

Learning face reconstruction in the wild, with 3DMM

SEMINAR ON 9/23

YUDA QIU

Task

◆ Input: in-the-wild face image/ video/ image set

◆ Output: face geometry & face texture

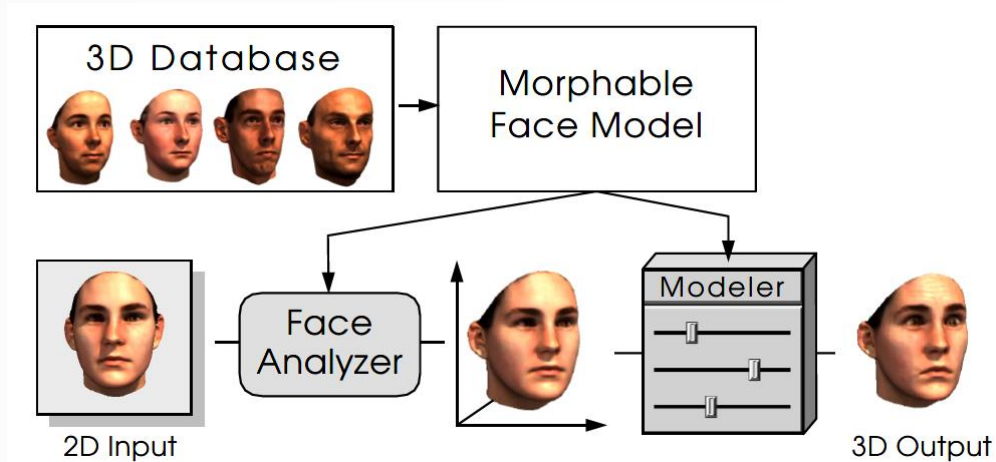
□ Lack of dense ground truth

□ Ill-posed



Model-based self-supervision (model-based decoder)

3DMM representation



$$S_{model} = \bar{S} + \sum_{i=1}^{m-1} \alpha_i s_i, \quad T_{model} = \bar{T} + \sum_{i=1}^{m-1} \beta_i t_i, \quad (1)$$

$\vec{\alpha}, \vec{\beta} \in \mathbb{R}^{m-1}$. The probability for coefficients $\vec{\alpha}$ is given by

$$p(\vec{\alpha}) \sim \exp\left[-\frac{1}{2} \sum_{i=1}^{m-1} (\alpha_i / \sigma_i)^2\right], \quad (2)$$

- PCA to build statistics model
- Analysis-by-synthesis manner

$$\mathbf{x} = (\underbrace{\alpha, \delta, \beta}_{\text{face}}, \underbrace{\mathbf{T}, \mathbf{t}, \gamma}_{\text{scene}}).$$

A Morphable Model For The Synthesis Of 3D Faces

Volker Blanz

Thomas Vetter

Max-Planck-Institut für biologische Kybernetik,
Tübingen, Germany*

Paper lists

- (CVPR 2017) Regressing robust and discriminative 3d morphable models with a very deep neural network
- (ICCV 2017) MoFA: Model-based Deep Convolutional Face Autoencoder for Unsupervised Monocular Reconstruction
- (CVPR 2018) Unsupervised Training for 3D Morphable Model Regression.
- (CVPR 2018) Self-supervised multi-level face model learning for monocular reconstruction at over 250 HZ
- (CVPRW 2019, Best award) Accurate 3D Face Reconstruction with Weakly-Supervised Learning: From Single Image to Image Set.**
- (CVPR 2019) GANFIT: Generative Adversarial Network Fitting for High Fidelity 3D Face Reconstruction.
- (CVPR 2019, oral) FML: Face Model Learning from Videos.

Regressing Robust and Discriminative 3D Morphable Models with a very Deep Neural Network

Anh Tuấn Trần¹, Tal Hassner^{2,3}, Iacopo Masi¹, and Gérard Medioni¹

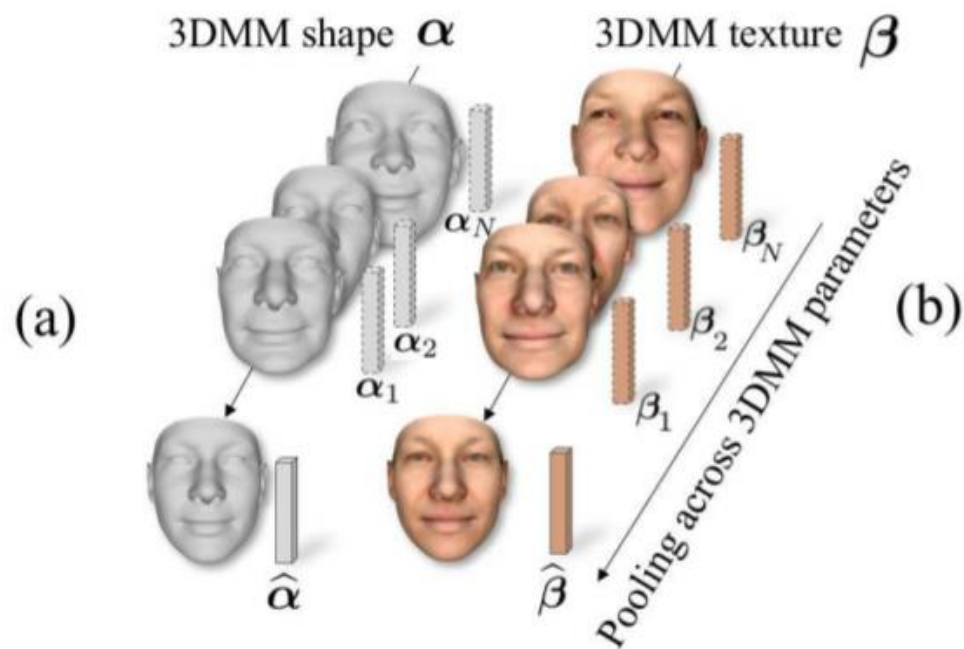
¹ Institute for Robotics and Intelligent Systems, USC, CA, USA

² Information Sciences Institute, USC, CA, USA

³ The Open University of Israel, Israel



Faces in the wild
to train the system



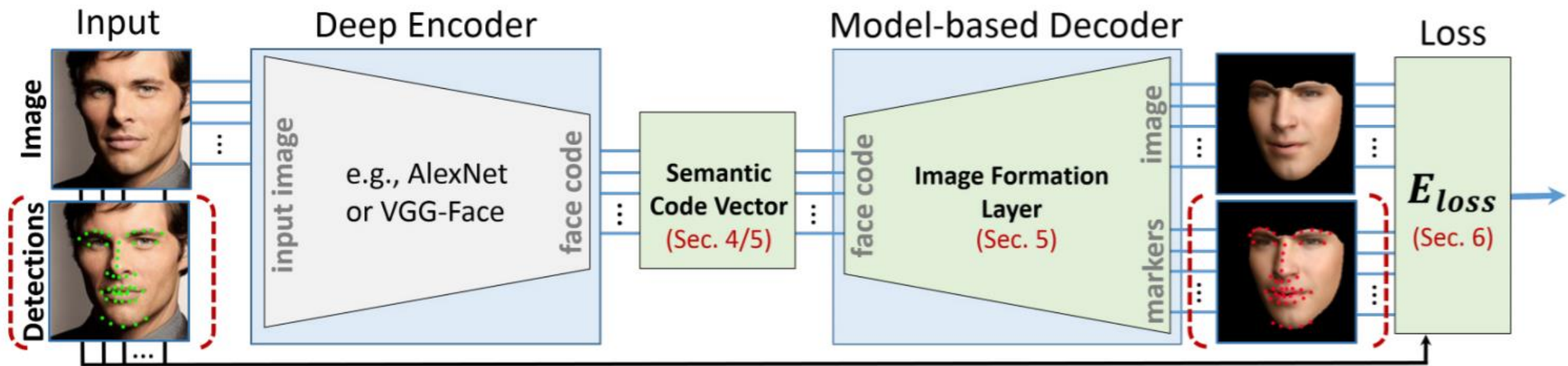
$$\hat{\gamma} = \sum_{i=1}^N w_i \cdot \gamma_i \quad \text{and} \quad \sum_{i=1}^N w_i = 1, \quad (2)$$

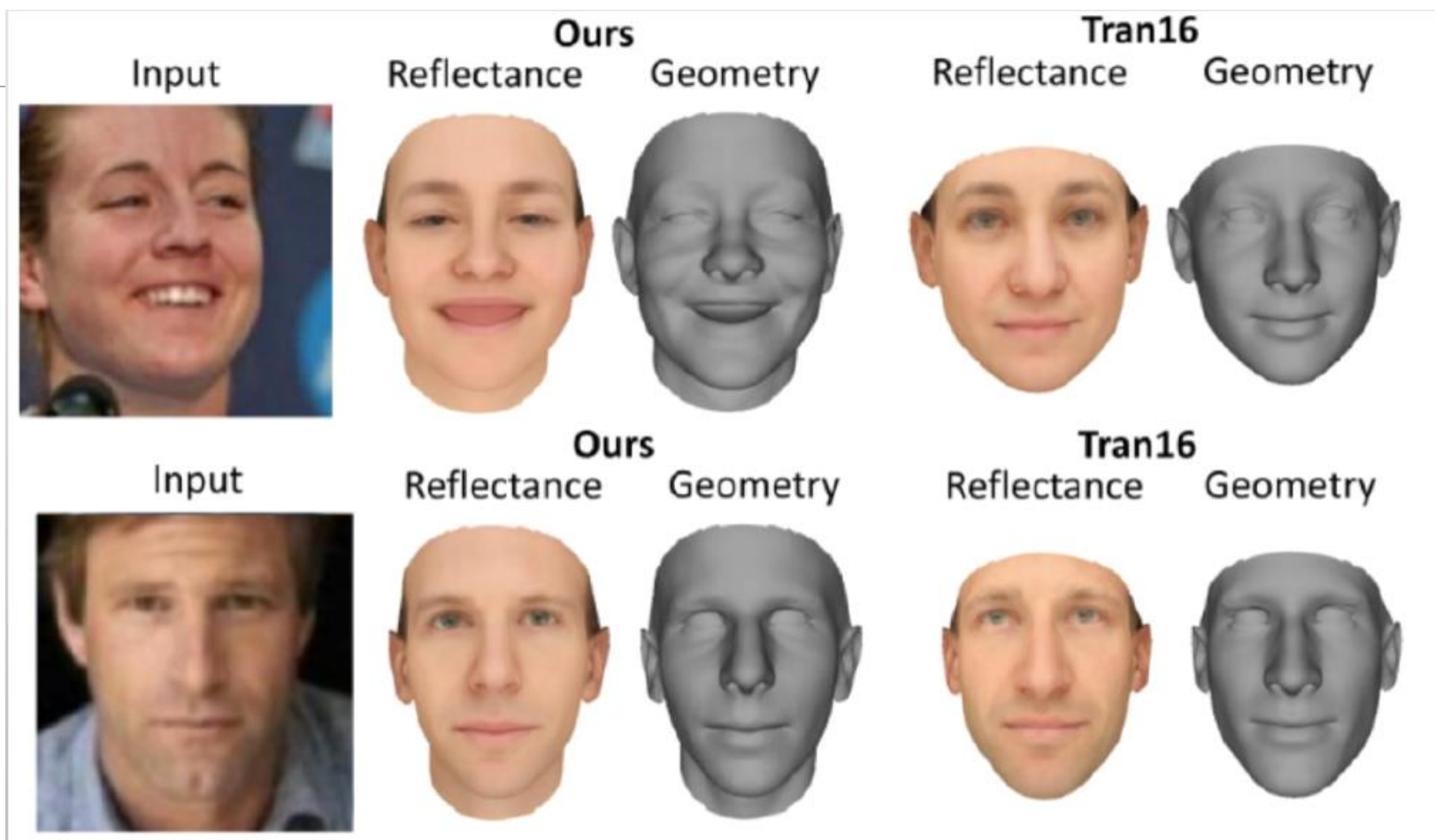
where w_i are normalized per-image confidences provided by the CLNF facial landmark detector.

MoFA: Model-based Deep Convolutional Face Autoencoder for Unsupervised Monocular Reconstruction

Ayush Tewari¹ Michael Zollhöfer¹ Hyeonwoo Kim¹ Pablo Garrido¹
Florian Bernard^{1,2} Patrick Pérez³ Christian Theobalt¹

¹Max-Planck-Institute for Informatics ² LCSB, University of Luxembourg ³Technicolor



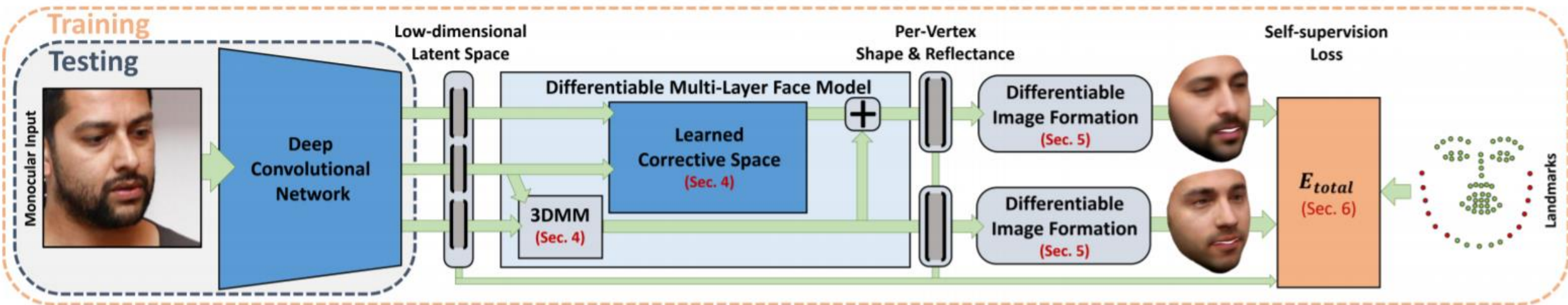


Self-supervised Multi-level Face Model Learning for Monocular Reconstruction at over 250 Hz

CVPR 2018 (Oral)

A. Tewari^{1,2} M. Zollhöfer^{1,2,3} F. Bernard^{1,2} P. Garrido^{1,2} H. Kim^{1,2} P. Perez⁴ C.Theobalt^{1,2}

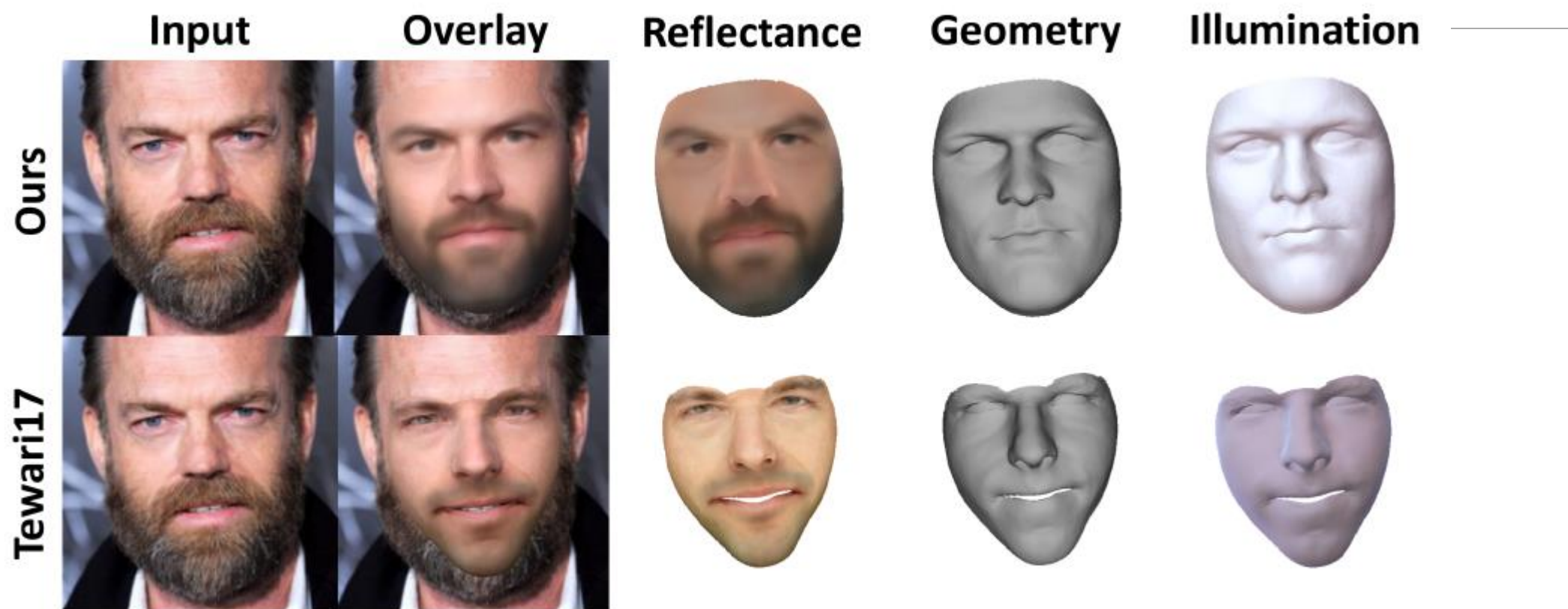
¹MPI Informatics ²Saarland Informatics Campus ³Stanford University ⁴Technicolor



$$\mathbf{v}^f(\mathbf{x}_g) = \mathbf{v}^b(\boldsymbol{\alpha}) + \mathcal{F}_g(\boldsymbol{\delta}_g | \Theta_g) \in \mathbb{R}^{3N} \text{ (geometry),}$$

$$\mathbf{r}^f(\mathbf{x}_r) = \mathbf{r}^b(\boldsymbol{\beta}) + \mathcal{F}_r(\boldsymbol{\delta}_r | \Theta_r) \in \mathbb{R}^{3N} \text{ (reflectance),}$$

$$\mathbf{x} = (\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\delta}_g, \boldsymbol{\delta}_r, \mathbf{R}, \mathbf{t}, \boldsymbol{\gamma}, \Theta_g, \Theta_r) \in \mathbb{R}^{257+2C+|\Theta_g|+|\Theta_r|}.$$



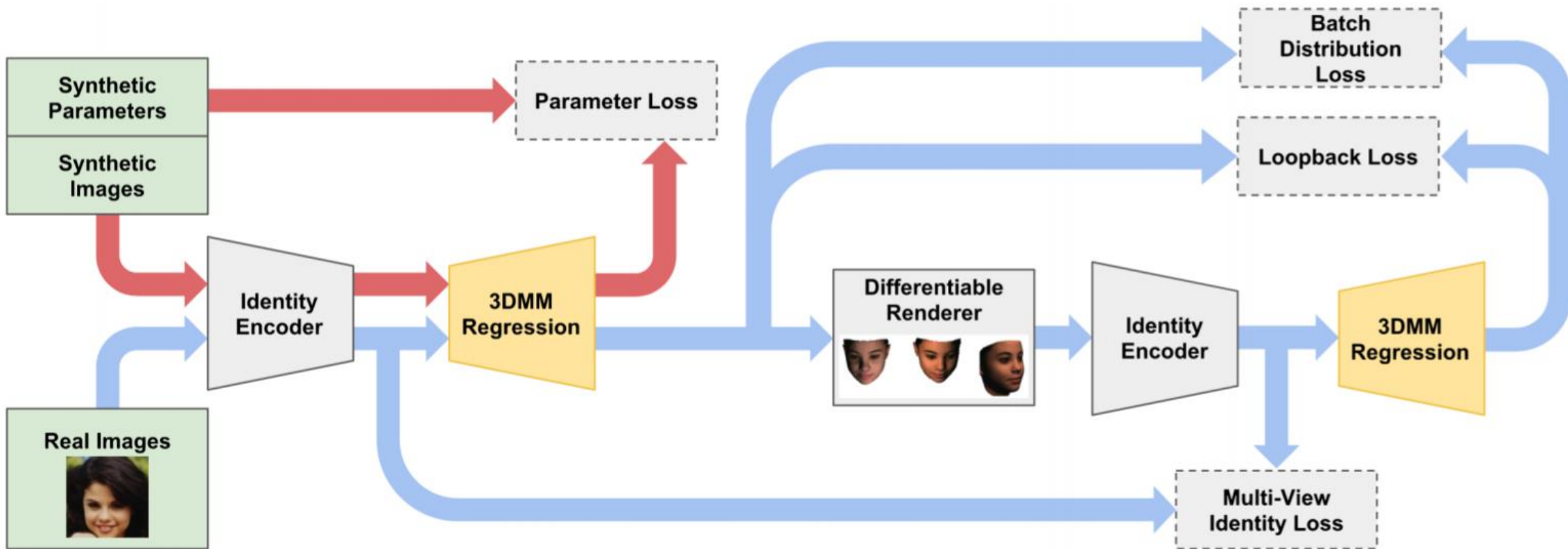
Unsupervised Training for 3D Morphable Model Regression

Kyle Genova^{1,2} Forrester Cole² Aaron Maschinot² Aaron Sarna² Daniel Vlasic² William T. Freeman^{2,3}

¹Princeton University

²Google Research

³MIT CSAIL



$$L = L_{param} + L_{id} + \omega_{batch} L_{batch} + \omega_{loop} L_{loop}$$

Input



No Batch Dist.



Full Model



No Loopback



No Multi-View



Input



Ours



Tran[29]



MoFA[27]





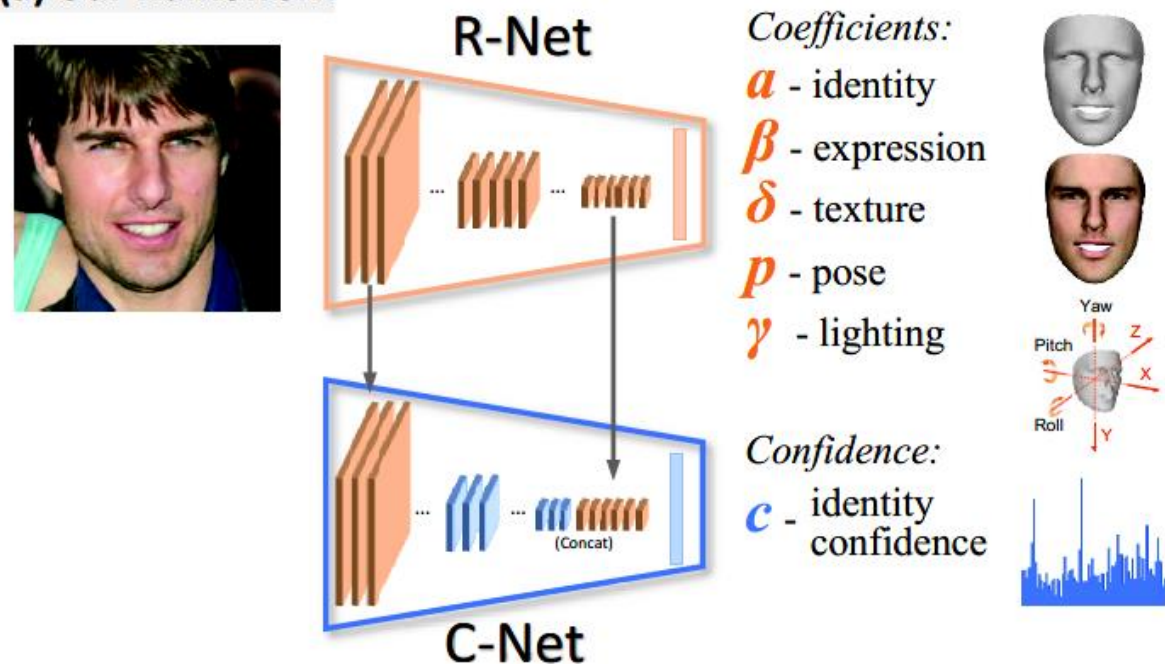
Accurate 3D Face Reconstruction with Weakly-Supervised Learning: From Single Image to Image Set

Yu Deng^{*1,2} Jiaolong Yang² Sicheng Xu^{3,2} Dong Chen² Yunde Jia³ Xin Tong²

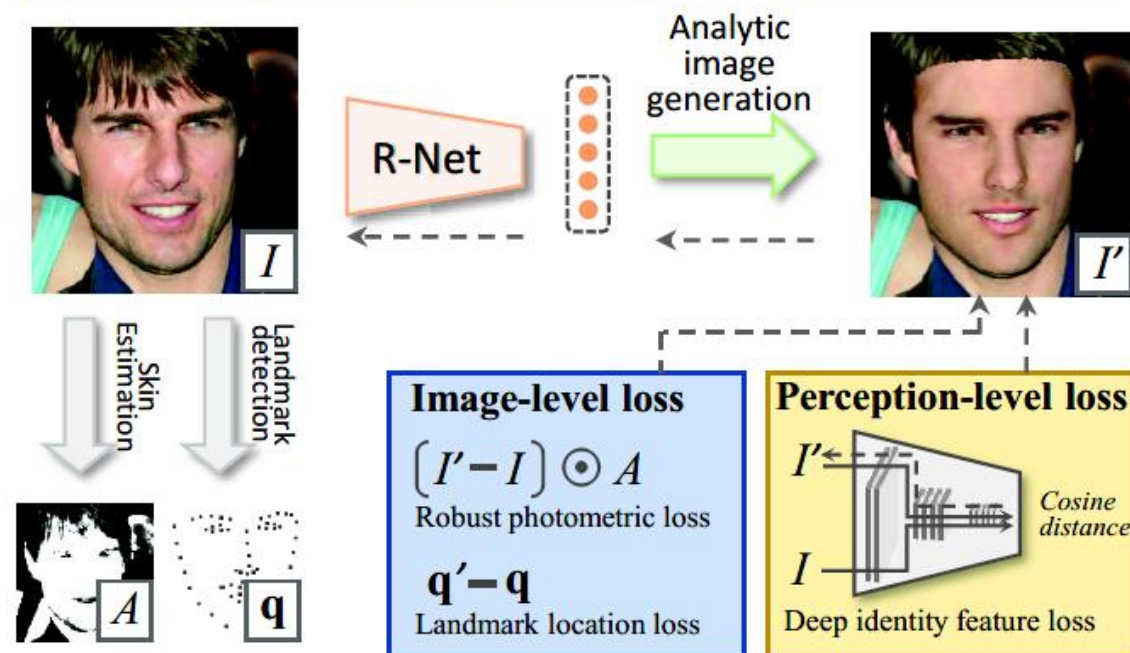
¹Tsinghua University ²Microsoft Research Asia ³Beijing Institute of Technology

{v-denyu, jiaoyan, doch, xtong}@microsoft.com, {xusicheng, jiaayunde}@bit.edu.cn

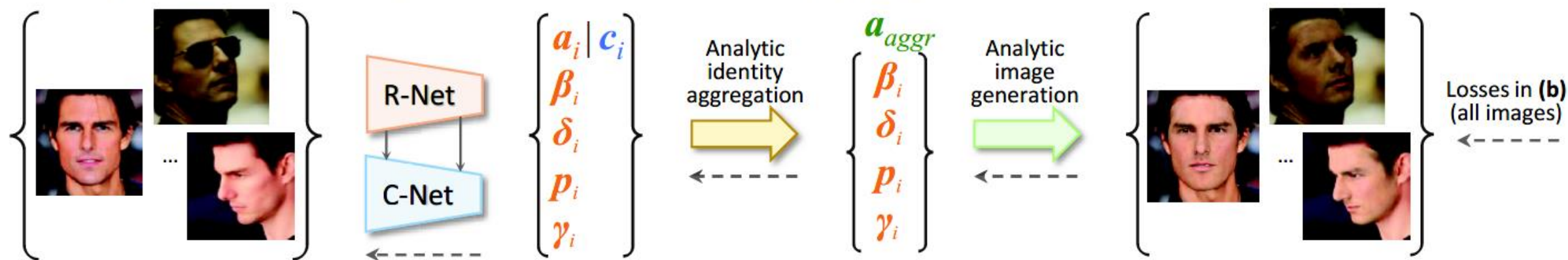
(a) Our framework



(b) Training pipeline for single image 3D face reconstruction



(c) Training pipeline for multi-image 3D face reconstruction with shape aggregation





Input *Genova et al.* Ours

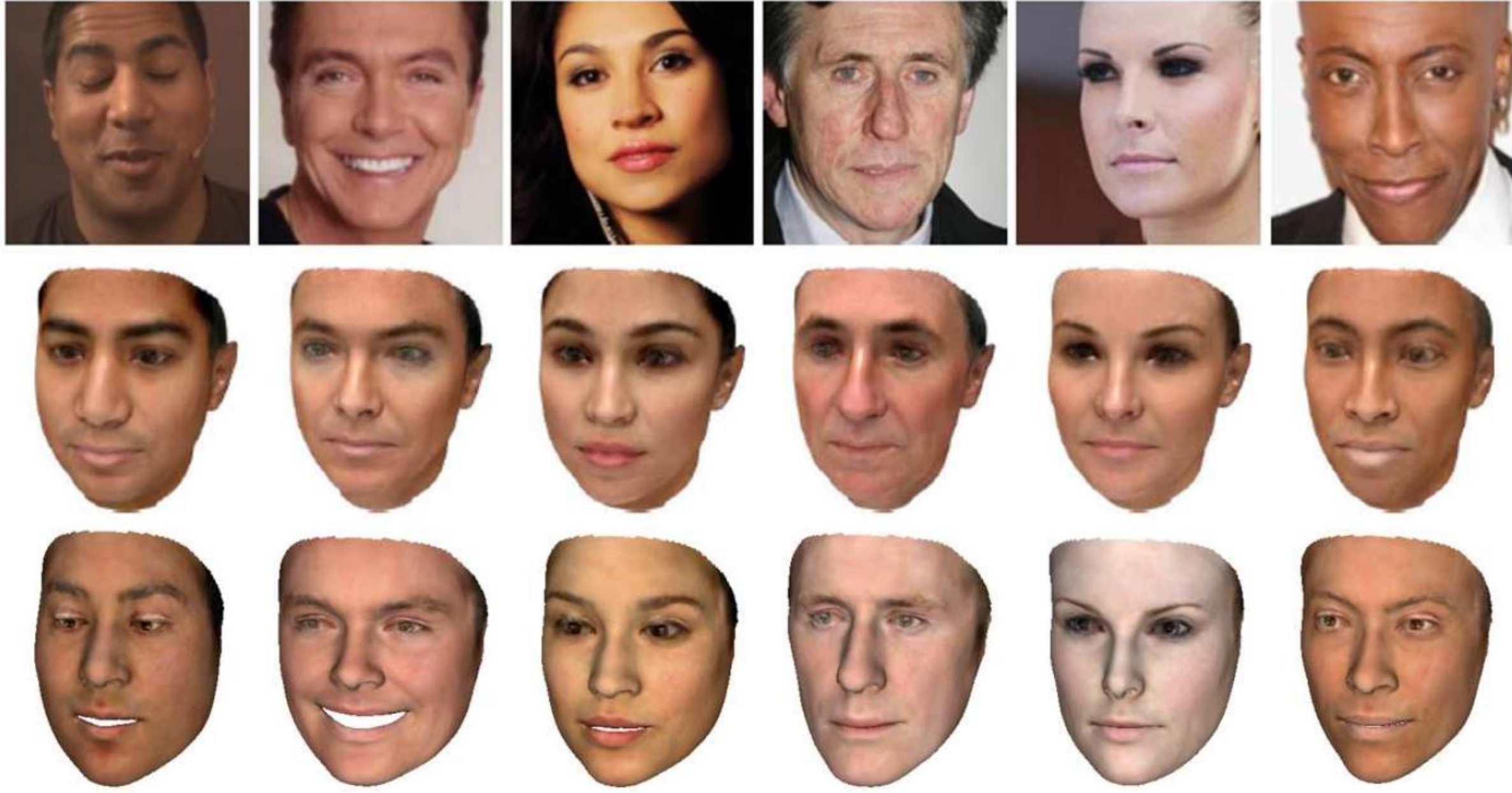


Table 4. Multi-image reconstruction errors on MICC rendered images with different aggregation strategies (see text for details).

Shape error mean	1.97 ± 0.70
Shape averaging	1.78 ± 0.59
Our S1: Global Aggr. with c^j	1.71 ± 0.56
Our S2: Global Aggr. with $\sum_i c_i^j$	1.70 ± 0.55
Our S3: Max Conf. $j = \arg \max_j \sum_i c_i^j$	1.71 ± 0.50
Our S4: Elementwise Aggr. with c^j	1.67 ± 0.54

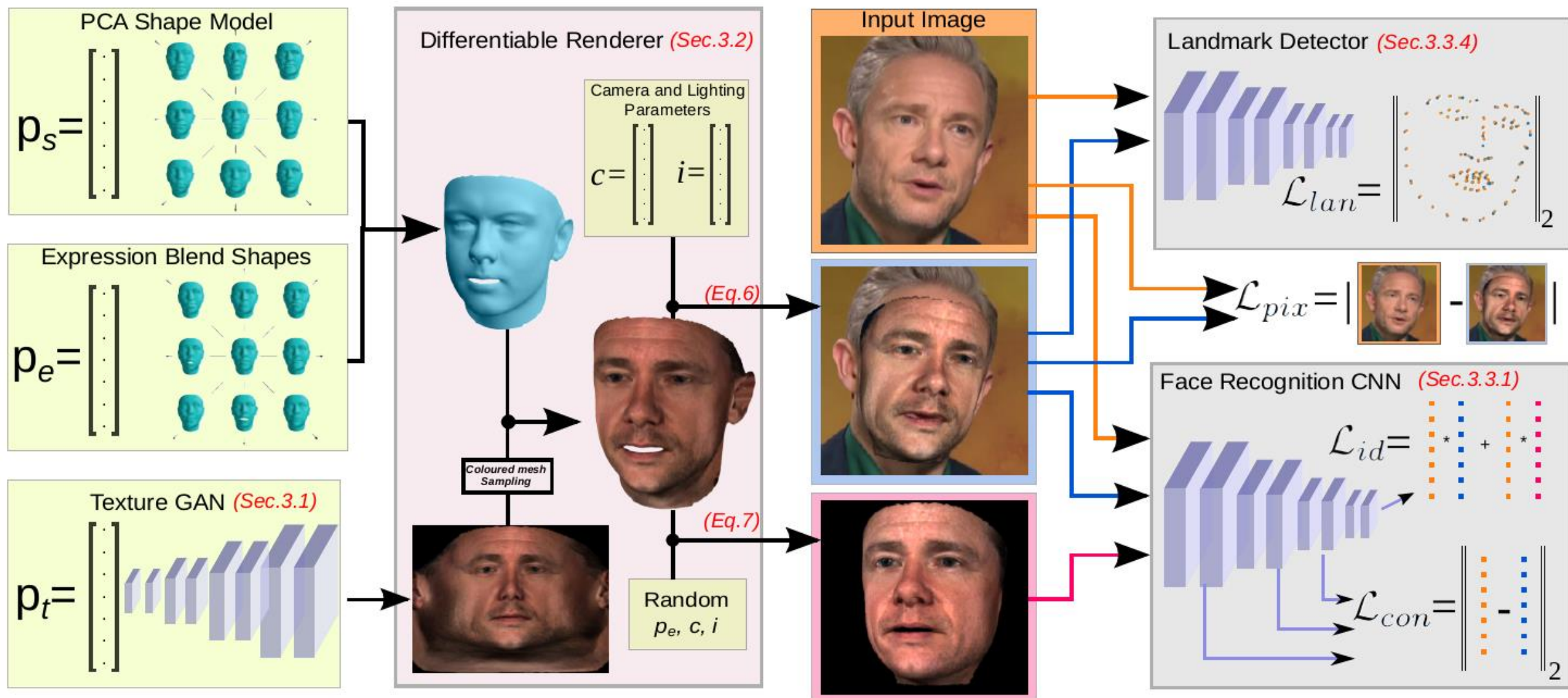
GANFIT: Generative Adversarial Network Fitting for High Fidelity 3D Face Reconstruction

Baris Gecer^{1,2}, Stylianos Ploumpis^{1,2}, Irene Kotsia³, and Stefanos Zafeiriou^{1,2}

¹Imperial College London

²FaceSoft.io

³University of Middlesex



Input Images



Ours



Genova
[16]



A.T. Tran et al.
[39]




Tewari et al.
[36]

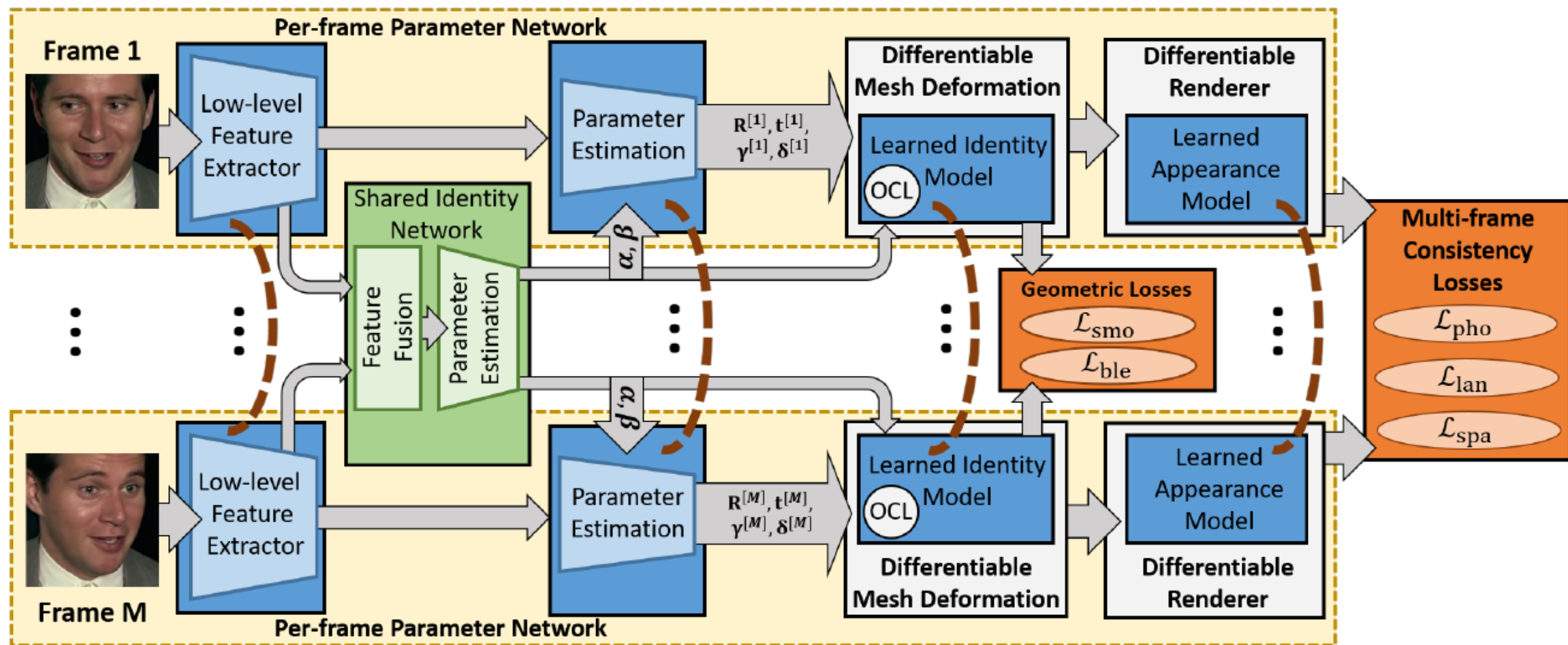


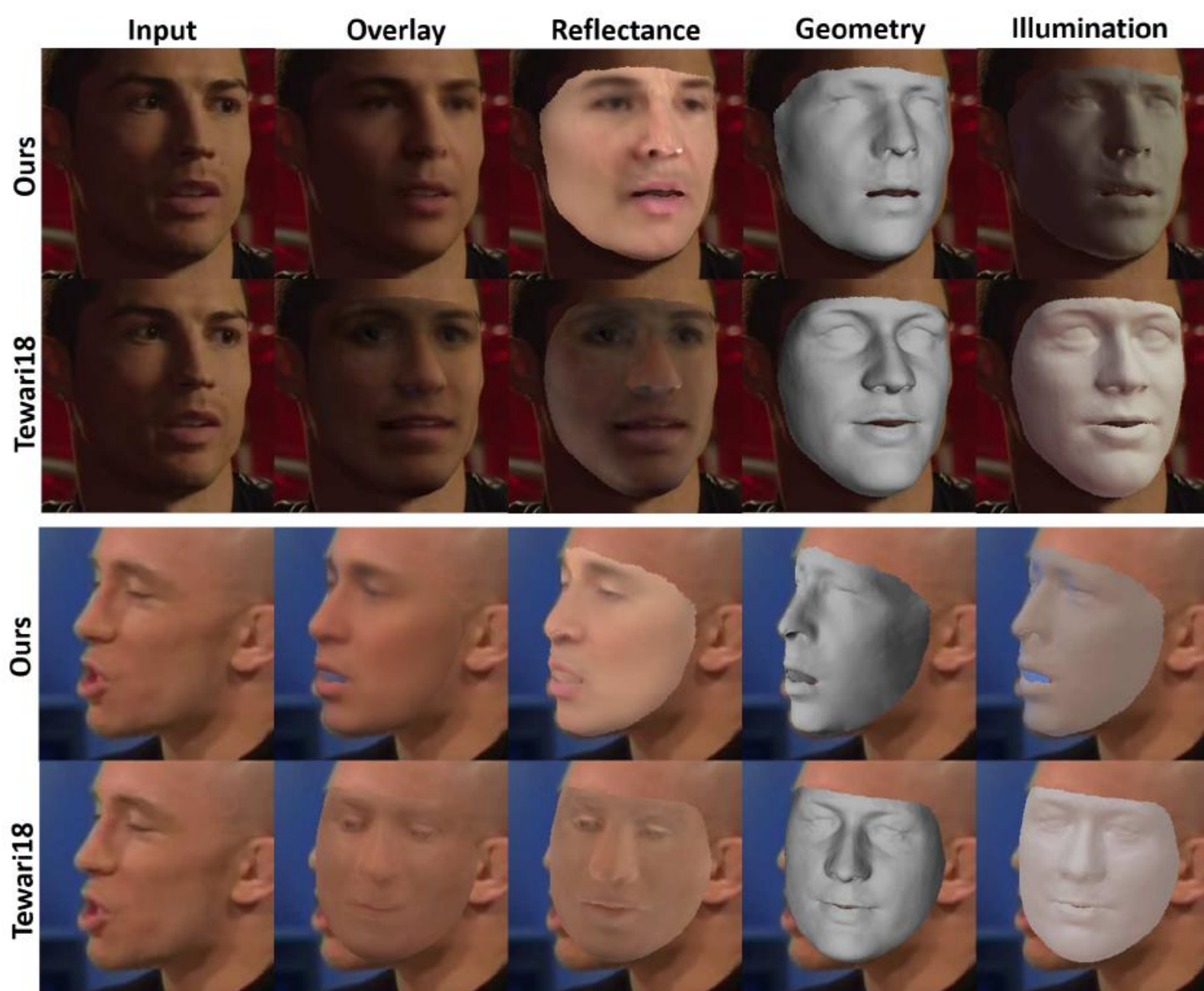
FML: Face Model Learning from Videos

Ayush Tewari¹ Florian Bernard¹ Pablo Garrido² Gaurav Bharaj² Mohamed Elgharib¹
Hans-Peter Seidel¹ Patrick Pérez³ Michael Zollhöfer⁴ Christian Theobalt¹

¹MPI Informatics, Saarland Informatics Campus ²Technicolor ³Valeo.ai ⁴Stanford University







Q&A
