn} , and E_{ldmks} that regularizes our model against tracked facial feature landmarks [1]. λ controls their relative importance.

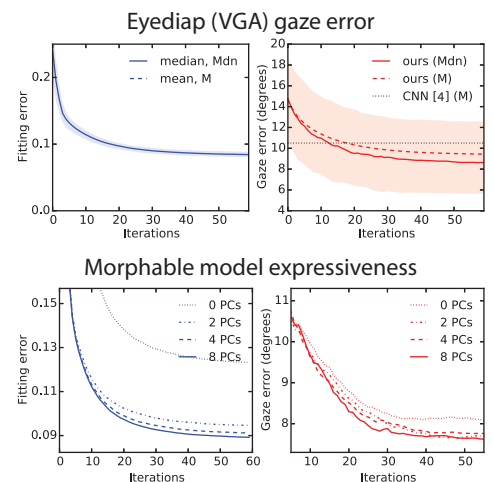
$$E(\Phi) = E_{img}(\Phi) + \lambda \cdot E_{ldmks}(\Phi, L)$$

We minimize $E(\Phi)$ using gradient descent with numerical central derivatives, solving for all parameters simultaneously.



Gaze estimation results

We evaluated on the Columbia [2] and Eyediap [3] datasets. We (**9.44°** gaze error) out-performed a state of the art convolutional neural network method [4] (10.5°). Our model based approach also generalizes to extreme gaze angles not covered in their limited training set.



Top: a comparison against a state-of-the-art CNN method [4] Bottom: more PCs lets our model fit better, and gives better gaze error.

[1] T. Baltrušaitis, P. Robinson, and L.-P. Morency, Constrained local neural fields for robust facial landmark detection in the wild, ICCVW, 2013

[2] Smith, B., Yin, Q., Feiner, S., Nayar, S. Gaze Locking: Passive Eye Contact Detection for Human-Object Interaction, UIST 2013

[3] Funes Mora, K.A., Monay, F., Odobez, J.M. EYEDIAP, ETRA 2014

[4] X. Zhang, Y. Sugano, M. Fritz, and A. Bulling, Appearance-Based Gaze Estimation in the Wild, CVPR, 2015.