

Reihe 10

Informatik/  
Kommunikation

Nr. 868

Stella Graßhof, M.Sc.,  
Kopenhagen

## Expressive Personalized 3D Face Models from 3D Face Scans



**Institut für Informationsverarbeitung**  
[www.tnt.uni-hannover.de](http://www.tnt.uni-hannover.de)

<https://doi.org/10.51201/0003186868107-1>

Generiert durch IP '18.116.19.183', am 19.04.2024, 17:42:00.

Das Erstellen und Weitergeben von Kopien dieses PDFs ist nicht zulässig.



# **Expressive Personalized 3D Face Models from 3D Face Scans**

Von der Fakultät für Elektrotechnik und Informatik  
der Gottfried Wilhelm Leibniz Universität Hannover  
zur Erlangung des akademischen Grades

**Doktor-Ingenieurin**

(abgekürzt: Dr.-Ing.)

genehmigte

**Dissertation**

von

**Stella Graßhof, M. Sc.**

geboren am 20. August 1985 in Hannover.

**2019**

Hauptreferent: Prof. Dr.-Ing. Ostermann  
Korreferent: Prof. Dr.-Ing. Rohs  
Vorsitzender: Prof. Dr.-Ing. Rosenhahn  
Tag der Promotion: 08.11.2019

# Fortschritt-Berichte VDI

Reihe 10

Informatik/  
Kommunikation

Stella Graßhof, M.Sc.,  
Kopenhagen

Nr. 868

Expressive Personalized  
3D Face Models from  
3D Face Scans



**Institut für Informationsverarbeitung**  
www.tnt.uni-hannover.de

Graßhof, Stella

## **Expressive Personalized 3D Face Models from 3D Face Scans**

Fortschr.-Ber. VDI Reihe 10 Nr. 868. Düsseldorf: VDI Verlag 2020.

216 Seiten, 57 Bilder, 6 Tabellen.

ISBN 978-3-18-386810-0, ISSN 0178-9627,

€ 76,00/VDI-Mitgliederpreis € 68,40.

**Keywords:** 3D face scans – nonrigid registration – correspondence estimation – expression intensity – tensor – factorization – statistical models – expression transfer – 3D reconstruction

**Für die Dokumentation:** 3D-Gesichts-Scan – nicht-rigide Registrierung – Korrespondenzschätzung – Intensität von Gesichtsausdrücken – Tensor – statistische Modelle – Transfer von Gesichtsausdrücken – 3D-Rekonstruktion

In this work, different methods are presented to create 3D face models from databases of 3D face scans. The challenge in this endeavour is to balance the limited training data with the high demands of various applications.

The 3D scans stem from various persons showing different expressions, with varying number of points per 3D scan and different numbers of scans per person. This data of posed facial expressions revealed substructures, which are utilised to improve the proposed model. In the process of creating and using the models, for each specific application objective quality criteria are carefully designed tailored to the task to quantify the quality.

In total four face models built from three databases are compared based on: 3D face synthesis, 3D approximation, person and expression transfer, and 3D reconstruction from 2D.

### **Bibliographische Information der Deutschen Bibliothek**

Die Deutsche Bibliothek verzeichnet diese Publikation in der Deutschen Nationalbibliographie; detaillierte bibliographische Daten sind im Internet unter [www.dnb.de](http://www.dnb.de) abrufbar.

### **Bibliographic information published by the Deutsche Bibliothek**

(German National Library)

The Deutsche Bibliothek lists this publication in the Deutsche Nationalbibliographie (German National Bibliography); detailed bibliographic data is available via Internet at [www.dnb.de](http://www.dnb.de).

© VDI Verlag GmbH · Düsseldorf 2020

Alle Rechte, auch das des auszugsweisen Nachdruckes, der auszugsweisen oder vollständigen Wiedergabe (Fotokopie, Mikrokopie), der Speicherung in Datenverarbeitungsanlagen, im Internet und das der Übersetzung, vorbehalten.

Als Manuskript gedruckt. Printed in Germany.

ISSN 0178-9627

ISBN 978-3-18-386810-0

<https://doi.org/10.51202/9783186868107-1>

Generiert durch IP '18.116.19.183', am 19.04.2024, 17:42:00.

Das Erstellen und Weitergeben von Kopien dieses PDFs ist nicht zulässig.

## Acknowledgment

First and foremost, I thank my doctoral advisor Prof. Dr.-Ing. Ostermann for giving me the opportunity to start and finish this work under his supervision at the Institut für Informationsverarbeitung (TNT), Leibniz University of Hannover, Germany. I appreciate the time dedicated for discussions, and guidance. I thank Prof. Dr.-Ing Rohs, who took the time to be my second supervisor and who contributed valuable hints to improve the final thesis. At the TNT I valued the discussion which I had with Prof. Dr.-Ing. Bodo Rosenhahn, whom I also thank for being the chair of my defense committee.

At the TNT several people helped me with software, hardware, and bureaucracy, which made my life easier and enabled me to focus on more relevant things. For lifting the burden of dealing with formalities, I thank Martin Pahl, and Thomas Wehberg, as well as the secretaries Doris Jasper-Göring, Hilke Brodersen, Pia Bank, and Melanie Huch. Many thanks go to Matthias Schuh for his hardware support, Marco Munderloh, Martin Pahl, Arne Ehlers, and Holger Meuel for resolving many software- and operating system-related issues, and who sometimes provided support during some unconventional hours.

During my time at the TNT I had the great honor to learn a lot from different fields from the nicest colleagues I could have wished for. Thank you all for making this a great place to work. Particularly I thank my former office mates and discussion partners Matthias Reso, Felix Kuhnke, Benjamin Spitschan, and Hanno Ackermann.

Finally, very special thanks go to my family, especially my parents, Frank and Veronika Graßhof, who supported me and encouraged me to pursue a career path, which feels right for me. Last but not least, I am grateful for my husband, Patrick Fließ, who believed in me, when I did not, and supported me with his patience and understanding.





# Contents

<b>Abbreviations and Nomenclature</b>	<b>XII</b>
<b>1 Introduction</b>	<b>1</b>
1.1 The Difficulty of Quality Assessment . . . . .	2
1.2 Face Models . . . . .	5
1.3 Data Preprocessing and Alignment . . . . .	8
1.4 Summary of Contributions . . . . .	10
1.5 Thesis Overview . . . . .	11
<b>2 Fundamentals</b>	<b>15</b>
2.1 Camera Models . . . . .	15
2.1.1 Orthographic Camera Model . . . . .	15
2.1.2 Weak-Perspective Camera Model . . . . .	15
2.1.3 Projective Camera Model . . . . .	16
2.2 Estimation of Camera Parameters . . . . .	17
2.3 Factorization . . . . .	20
2.3.1 Principal Component Analysis . . . . .	20
2.3.2 Whitening . . . . .	22
2.3.3 Correlation vs. Dependence . . . . .	22
2.3.4 Independent Component Analysis . . . . .	23
2.3.5 Projection Pursuit . . . . .	26
2.4 Tensor Algebra . . . . .	27
2.4.1 Notation . . . . .	27
2.4.2 High-Order Singular Value Decomposition . . . . .	30
2.5 Numerical Optimization . . . . .	31
2.5.1 Definitions . . . . .	31
2.5.2 Line-Search based Methods . . . . .	33
2.6 Generalized Canonical Time Warping . . . . .	37
<b>3 Face Databases</b>	<b>40</b>
3.1 Overview . . . . .	41

3.2	Selected Databases . . . . .	42
3.2.1	BU3DFE . . . . .	42
3.2.2	BU4DFE . . . . .	44
3.2.3	Bosphorus . . . . .	46
3.2.4	Facewarehouse . . . . .	54
3.2.5	MMI . . . . .	57
3.2.6	ADFES . . . . .	58
3.3	Conclusion . . . . .	58
<b>4</b>	<b>From 3D Face Scans to Aligned Faces</b>	<b>59</b>
4.1	Preprocessing . . . . .	60
4.1.1	Rigid Global Alignment . . . . .	60
4.1.2	Detection of Outliers . . . . .	60
4.1.3	Removing Points outside of the Face Region . . . . .	61
4.2	Spatial Alignment by nonrigid Registration . . . . .	66
4.2.1	Correspondence between Point Sets . . . . .	66
4.2.2	Nonrigid 3D Registration . . . . .	68
4.2.3	Quantifying Quality . . . . .	78
4.2.4	Experiments and Evaluation . . . . .	84
4.3	Temporal Alignment . . . . .	100
4.3.1	Quantifying Expression Intensity . . . . .	101
4.3.2	Alignment of Expression Intensities . . . . .	107
4.3.3	Applications for Proposed Expression Intensities . . . . .	110
<b>5</b>	<b>Face Models</b>	<b>114</b>
5.1	Surrey's 3D Morphable Face Model . . . . .	114
5.2	Sela's Neural Network for detailed 3D Face Reconstruction . . . . .	115
5.3	Proposed Tensor Face Models . . . . .	116
5.3.1	The Expression Space and the Apathy Mode . . . . .	117
5.3.2	Model 1: Basic Model . . . . .	124
5.3.3	Model 2: Subspace-aware Parameterization . . . . .	128
5.3.4	Model 3: Projection Pursuit in Expression Space . . . . .	132
5.3.5	Model 4: Four-Way Model including Expression Strength . . . . .	134
5.3.6	Overview of Presented Tensor Face Models . . . . .	141
5.4	Quality of Face Models . . . . .	142
<b>6</b>	<b>Experiments</b>	<b>146</b>
6.1	Facial Animation by Improved Synthesis Using Apathy . . . . .	146

---

6.2	3D Approximation, Person and Expression Transfer . . . . .	149
6.2.1	Evaluation . . . . .	153
6.3	Dense 3D Reconstruction from sparse 2D . . . . .	156
6.3.1	3D Reconstruction With Ground Truth . . . . .	156
6.3.2	3D Reconstruction Without Ground Truth . . . . .	170
6.3.3	Summary . . . . .	170
<b>7</b>	<b>Summary and Conclusions</b>	<b>174</b>
7.1	Future Work . . . . .	178
<b>Appendix</b>		<b>179</b>
A	3D Rotations and Computing Optimal Angles . . . . .	179
B	Normal Vector of 3D Points . . . . .	182
C	Parameterization of Lines along Principal Axis . . . . .	182
D	Apathy Estimation - How to Find the Closest Point . . . . .	183
E	Examples of Dense 3D Reconstruction of Bosphorus Database	185
<b>Literature</b>		<b>189</b>
<b>Index</b>		<b>199</b>

## Abstract

The creation of versatile 3D face models from limited training data has been a long-standing goal in facial animation. These models need to fully represent each individual face shape, including changes of facial expressions without the loss of individual facial features.

This difficulty is especially well-known in the movie industry, where even nowadays extensive manual work is necessary to achieve a natural representation of a human face with convincing expressive performance. This process is already challenging if sufficient high-quality 3D material of one person is available, but is considerably more difficult in the case of low-quality input data caused by limited hardware. In this work, different methods are presented to create 3D face models from databases of 3D face scans. The databases contain scans of various persons showing different expressions, a variety of points per 3D scan and different numbers of scans per person. Throughout this work objective quality criteria are carefully designed to quantify the quality requirements of each specific application.

In the first part of this work a preprocessing pipeline is presented, followed by a procedure to achieve dense meaningful correspondences between the 3D face scans. Then, based on the assumption of a shared motion pattern, a temporal alignment is estimated, which provides the same number of scans per person, such that facial motions are performed in synchrony. In this process, a robust descriptor for expression intensity is proposed, for which additional applications are presented, e.g. person-specific emotion cluster unveiling a variance in performance for each emotion between persons.

Since the resulting 3D faces are in full dense point-wise correspondence and their temporal facial movements are synchronized, they are aligned in space and time. Therefore, the processed 3D face scans can be arranged into a single data structure representing a 3D cube with axes corresponding to the number of 3D points, subjects and expressions, respectively. This data structure is referred to as *tensor* and outperforms the separation of different modes of individual shape and expression compared to traditional approaches based on 2D data structures, i.e. matrices. A 3D face model

is created from the data tensor by factorization into different modes. In contrast to former methods, the structures of the expression subspace are employed to derive reasonable constraints. In the expression subspace, it is observed that the six basic emotions (anger, disgust, fear, happiness, sadness, surprise) performed in different strengths, each form linear trajectories within the subspace. These six lines, each corresponding to one emotion, converge at a point which defines the natural origin of all expressions. It appears that this specific expression is not part of the database and that it differs from the neutral expression. Due to the fact that the database is based on posed instead of spontaneously performed expressions, the expression labeled as neutral differs from the fully relaxed face which would represent the expected case. Therefore the newly discovered origin is referred to as *apathetic*, which corresponds to an expression with fully relaxed facial muscles. It can be used for various applications: (1) to neutralize faces and replace the face with the original label neutral in the database, thereby improving the quality of the originally posed data without the need for new recordings, (2) to synthesize more convincing facial animations with an improved separation of distinct emotions, and (3) to adapt the statistical face model to render it more compact and robust, thereby enabling to perform stable expression and person transfer.

In this work four different face tensor models based on three databases are presented and compared for different applications: (1) 3D face synthesis, (2) 3D approximation, person and expression transfer, and (3) 3D reconstruction from 2D. The experiments show that dense 3D face reconstructions from sparse 2D landmarks based on the proposed models outperform those of the two state-of-the-art methods, although they employ more information from the original image.

**Keywords:** 3D face scans, 3D faces, nonrigid registration, correspondence estimation, expression intensity, tensor, factorization, statistical models, expression transfer, 3D reconstruction

## Kurzfassung

Schon lange arbeiten Menschen daran aus begrenzten Trainingsdaten vielseitig einsetzbare 3D-Gesichtsmodelle zu erzeugen. Diese sollen einerseits das Gesicht gut repräsentieren und andererseits glaubhafte Änderungen des Gesichtsausdrucks ermöglichen. Insbesondere in der Filmindustrie ist das Problem bekannt ein überzeugendes Ergebnis eines menschlichen Gesichts mit natürlichen Gesichtsausdrücken zu erzeugen und erfordert noch heute viel manuelle Arbeit. Trotz der Verfügbarkeit hochqualitativer Daten ist dies weiterhin eine Herausforderung, insbesondere dort, wo limitierte Hardware weniger gute Ergebnisse liefert. In dieser Arbeit werden Ansätze präsentiert, um verschiedene 3D-Gesichtsmodelle zu erzeugen. Diese basieren auf Datenbanken mit 3D-Scans von Gesichtern, die jeweils verschiedene Personen und Gesichtsausdrücke enthalten, sich jedoch in der Anzahl der Punkte und Scans pro Person unterscheiden. Zudem werden in jedem Teil dieser Arbeit objektive Qualitätskriterien definiert, die jeweils Eigenschaften speziell für die jeweilige Anwendungen quantifizieren.

Im ersten Teil dieser Arbeit wird eine zielgerichtete Vorverarbeitung der Daten präsentiert, gefolgt von einem Ansatz, um sinnvolle Korrespondenzen zwischen den 3D Gesicht-Scans zu schätzen. Unter der Annahme, dass es ein gemeinsames Bewegungsmuster in mehreren Aufnahmen von Gesichtsausdrücken gibt, wird eine zeitliche Ausrichtung geschätzt, um dieselbe Anzahl von Scans pro Person zu erhalten, so dass Gesichtsbewegungen synchron erfolgen. Dabei werden ein Deskriptor für die Intensität des Gesichtsausdrucks definiert und weitere Anwendungen präsentiert.

Die verarbeiteten 3D-Gesichts-Scans sind nun in Zeit und Raum sinnvoll geordnet. Daher können diese in eine Datenstruktur sortiert werden, die einem Würfel entspricht, bei dem die drei Dimensionen folgende Informationen enthalten: Anzahl der 3D-Punkte, der Identitäten und der Gesichtsausdrücke. Diese Datenstruktur wird als Tensor bezeichnet und erleichtert die Trennung von individueller Gesichtsform und Gesichtsausdruck im Vergleich zu traditionellen Methoden, welche die Daten in eine 2D-Datenstruktur, d.h. Matrizen, einordnen. Basierend auf dieser Datenstruktur wird ein Gesichtsaus-

modell mit einem Faktorisierungssatz erstellt. Anders als vorangegangenen Arbeiten, werden hier die gefundenen Strukturen in den Unterräumen verwendet um sinnvolle Nebenbedingungen zu definieren. In einem Unterraum finden sich Strukturen, in denen die sechs prototypischen Emotionen (Ärger, Ekel, Angst, Glück, Trauer, Überraschung), die in unterschiedlicher Stärke ausgeführt wurden, jeweils eine Gerade in dem Unterraum bilden. Diese sechs Geraden, jeweils zugehörig zu einer Emotion, treffen sich in einem gemeinsamen Punkt, welcher dem natürlichen Ursprung aller Gesichtsausdrücke entspricht. Es stellt sich heraus, dass dieser Gesichtsausdruck nicht Teil der Datenbank ist und sich von dem als neutral gekennzeichneten Gesichtsausdruck unterscheidet. Basierend darauf, dass die Datenbank aus Aufnahmen von gestellten und nicht aus spontan ausgeführten Gesichtsausdrücken besteht, folgt, dass Gesichter, die als neutral gekennzeichnet sind, individuelle Merkmale enthalten, die nicht immer dem erwarteten neutralen Gesichtsausdruck entsprechen, nämlich einem entspannten Gesichtsausdruck. Daher wird der gefundene neue Ursprung als *apathischer* Gesichtsausdruck bezeichnet, da er einem Ausdruck entspricht, bei dem alle Gesichtsmuskeln vollständig entspannt sind. Dieser kann für verschiedene Anwendungen genutzt werden: (1) Nachträgliche *Neutralisierung* des Gesichtsausdrucks, um die originalen Daten mit dem Label *neutral* zu ersetzen, wobei die Qualität der Daten verbessert werden kann ohne neue Aufnahmen zu benötigen. (2) Synthese überzeugender Gesichtsanimationen, welche die Vermischung verschiedener Emotionen verhindert. (3) Darüber hinaus wird demonstriert wie statistische Gesichtsmodelle robuster gemacht werden können, so dass der Austausch von Gesichtsausdrücken und Personen stabilisiert wird.

In dieser Arbeit werden vier verschiedene tensorbasierte Gesichtsmodelle, erstellt aus drei Datenbanken, vorgestellt und anhand verschiedener Anwendungen verglichen: (1) Synthese von 3D-Gesichtern, (2) 3D-Approximation, Transfer von Gesichtsausdrücken und Identität, und (3) 3D-Rekonstruktion aus 2D-Input. Es wird gezeigt, dass die präsentierten 3D Rekonstruktionen basierend auf wenigen 2D-Landmarken, die durch die vorgestellten Modelle erzeugt wurden, bessere Ergebnisse liefern als zwei State-of-the-Art-Methoden, obwohl diese mehr Informationen aus den Bildern verwenden.

**Stichworte:** 3D-Gesichts-Scan, 3D-Gesichter, nicht-rigide Registrierung, Korrespondenzschätzung, Intensität von Gesichtsausdrücken, Tensor, statistische Modelle, Transfer von Gesichtsausdrücken, 3D-Rekonstruktion

# Abbreviations and Nomenclature

## Abbreviations

<b>3DMM</b>	3D Morphable Model
<b>AU</b>	Facial Action Unit
<b>CPD</b>	Coherent Point Drift
<b>DTW</b>	Dynamic Time Warping
<b>ECPD</b>	Extended Coherent Point Drift
<b>EM</b>	Expectation Maximization
<b>FACS</b>	Facial Action Coding System
<b>ffp</b>	facial feature point
<b>FAP</b>	Facial Animation Parameter
<b>FAU</b>	Facial Action Unit
<b>GCTW</b>	Generalized Canonical Time Warping
<b>GMM</b>	Gaussian Mixture Model
<b>HOSVD</b>	High-Order Singular Value Decomposition
<b>ICA</b>	Independent Component Analysis
<b>ICP</b>	Iterative Closest Point
<b>LFAU</b>	Lower Facial Action Unit
<b>MAP</b>	Maximum A Posteriori
<b>ML</b>	Maximum Likelihood



---

<b>PCA</b>	Principal Component Analysis
<b>pdf</b>	probability density function
<b>SVD</b>	Singular Value Decomposition
<b>UFAU</b>	Upper Facial Action Unit

## Nomenclature

$\mathbb{C}$	complex numbers
$\mathbb{R}$	real numbers
$\mathbb{N}$	natural numbers
$\mathbb{E}(y)$	expectation value of random variable $y$
$s$	lower case italic letters define scalar values: $s \in \mathbb{R}$ .
$\mathbf{v}$	lower case bold letters define column vectors: $\mathbf{v} \in \mathbb{R}^{N \times 1}$
$\mathbf{M}$	upper case bold letters define matrices: $\mathbf{M} \in \mathbb{R}^{M \times N}$
$\mathcal{T}$	upper case slanted letters define sets or tensors, e.g. a 3D tensor $\mathcal{T} \in \mathbb{R}^{L \times M \times N}$
$v_i$	$i$ th element of vector $\mathbf{v}$
$m_{ij}$	element of $i$ th row and $j$ th column of matrix $\mathbf{M}$
$\mathbf{M}(i, :)$	$i$ th row of matrix $\mathbf{M}$
$\mathbf{M}(:, j)$	$j$ th column of matrix $\mathbf{M}$
$\mathbf{M}(:)$	concatenate columns of matrix $\mathbf{M}$ to vector
$\mathbf{v}^T, \mathbf{M}^T$	transposed vector $\mathbf{v}$ and matrix $\mathbf{M}$
$\mathbf{M}^{-1}$	inverse matrix
$\mathbf{1}_n$	vector of $n$ ones: $\mathbf{1}_n \in \mathbb{R}^n$
$\mathbf{I}_n$	unity matrix $\mathbf{I}_n \in \mathbb{R}^{n \times n}$
$ a $	absolute value of scalar value
$\ \mathbf{a}\ $	length of vector $\mathbf{a}$

---

$ \mathbf{A} $	determinant of matrix $\mathbf{A}$
$\ \mathbf{v}\ _0$	0-norm of vector $\mathbf{v}$
$\ \mathbf{v}\ _1$	1-norm of vector $\mathbf{v}$
$\ \mathbf{v}\ _2$	Euclidean norm of vector $\mathbf{v}$
$\ \mathbf{M}\ _F$	Frobenius norm of matrix $\mathbf{M}$
$\mathbb{1}(A)$	the indicator function evaluates to either 0 if $A$ is false and 1 if $A$ is true
$\delta_{ij}$	Kronecker delta is 0 if $i \neq j$ or 1 if $i = j$
$\otimes$	Kronecker product of vectors or matrices
$\mathcal{N}(\mu, \sigma^2)$	Gaussian distribution with expectation value $\mu$ and variance $\sigma^2$
$\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$	multivariate Gaussian distribution with vector-valued expectation value $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$
$\mathcal{T} \times_k \mathbf{M}$	mode- $k$ tensor product between tensor $\mathcal{T}$ and matrix $\mathbf{M}$
$\ln(\cdot)$	logarithm to base $e$
$\text{diag}(\mathbf{v})$	returns diagonal matrix with input vector on the diagonal
$\text{diag}(\mathbf{M})$	extracts diagonal from matrix $\mathbf{M}$ as vector
$\text{tr}(\mathbf{M})$	trace of matrix $\mathbf{M}$