

Learning to Generate Facial Depth Maps

Stefano Pini, Filippo Grazioli, Guido Borghi, Roberto Vezzani, Rita Cucchiara

University of Modena and Reggio Emilia, Italy

name.surname@unimore.it

Abstract



In this paper, an adversarial architecture for **facial depth map estimation** from monocular intensity images is presented.



By following an *image-to-image* approach, we combine the advantages of **supervised learning** and **adversarial training**, proposing a **conditional Generative Adversarial Network** that effectively learns to translate intensity face images into the corresponding depth maps.

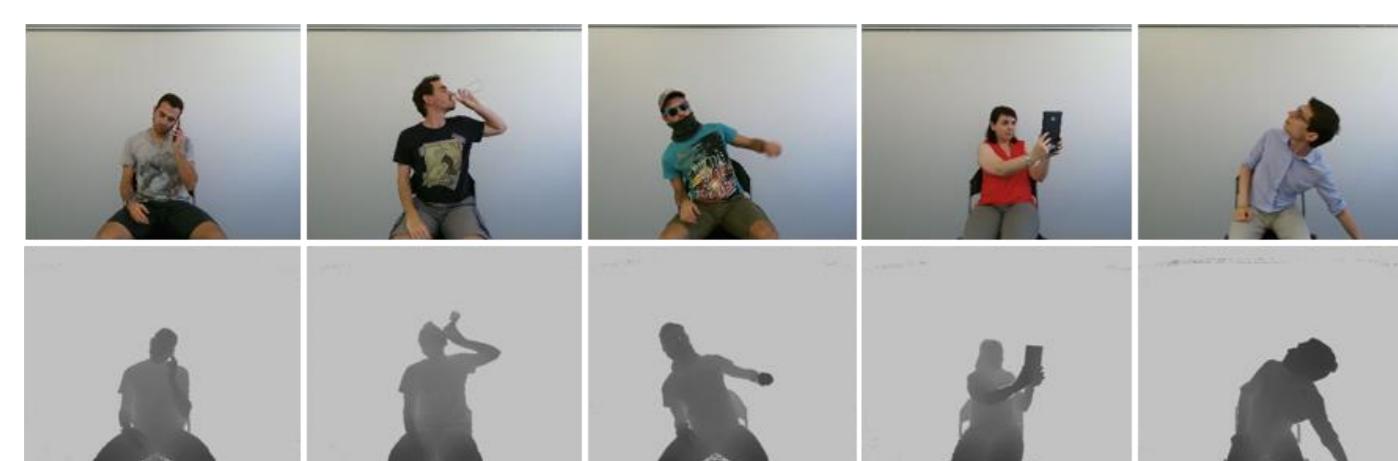
Adversarial Network that effectively learns to translate intensity face images into the corresponding depth maps. Furthermore, we show that the model is capable of predicting distinctive facial details by testing the generated depth maps through a deep model trained on authentic depth maps for the **face verification** task.

Datasets

Pandora dataset

This dataset has been presented in [1] and it has the following features:

- 250k frames from 110 sequences
- 22 subjects (10 males and 12 females).
- Both **depth** and **RGB** frames and **skeleton** annotations frame by frame.
- Acquired by **Microsoft Kinect One v2**
- Original tasks: **head** and **shoulder** pose **estimation**. Subjects can vary their head appearance wearing **garments** and **objects** like smartphones, tablets and plastic bottles that can generate **head and body occlusions**.



Biwi Kinect Head Pose

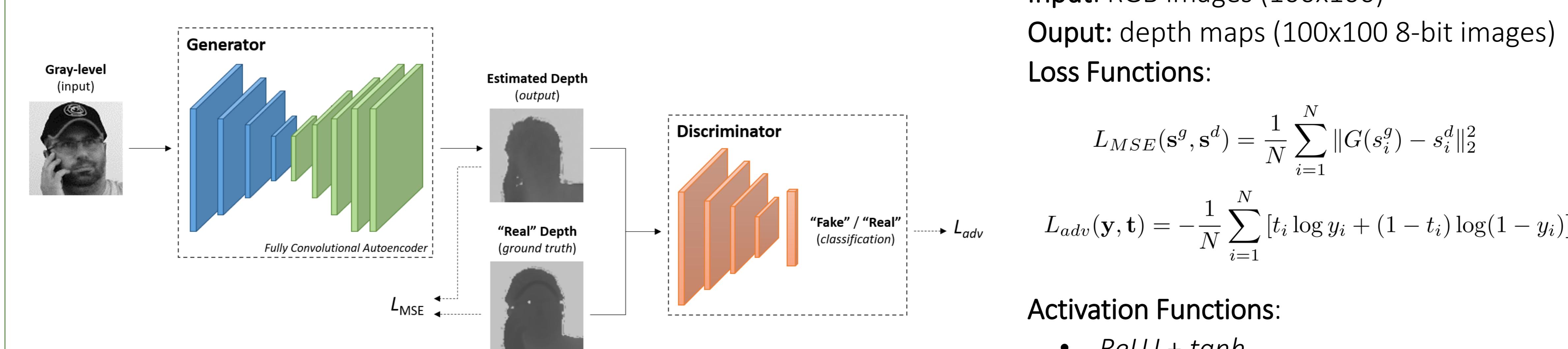
This dataset [2] contains:

- 15k frames
- 20 subjects (14 males and 6 females)
- Both **depth** and **RGB** frames
- Acquired by **Microsoft Kinect v1**
- Original task: head pose estimation only



Proposed Method

We propose a **Conditional Generative Adversarial Network** that effectively learns to translate intensity face images into the corresponding depth maps.



Dynamic Face Crop → a new window of size (w_H, h_H) for each frame

$f_{x,y}$: focal lengths
 $R_{x,y}$: face size

D : estimated distance

$$w_H = \frac{f_x \cdot R_x}{D} \quad h_H = \frac{f_y \cdot R_y}{D}$$

Input: RGB images (100x100)

Output: depth maps (100x100 8-bit images)

Loss Functions:

$$L_{MSE}(s^g, s^d) = \frac{1}{N} \sum_{i=1}^N \|G(s_i^g) - s_i^d\|_2^2$$

$$L_{adv}(\mathbf{y}, \mathbf{t}) = -\frac{1}{N} \sum_{i=1}^N [t_i \log y_i + (1 - t_i) \log(1 - y_i)]$$

Activation Functions:

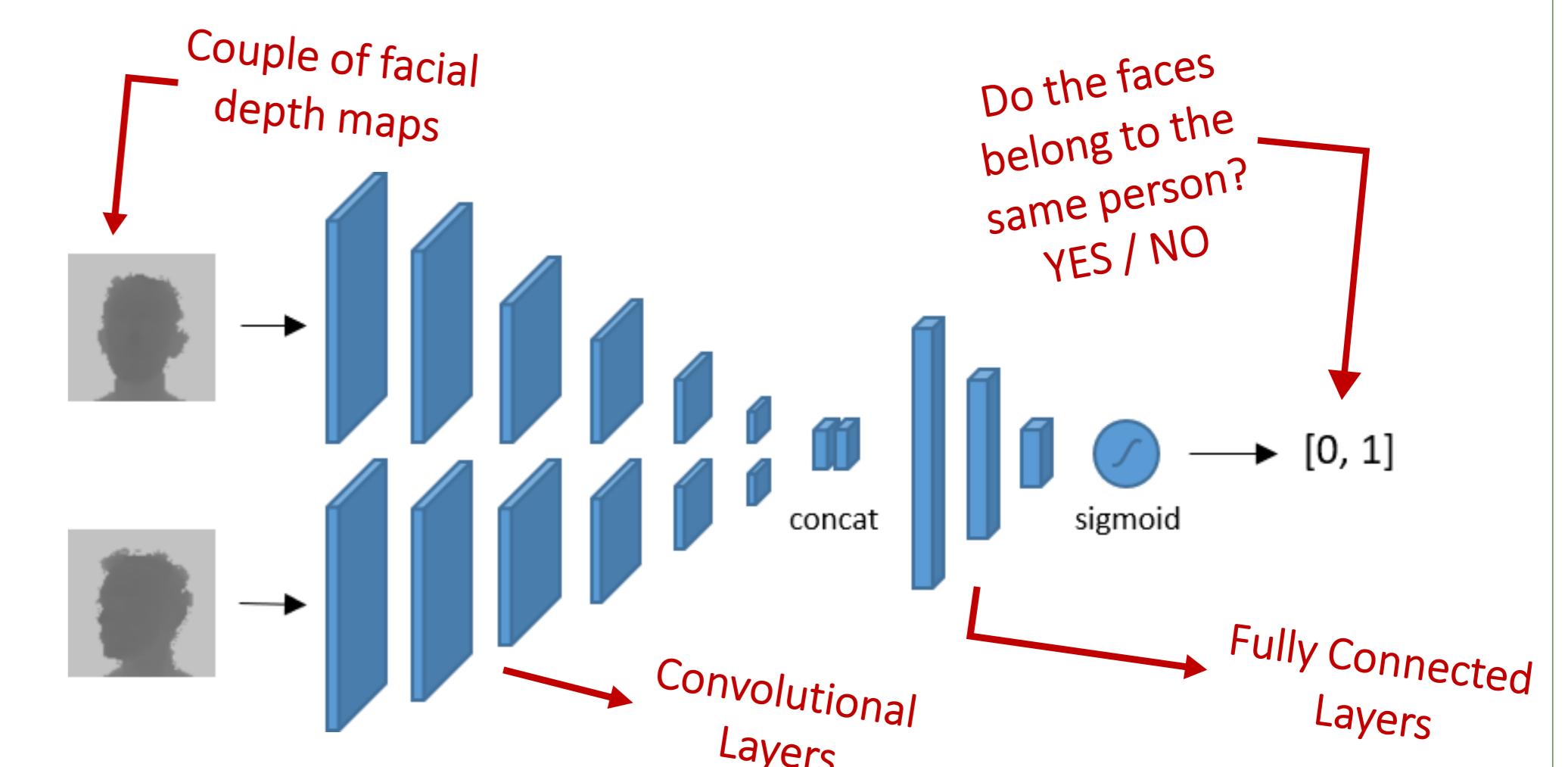
- *ReLU + tanh*

Network Architecture:

- Generator: *Fully Convolutional Network*
- Discriminator: *conv + fc* network

Face Verification test

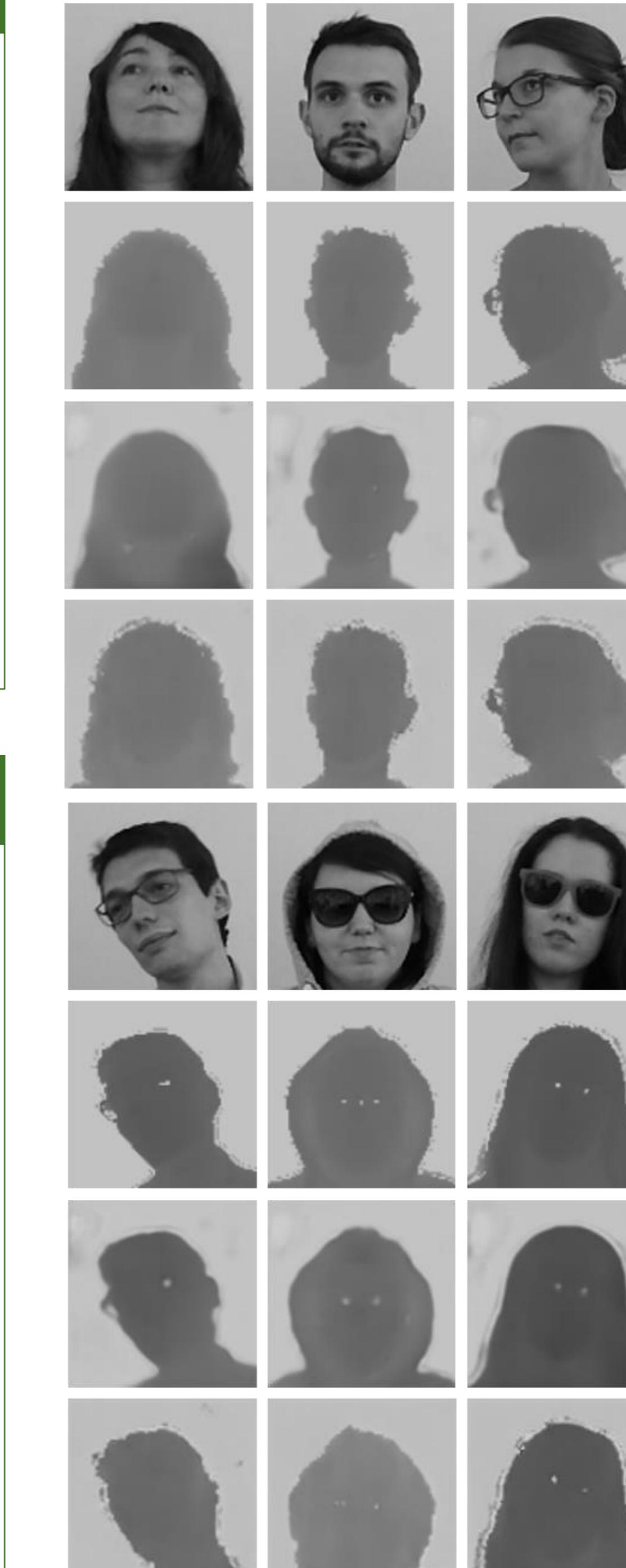
We introduce a **Face Verification Siamese** architecture [3], trained on the original face depth images, to check if the generated images maintain the **facial distinctive features** of the original subjects, not only when visually inspected by humans, but also when processed by deep convolutional network.



Limitations

- The detail accuracy of the proposed model is quite good, compared with the tested competitors;
- δ -metrics commonly used ($\delta < 1.25$, $\delta < 1.25^2$, $\delta < 1.25^3$) [6, 7], are effective to check the overall quality of depth maps generated from landscapes or wide-angle scenes, but the threshold value is too high to take fine details into account;
- We introduce a new set of δ -metrics ($\delta < 1.25^{\frac{1}{2}}$, $\delta < 1.25^{\frac{1}{3}}$, $\delta < 1.25^{\frac{1}{4}}$) with harder thresholds

Detail Accuracy is still an open problem with Conditional GANs!



Sample outputs of the proposed method

[RGB, Depth, AE and cGAN]



Check the Code
and Models



Download
Pandora Dataset



Acknowledgements

This work has been carried out within the project "FAR2015 - Monitoring the car driver's attention with multisensory systems, computer vision and machine learning" funded by the University of Modena and Reggio Emilia. We also acknowledge the CINECA award under the ISCRA initiative, for the availability of high performance computing resources and support

References

- [1] G. Borghi, M. Venturelli, R. Vezzani, and R. Cucchiara "Poseidon: Face-from-depth for driver pose estimation" (CVPR 2017)
- [2] G. Fanelli, M. Dantone, J. Gall, A. Fossati, L. Van Gool "Random Forests for Real Time 3D Face Analysis" (IJCV 2013)
- [3] G. Borghi, S. Pini, F. Grazioli, R. Vezzani, R. Cucchiara "Face Verification from Depth using Privileged Information" (BMVC 2018)
- [4] P. Isola, J. Zhu, T. Zhou, and A. A. Efros. "Image-to-image translation with conditional adversarial networks." (CVPR 2017)
- [5] J. Cui, H. Zhang, H. Han, S. Shan, and X. Chen. "Improving 2d face recognition via discriminative face depth estimation" (ICB 2018)
- [6] A. Atapour-Abarghouei and T. Breckon. "Real-time monocular depth estimation using synthetic data with domain adaptation via image style transfer" (CVPR 2018)
- [7] D. Xu, W. Wang, H. Tang, H. Liu, N. Sebe, and E. Ricci. "Structured attention guided convolutional neural fields for monocular depth estimation" (CVPR 2018)