


graphics and vision **gravis**  University of Basel

# 2D Face Image Analysis


Probabilistic Morphable Model Fitting  
Basel2019

University of Basel

1

UNIVERSITÄT BASEL  
> DEPARTMENT OF MATHEMATICS AND COMPUTER SCIENCE GRAVIS2019: BASEL

## Modeling of 2D Images



# Morphable Models for Image Registration



$$= R_{\rho} \left\{ \begin{array}{l} \alpha_1 \text{ (face)} + \alpha_2 \text{ (face)} + \alpha_3 \text{ (face)} + \dots \\ \beta_1 \text{ (face)} + \beta_2 \text{ (face)} + \beta_3 \text{ (face)} + \dots \end{array} \right\}$$

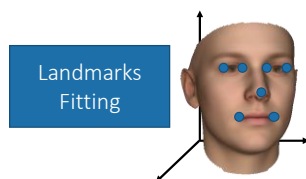
$R$  = Rendering Function

$\rho$  = Parameters for Pose, Illumination, ...

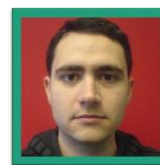
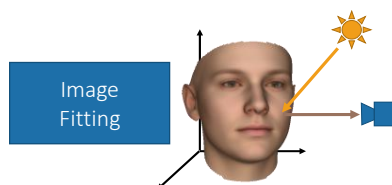
**Optimization Problem:** Find optimal  $\alpha$ ,  $\beta$ ,  $\rho$  !



## Contents



Observed Landmarks in 2D



Observed Image

## 2D Face Image Analysis

Morphable Model adaptation to explain image  
*Bayesian Inference Setup*

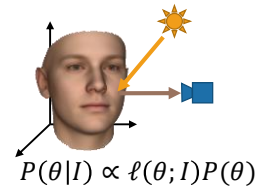
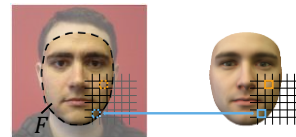


Image Likelihood  
*Image as observation*



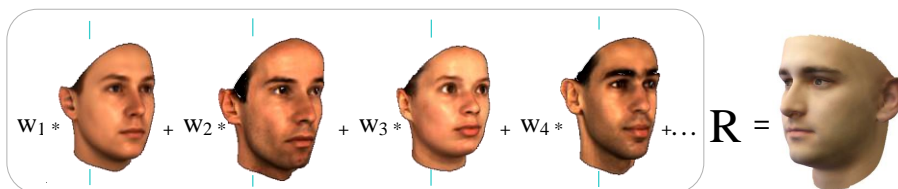
Face & Feature point detection  
*Integration of fast bottom-up methods*



## Computer Graphics: Rendering Faces

3D Face Scans

2D Images



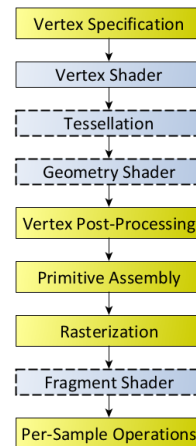
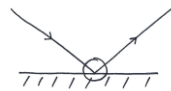
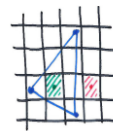
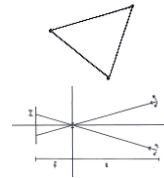
Faces: GP models for shape & color:

$$s_\alpha = \mu + UD\alpha \quad \alpha \sim N(0, I_d)$$

$$c_\beta = \mu + UD\beta \quad \beta \sim N(0, I_d)$$

# Computer Graphics Overview

- **Geometry** (result of shape modelling)
- **Camera & Projection**  
Transformations in space and projection  
Maps 3D space and 2D image plane
- **Rasterization**  
Correspondence: image pixels ↔ surface  
Z-Buffer: Hidden surface removal
- **Shading**  
Illumination simulation models
  - **Illumination**  
Phong: Ambient, diffuse & specular  
Global Illumination



7

# Face-to-Image Transformations

- **Model-View**

$$T_{MV}(x) = R_{\varphi, \psi, \vartheta}(x) + t$$

- **Projection**

$$\mathcal{P}(x) = \frac{f}{z} \begin{bmatrix} x \\ y \end{bmatrix}$$

- **Viewport**

$$T_{VP}(x) = \begin{bmatrix} \frac{w}{2}(x+1) \\ \frac{h}{2}(1-y) \end{bmatrix} + t_{pp}$$

- 9 Parameters:

- (3) Translation  $t$
- (3) Rotation  $\varphi, \psi, \vartheta$
- (1) Focal length  $f$
- (2) Image Offset  $t_{pp}$

- 2 Constants:

- (2) Image size / sampling

8

## Perspective Effect

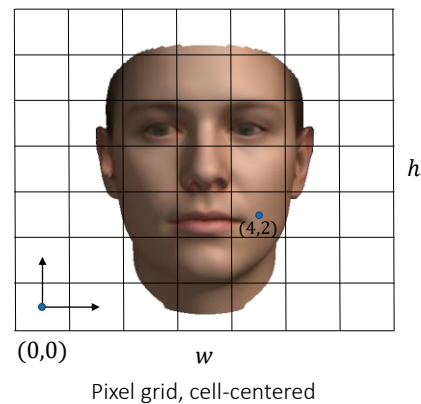
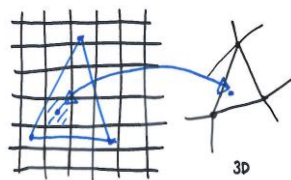
- Perspective division distorts image non-linearly
- Effect depends on relation of object depth and camera distance



9

## Rasterization

- Camera:  $3D \rightarrow 2D$  transformation for *points*
- Raster Image in image plane
- Establishes correspondence to 3D surface for each *pixel*
- Basis: geometric *primitives*

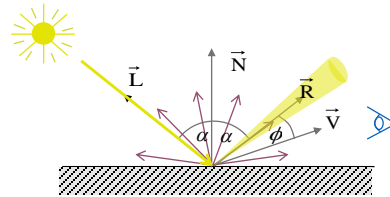


10

## Phong Illumination Model

- Combination of three illumination contributions:

- Lambert (diffuse)  $k_{\text{diff}} * I_L * \cos(L, N)$
- Specular  $k_{\text{spec}} * I_L * \cos(R, V)^n$
- Ambient (global)  $k_{\text{amb}} * I_A$



- Ambient is a scene *average* light intensity  $I_A$
- Lambert and specular part for each light source

$$I' = k_{\text{amb}} * I_A + k_{\text{diff}} * I_L * \cos(L, N) + k_{\text{spec}} * I_L * \cos(R, V)^n$$

usually colored

11

## Phong Illumination Model

- Combination of three illumination contributions:

- Lambert (diffuse)  $k_{\text{diff}} * I_L * \cos(L, N)$
- Specular  $k_{\text{spec}} * I_L * \cos(R, V)^n$
- Ambient (global)  $k_{\text{amb}} * I_A$

- Ambient is a scene *average* light intensity  $I_A$
- Lambert and specular part for each light source

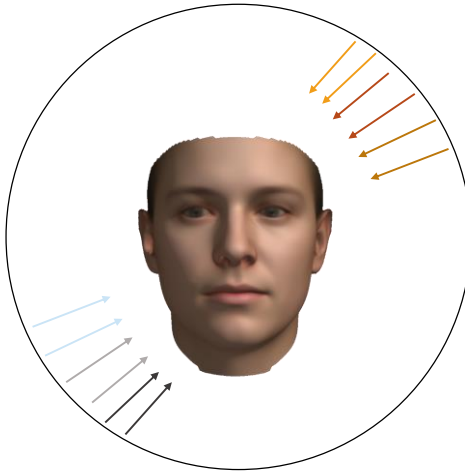


$$I' = k_{\text{amb}} * I_A + k_{\text{diff}} * I_L * \cos(L, N) + k_{\text{spec}} * I_L * \cos(R, V)^n$$

usually colored

12

# Environment Maps

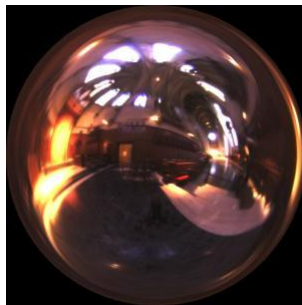


- Mapping of incoming light intensity from every direction  

$$I_L^{\text{RGB}}(\theta, \varphi)$$
- Modeled at infinity
- Typically *empirically* captured
- Shading with environment maps requires *integration* over all incoming directions

13

# Environment Maps



Grace Cathedral (San Francisco)  
P. Debevec

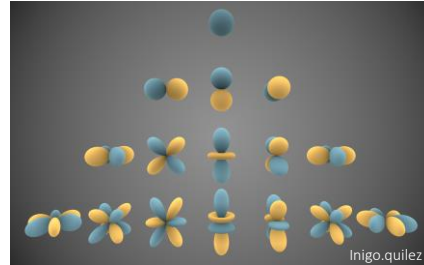


White surface in Grace Cathedral

14

## Spherical Harmonics Illumination

- Expand map  $I_L^{RGB}(\theta, \varphi)$  with *basis* functions
- Choose *Spherical Harmonics*: Eigenfunctions of Laplace operator on sphere surface  $Y_{lm}(\theta, \varphi)$
- Corresponds to Fourier transform
- Integration becomes multiplication of coefficients ( $\rightarrow$  *fast convolution*)
- Low frequency part is sufficient for Lambertian reflectance



Ramamoorthi, Ravi, and Pat Hanrahan. "An efficient representation for irradiance environment maps." Proceedings of the 28th annual conference on Computer graphics and interactive techniques. ACM, 2001.

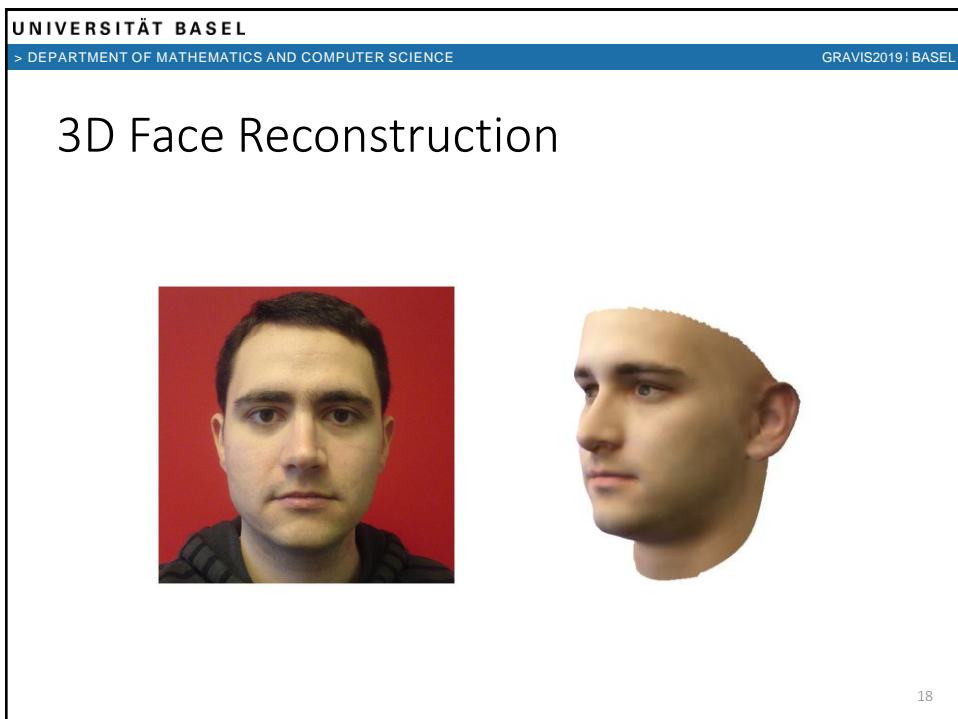
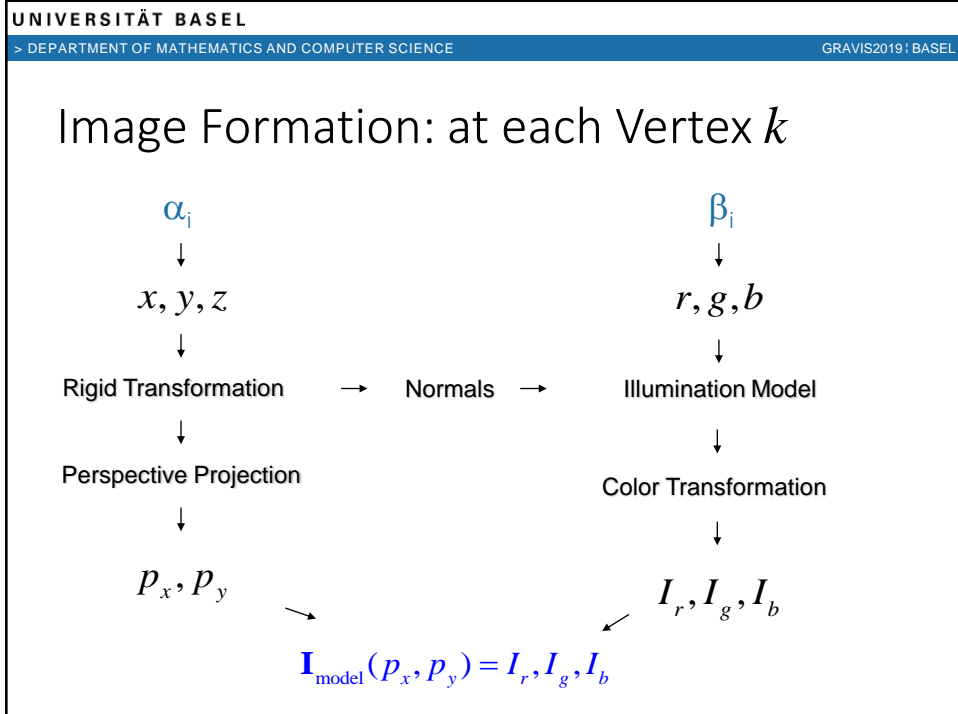
15

## Environment Map Illumination



16





## Probabilistic Inference for Image Registration

- Generative image explanation: How to find  $\theta$  explaining  $I$ ?

$$p(\theta|I) = \frac{\ell(\theta; I) p(\theta)}{N(I)} \quad N(I) = \int \ell(\theta; I) p(\theta) d\theta$$

-----> Normalization intractable in our setting

- What can be done:
  1. Accept MAP as the only option
  2. Approximate posterior distribution (e.g. use sampling methods)

## MH Inference of the 3DMM

- Target distribution is our “posterior”:

$$P: \tilde{P}(\theta|I) = \ell(\theta; I) P(\theta)$$

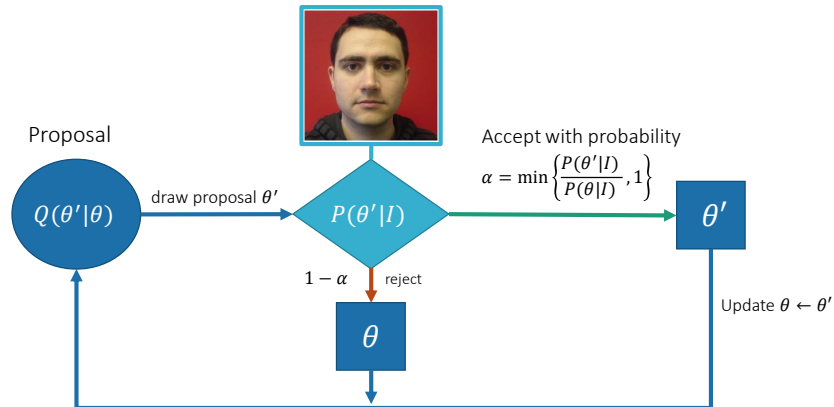
- Unnormalized
- Point-wise evaluation only

- **Parameters**

- Shape: 50 – 200, low-rank parameterized GP shape model
- Color: 50 – 200, low-rank parameterized GP color model
- Pose/Camera: 9 parameters, pin-hole camera model
- Illumination: 9\*3 Spherical Harmonics illumination/reflectance

$\approx$  300 dimensions (!!)

# Metropolis Algorithm



- Asymptotically generates samples  $\theta_i \sim P(\theta|I)$ :  $\theta_1, \theta_2, \theta_3, \dots$
- Markov chain Monte Carlo (MCMC) Method
- Works with *unnormalized*, point-wise posterior

23

# Proposals

- Choose simple Gaussian random walk proposals (Metropolis)  
 $Q(\theta'|\theta) = N(\theta'|\theta, \Sigma_\theta)$

- Normal *perturbations* of current state

- Block-wise to account for different parameter types

- Shape  $N(\alpha'|\alpha, \sigma_s^2 E_s)$
- Color  $N(\beta'|\beta, \sigma_c^2 E_c)$
- Camera  $\Sigma_c N(\theta'_c|\theta_c, \sigma_c^2)$
- Illumination  $\Sigma_i N(\theta'_L|\theta_L, \sigma_{L,i}^2 E_L)$



In practice, we often add more complicated proposals, e.g. shape scaling, a direct illumination estimation and decorrelation

- Large mixture distributions, e.g.

$$\frac{2}{3} Q_P(\theta'|\theta) + \frac{1}{3} \sum_i \lambda_i Q_i^L(\theta'|\theta)$$

24

UNIVERSITÄT BASEL  
 > DEPARTMENT OF MATHEMATICS AND COMPUTER SCIENCE  
 GRAVIS2019: BASEL

## Landmarks Fitting

Face Model

Prior  $P(\theta)$

Projection

Variable Parameters

- Pose
- Shape

Rendered Landmarks

Target Landmarks

Likelihood  $\ell(\theta; \tilde{\mathbf{x}}) \propto P(\tilde{\mathbf{x}}|\mathbf{x}(\theta))$

UNIVERSITÄT BASEL  
 > DEPARTMENT OF MATHEMATICS AND COMPUTER SCIENCE  
 GRAVIS2019: BASEL

## 3DMM Landmarks Likelihood

Simple models: **Independent Gaussians**

- Observation of landmark locations in image
  - Single landmark position model:

$$\mathbf{x}_i^{2D}(\theta) = (\mathbf{T}_{VP} \circ \text{Pr} \circ \mathbf{T}_{MV})(\mathbf{x}_i^{3D})$$

$$\ell_i(\theta; \tilde{\mathbf{x}}_i^{2D}) = N(\tilde{\mathbf{x}}_i^{2D} | \mathbf{x}_i^{2D}(\theta), \sigma_{LM}^2)$$

$$\mathbf{T}_{MV}(\mathbf{x}) = \mathbf{R}_{\varphi, \psi, \theta}(\mathbf{x}) + \mathbf{t}$$

$$(\mathbf{T}_{VP} \circ \text{Pr})(\mathbf{x}) = \begin{bmatrix} \frac{w}{2} * \frac{x}{z} \\ h \\ -\frac{w}{2} * \frac{y}{z} \end{bmatrix} + \mathbf{t}_{vp}$$

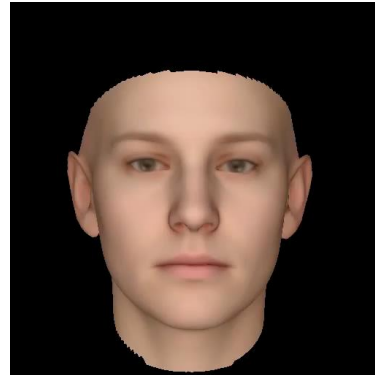
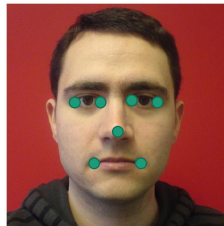
- Independent model

$$\ell(\theta; \{\tilde{\mathbf{x}}_i^{2D}\}_i) = \prod_i \ell(\theta; \tilde{\mathbf{x}}_i^{2D})$$

Independence and Gaussian are just *simple models* (questionable)

27

## Landmarks: Samples



28

## Results: 2D Landmarks

- Landmarks posterior:  
Manual labelling:  $\sigma_{LM} = 4\text{pix}$   
Image: 512x512
- Certainty of pose fit?
  - Influence of ear points?
  - Frontal better than side-view?

- Landmarks posterior:  
Manual labelling:  $\sigma_{LM} = 4\text{pix}$   
Image: 512x512
- Certainty of pose fit?
  - Influence of ear points?
  - Frontal better than side-view?

Yaw, $\sigma_{LM} = 4\text{pix}$	with ears	w/o ears
Frontal	$1.4^\circ \pm 0.9^\circ$	$-0.8^\circ \pm 2.7^\circ$
Side view	$24.8^\circ \pm 2.5^\circ$	$25.2^\circ \pm 4.0^\circ$

UNIVERSITÄT BASEL  
 > DEPARTMENT OF MATHEMATICS AND COMPUTER SCIENCE  
 GRAVIS2019: BASEL

## Face Model Fitting

Reconstruction: Analysis-by-Synthesis

Face Model

Parametric face model

Rendered Image  $I(\theta)$

Target Image  $I$

Likelihood  $\ell(\theta; I) \propto P(I | I(\theta))$

$\theta = (\vartheta, \alpha, \beta)$ :  $\vartheta$  Scene Parameters,  $\alpha$  Face shape,  $\beta$  Face color

31

UNIVERSITÄT BASEL  
 > DEPARTMENT OF MATHEMATICS AND COMPUTER SCIENCE  
 GRAVIS2019: BASEL

## Independent Pixels Likelihood

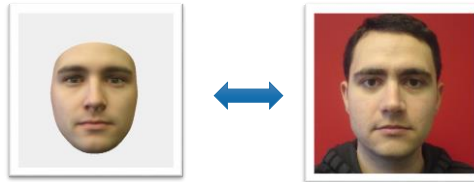
Standard choice  
 Corresponds to least squares fitting

$F$

$$\ell(\theta; \tilde{I}) = \mathcal{N}(\tilde{I} | I(\theta), \sigma^2 I_3) * \mathcal{N}(\tilde{I} | I(\theta), \sigma^2 I_3) * \dots$$

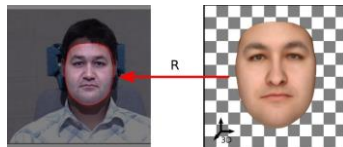
$$\ell(\theta; \tilde{I}) = \prod_{i \in F} \mathcal{N}(\tilde{I}_i | I_i(\theta), \sigma^2 I_3)$$

# Image Likelihood



## Background model is required

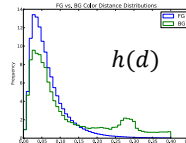
The face model does not cover the complete target image and shows self-occlusion.



## Collective likelihood model

Pixels are not independent. We can also model the empirical distribution of image distance

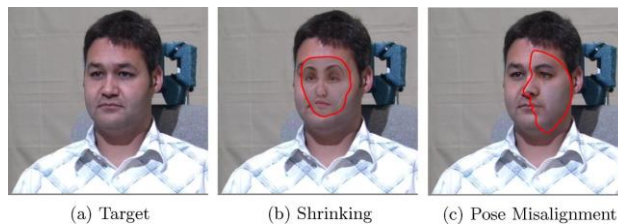
$$d = \| \text{face}_1 - \text{face}_2 \|$$



# Background Model

Schönborn et al. 2015  
«Background modeling for generative image models»,  
Computer Vision and Image Understanding, Volume 136

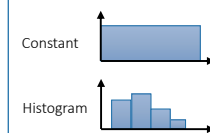
Face covers only parts of the image – background must not be ignored



- Variable alignment of model with the image
  - Projected size and self-occlusion
  - Shrinking or misalignment
- Model background pixels explicitly

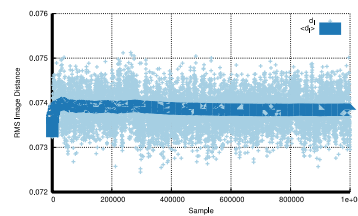
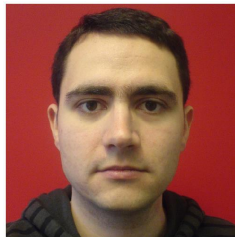
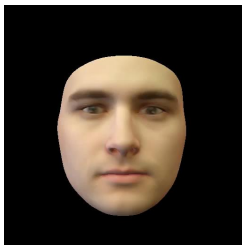
$$\ell(\theta; \tilde{I}) = \prod_{i \in F} \ell_F(\theta; \tilde{I}_i) \prod_{j \in B} b_{BG}(\tilde{I}_j)$$

Arbitrary background: The explicit background model needs to be based on *generic* and *simple* assumptions:



## Posterior Samples: Fitting Result

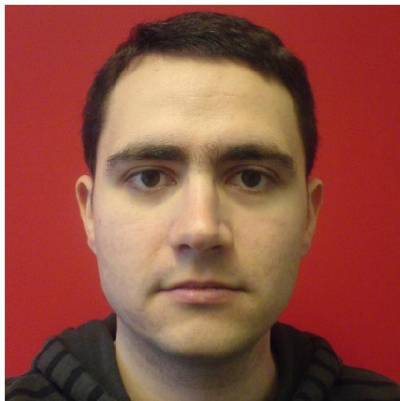
- Model instances with comparable reconstruction quality
- Remaining uncertainty of model representation
- Integration of uncertain detection directly into model adaptation



Posterior using *collective likelihood*

35

## Results: Image

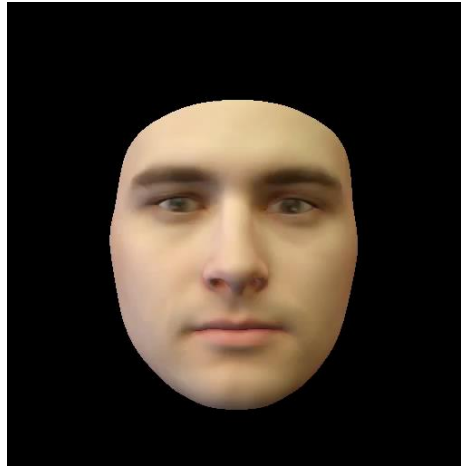


Yaw angle:  $1.9^\circ \pm 0.2^\circ$

36

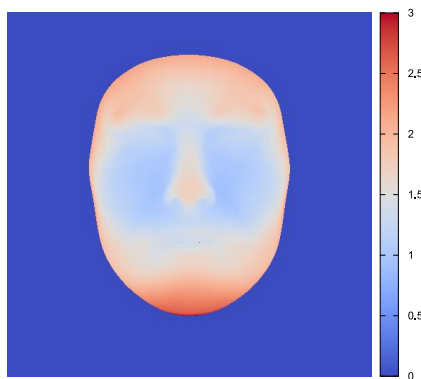


## Image: Samples



37

## Posterior Shape Variation



Landmarks posterior,  
sd[mm]

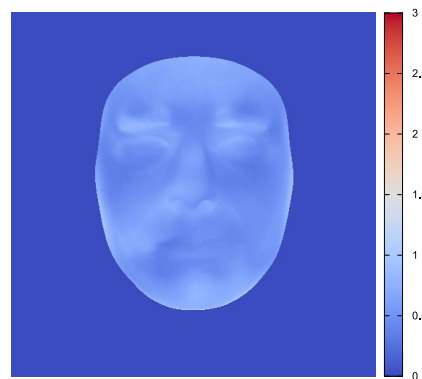
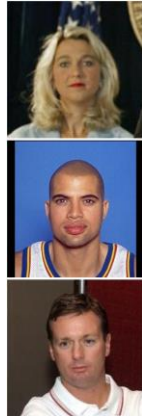


Image posterior,  
sd[mm]

38

## Fitting Results



LFW

Images from: Huang, Gary B., et al. *Labeled faces in the wild: A database for studying face recognition in unconstrained environments*. Vol. 1. No. 2. Technical Report 07-49, University of Massachusetts, Amherst, 2007.



AFLW

Images from: Köstinger, Martin, et al. "Annotated facial landmarks in the wild: A large-scale, real-world database for facial landmark localization." *Computer Vision Workshops (ICCV Workshops)*, 2011 IEEE International Conference on. IEEE, 2011.

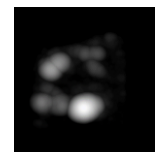
39

## Automatic Fitting

- Detection of *face* and *feature points*
  - Scanning window & classifier
  - Uncertain results
  - Feed-forward: early *hard* decisions
- Integration concept
  - Bayesian integration
    - **Filtering**
  - Metropolis sampling
    - **Propose & verify**



Which box contains the face?

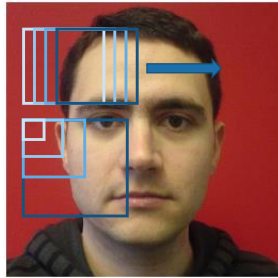


Schönborn, Sandro, et al. "Markov Chain Monte Carlo for Automated Face Image Analysis." *International Journal of Computer Vision* (2016): 1-24.

40

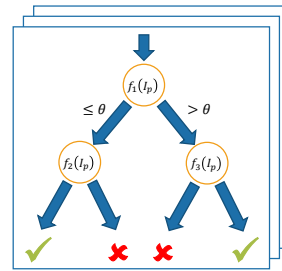
## Random Forest Detection

- Scanning Window



- Classify each patch: face or not
- Search over image
- Search over scales
- Histogram equalization

- Random Forest Classifier



- Haar Features
- Information gain splitting
- Bagging many trees, depth ~16
- ~200k training patches (AFLW)

41

## Bayesian Integration

### Detection data



### Bayesian integration

Observation likelihood

$$\ell(\theta; F, D) = P(F|\theta)P(D|\theta)$$

Bayesian inference

$$P(\theta|F, D) = \frac{\ell(\theta; F, D)P(\theta)}{N(F, D)}$$

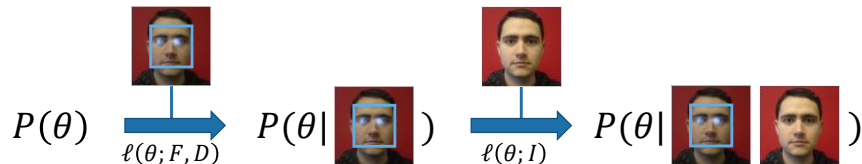
- Different *modality*
  - Box  $F$ : position & size
  - Landmarks  $D$ : certainty
- Detection is uncertain

- Likelihood* models
  - Detection is *observation*
  - Different observation models
- Conceptual uncertainty

42

## Integration by Filtering

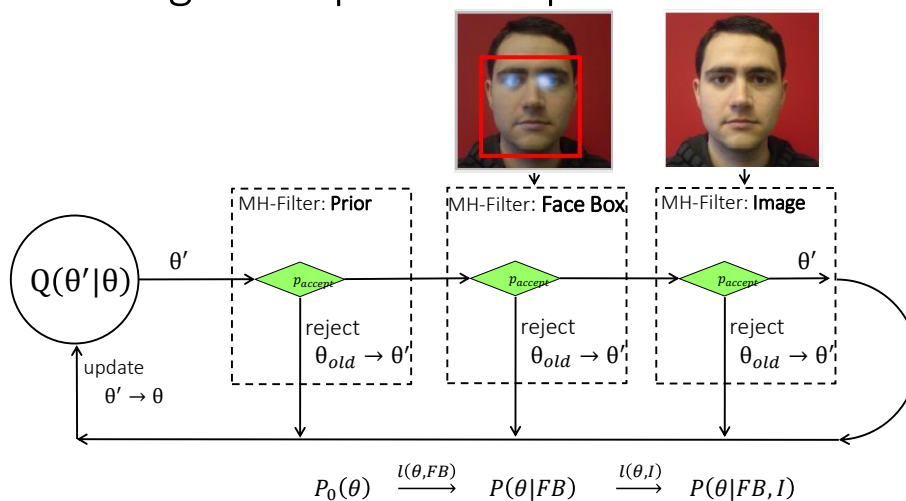
- Step-by-step Bayesian inference



- Condition on observations *one after the other*
- *Posterior* of first observation becomes *prior* for next step
  - Each step adds an observation through conditioning with its likelihood
- Equivalent to single-step Bayesian inference

44

## Filtering: Multiple Metropolis Decisions




- Step-wise Bayesian inference: Needs  $\ell(\theta)$  for each step
- Saves computation time if properly ordered

UNIVERSITÄT BASEL

> DEPARTMENT OF MATHEMATICS AND COMPUTER SCIENCE

GRAVIS2019: BASEL



by courtesy of keystone

49

UNIVERSITÄT BASEL

> DEPARTMENT OF MATHEMATICS AND COMPUTER SCIENCE

GRAVIS2019: BASEL

## Summary

- Fitting as *probabilistic inference*
- Probabilistic inference is often intractable
- Sampling methods *approximate by simulation*
- MCMC methods provide a powerful sampling framework
  - Markov Chain with target distribution as equilibrium distribution
  - General algorithms, e.g. Metropolis-Hastings
- Fitting of the 3DMM as a real inference problem
- MH algorithm to integrate information: Framework
  - *Filtering*: Uncertain information as observation, step-by-step
  - *Propose-and-verify*: Alternatives, multiple hypotheses, heuristics

50

## Occlusion-aware 3D Morphable Face Models

*Bernhard Egger, Sandro Schönborn, Andreas Schneider, Adam Kortylewski, Andreas Morel-Forster, Clemens Blumer and Thomas Vetter*  
*International Journal of Computer Vision, 2018*

51

## Face Image Analysis under Occlusion

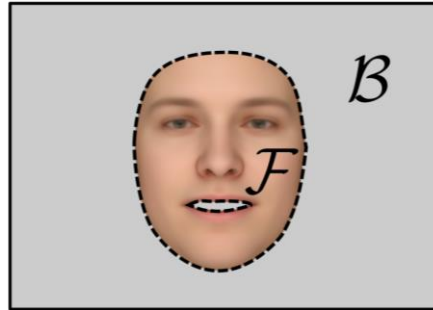


Source: AFLW Database

Source: AR Face Database

52

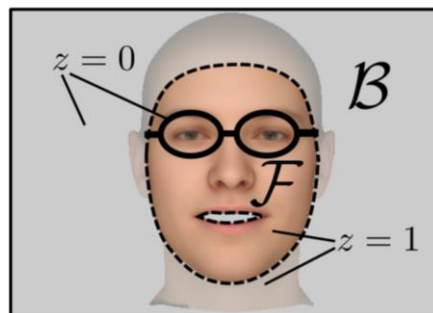
There is nothing like: no background model



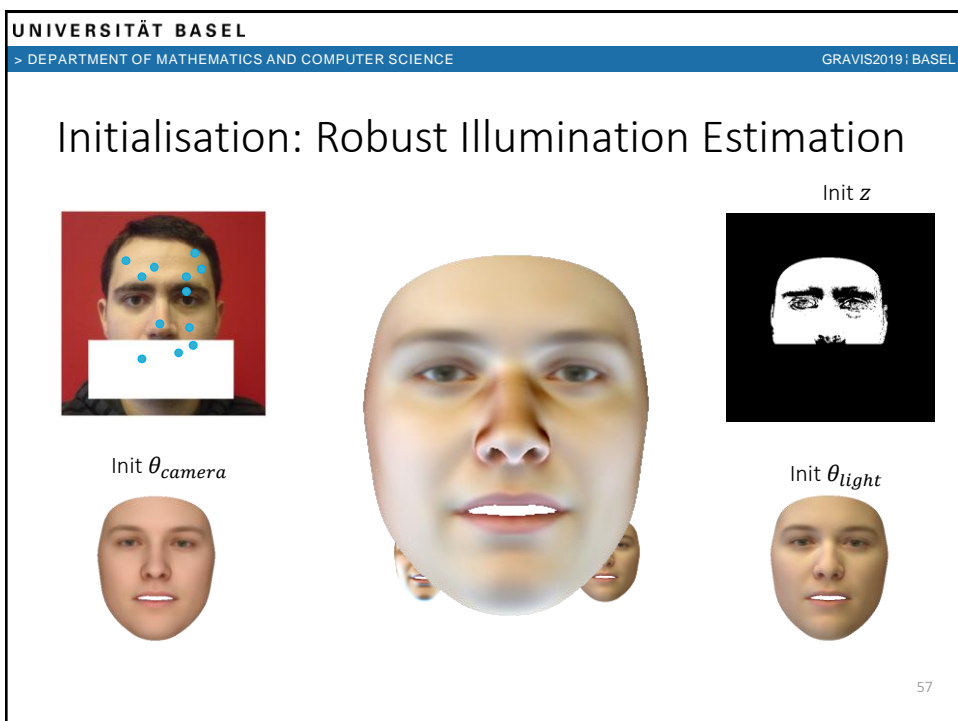
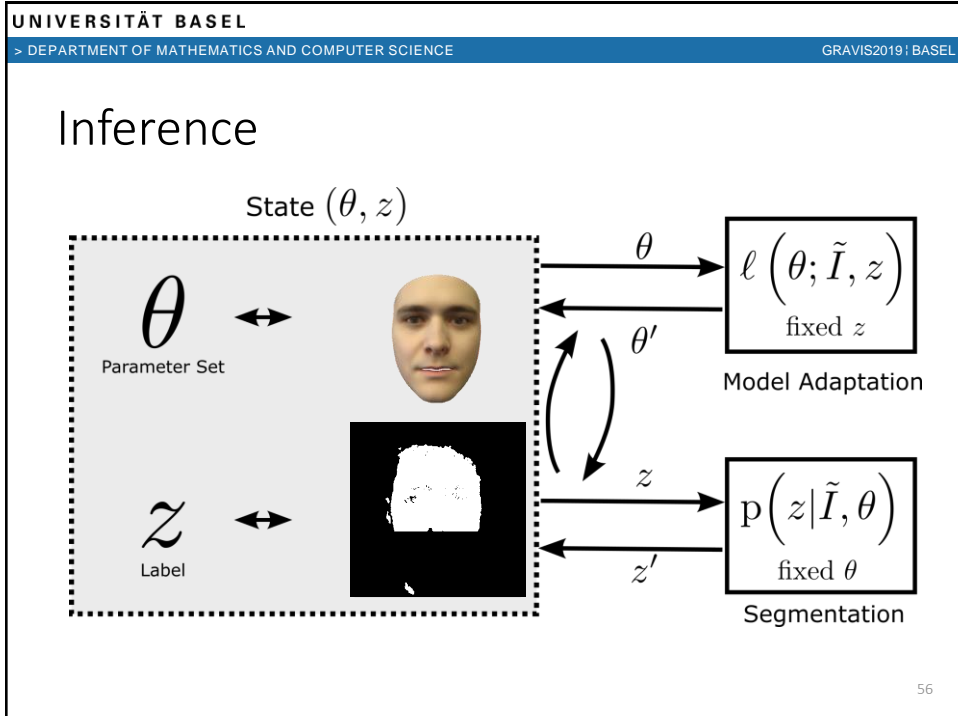
$$\ell(\theta; I) = \prod_{x \in I} \ell(\theta; I(x)) = \prod_{i \in F} l_{face}(\theta; \tilde{I}_i) \prod_{i \in B} b(\tilde{I}_i)$$

"Background Modeling for Generative Image Models"  
Sandro Schönborn, Bernhard Egger, Andreas Forster, and Thomas Vetter. Computer Vision and Image Understanding, Vol 113, 2015.

Occlusion-aware Model



$$l(\theta; \tilde{I}, z) = \prod_i l_{face}(\theta; \tilde{I}_i)^z \cdot l_{non-face}(\theta; \tilde{I}_i)^{1-z}$$

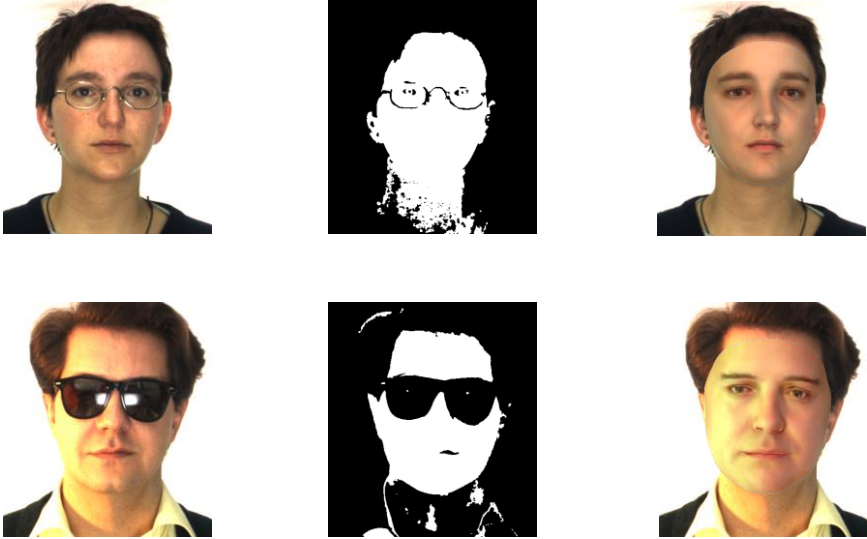




UNIVERSITÄT BASEL  
 > DEPARTMENT OF MATHEMATICS AND COMPUTER SCIENCE  
 GRAVIS2019 | BASEL

## Results: Qualitative

Source: AR Face Database

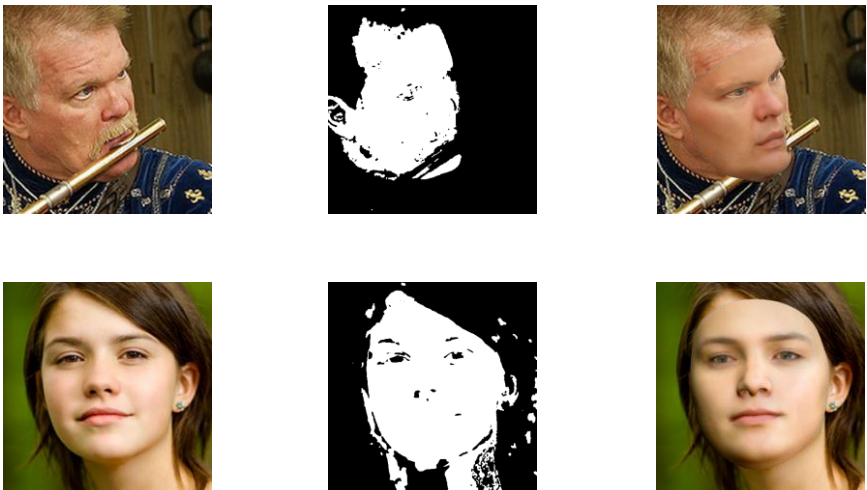


The figure displays a 2x3 grid of face images from the AR Face Database. Each row represents a different subject. The first column shows the original face image. The second column shows the segmented face image, where the face is white and the background is black. The third column shows the reconstructed face image, which is a slightly processed version of the original. The subjects are a woman with glasses and a man with sunglasses.

UNIVERSITÄT BASEL  
 > DEPARTMENT OF MATHEMATICS AND COMPUTER SCIENCE  
 GRAVIS2019 | BASEL

## Results: Qualitative

Source: AFLW Database



The figure displays a 2x3 grid of face images from the AFLW Database. Each row represents a different subject. The first column shows the original face image. The second column shows the segmented face image, where the face is white and the background is black. The third column shows the reconstructed face image, which is a slightly processed version of the original. The subjects are a man with a mustache and a woman with brown hair.

59

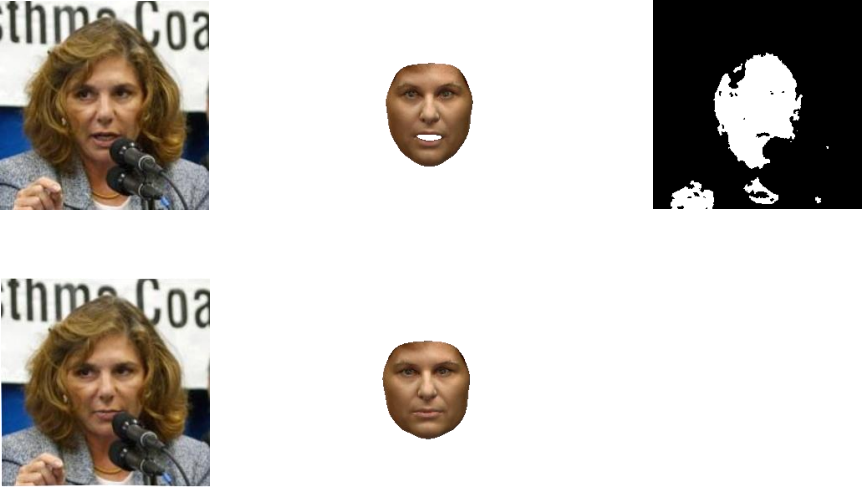
UNIVERSITÄT BASEL

> DEPARTMENT OF MATHEMATICS AND COMPUTER SCIENCE

GRAVIS2019 | BASEL

## Results: Applications

Source: LFW Database



The slide displays two rows of face images. Each row contains three images: an original photograph of a woman speaking at a microphone, a generated face image, and a binary mask of the face. The generated faces appear to be high-quality reconstructions of the original faces. The binary masks are black and white, highlighting the face area against the background.

60