

This is a pre print version of the following article:

Learning to Generate Facial Depth Maps / Pini, Stefano; Grazioli, Filippo; Borghi, Guido; Vezzani, Roberto; Cucchiara, Rita. - (2018). (Intervento presentato al convegno The 6th international conference on 3D Vision (3DV) tenutosi a Verona nel 5-8 September) [10.1109/3DV.2018.00078].

Terms of use:

The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

18/05/2024 21:40

(Article begins on next page)

Learning to Generate Facial Depth Maps

Stefano Pini Filippo Grazioli Guido Borghi Roberto Vezzani Rita Cucchiara
Department of Engineering “Enzo Ferrari”
University of Modena and Reggio Emilia
Via Vivarelli 10, Modena 41125, 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. Two public datasets, namely Biwi database and Pandora dataset, are exploited to demonstrate that the proposed model generates high-quality synthetic depth images, both in terms of visual appearance and informative content. 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.

1. Introduction

Depth estimation is a task at which humans naturally excel thanks to the presence of two high-quality stereo cameras (*i.e.* the human eyes) and an exceptional learning tool (*i.e.* the human brain).

What makes humans so excellent at estimating depth even from a single monocular image and how does this learning process happen?

One hypothesis is that we develop the faculty to estimate the 3D structure of the world through our past visual experience, which consists in an extremely large number of observations associated with tactile stimuli (for small objects) and movements (for wider spaces) [44].

This process allows humans to develop the capability to infer the structural model of objects and scenes they see, even from monocular images. Even though depth estimation is a natural human brain activity, the task is an ill-posed problem in the computer vision context, since the same 2D image may be generated by different 3D maps. Moreover, the translation between these two domains is demanding due to



Figure 1. Examples of the RGB face images (first row), ground-truth depth images (second row), and in the last row the depth maps estimated by the proposed model. Not only are generated face images visually realistic and pleasant, but they also preserve discriminative features for the face verification task.

the extremely different source of information that belong to intensity images and depth maps: texture and shape data, respectively.

Traditionally, the computer vision community has broadly addressed the problem of depth estimation in different ways, as *Stereo Cameras* [16, 41], *Structure from Motion* [4, 6], and *Depth from shading and light diffusion* [35, 37]. The mentioned methods suffer from different issues, like depth homogeneity and missing values, resulting in holes in depth images. Additional challenging elements are related to the camera calibration, setup, and post-processing steps that can be time consuming and computational expensive. Recently, thanks to the advances of deep neural networks, the research community has investigated the monocular depth estimation task from intensity images in order to overcome to the previously reported issues [1, 8, 10, 11, 44].

This paper presents a framework for the generation of depth maps from monocular intensity images of human faces. An adversarial approach [12, 28] is employed to

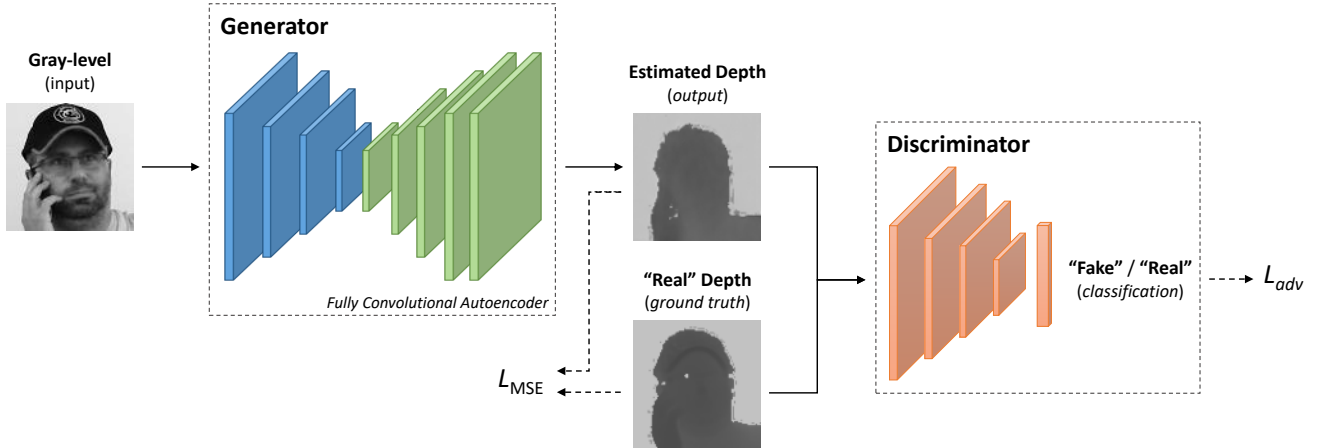


Figure 2. Overall view of the proposed Conditional GAN architecture. The architecture of the *Generator* is: k5n128s2 - k5n256s2 - k5n512s2 - k5n1024s2 - k5n512s2 - k5n256s2 - k5n128s2 - k5n64s2 - k5n1s1. The architecture of the *Discriminator* is: k5n128s2 - k5n256s2 - k5n512s2 - k5n1024s2 - fc1. k, n, s, fc correspond to the kernel size, the number of feature maps, the stride, and the number of fully connected units, respectively.

effectively train a fully convolutional encoder-decoder architecture that is able to estimate facial depth maps from the corresponding gray-level images. To train and test the proposed method, two public datasets, namely *Pandora* [3] and *Biwi Kinect Head Pose* [9] dataset, that consists of a great amount of paired depth and intensity images, are exploited. To the best of our knowledge, this is one of the first attempts to tackle this task through an adversarial approach that, differently from global depth scene estimation, involves small sized objects, full of details: the human faces.

Finally, we investigate how to effectively measure the performance of the system, introducing a variety of *pixel-wise* metrics. Besides, we introduce a *Face Verification* model, 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 networks.

2. Related Work

We consider related work within two distinct domains: *Monocular Depth Estimation* and *Domain Translation*, detailed in the following sections.

Monocular Depth Estimation. Facial depth estimation from monocular images has been investigated during the last decade. In [33], *Constrained Independent Component Analysis* is exploited for depth estimation from various pose of 2D faces. A nonlinear least-square model is employed in [34] to predict the 3D structure of the human face. Both methods rely on an initial face parameters detection based

on facial landmarks. Consequently, the final estimation is influenced by detection performance and, therefore, head pose angles. In [43], this task is tackled as a statistical learning problem, using the *Local Binary Patterns* as feature, but only frontal faces are taken into account. Reiter *et al.* [29] propose the use of canonical correlation analysis to predict depth information from RGB frontal face images. In [17], the face *Delaunay Triangulation* is exploited in order to estimate depth maps based on similarity. Cui *et al.* [5] propose a cascade FCN and a CNN architecture to estimate the depth information. In particular, the FCN aims to recover the depth from an RGB image, while the CNN is employed during the training phase to maintain the original subject's identity.

A wide body of literature addresses the monocular depth estimation task in the automotive context [1, 11] or in indoor scenarios [22, 39]. The *Markov Random Field* (MRF) and the linear regression are employed in [31] to predict depth values. An evolution of the MRF combines the 3D orientation and position of segmented patches within RGB images [32]. The main challenges to these works is that depth values are locally predicted therefore scene depth prediction lacks of global coherence. To improve the global scene depth prediction accuracy, sparse coding has been investigated in [2], while semantic labels are exploited in [18].

Recently, this research field has received a great improvement thanks to the introduction of Convolutional Neural Networks [7, 8, 23]. Several works propose the use of RGB samples paired with depth images as ground truth data in order to learn how to estimate depth maps by means of a supervised approach. In [7, 8], a two-scale network is trained on intensity images to produce depth values. The main issue is related to the limited size of publicly-available

training data and the overall low image quality [19, 21].

In this work, we aim to investigate the adversarial training approach [12] in order to propose a method that directly estimate facial depth maps, without a-priori facial feature detection, like facial landmarks or head pose angles.

Domain Translation. The *domain translation* task, which is often referred as *image translation* in the computer vision community, consists in learning a parametric mapping function between two distinct domains.

Image-to-image translation problems are often formulated as a per-pixel classification or regression [13, 20, 25, 38, 42]. Borghi *et al.* [3] propose an approach for computing the appearance of a face using the corresponding depth information based on a traditional CNN combining aspects of autoencoders [27] and Fully Convolutional Networks [25]. Recently, a consistent body of literature has addressed the image-to-image translation problem by exploiting *conditional Generative Adversarial Networks* (cGANs) [28] in order to learn a mapping between two image domains. Wang *et al.* [36] proposed a method, namely *Style GAN*, that renders a realistic image from a synthetic one. Isola *et al.* [14] demonstrated that their model, called *pix2pix*, is effective at synthesizing photos from semantic labels, reconstructing objects from edge maps, and colorizing images. In [24], a framework of coupled GANs, which is able to generate pairs of corresponding images in two different domains, was proposed.

We tackle the domain translation task in order to generate facial depth maps which are visually pleasant and contain enough discriminative information for the face verification task.

3. Proposed Method

In this section, we present the proposed model for depth estimation from face intensity images, detailing the cGAN architecture (Section 3.1.1), its training procedure (Section 3.1.2), and the adopted pre-processing face crop algorithm (Section 3.2). The implementation of the model follows the guidelines proposed in [12].

3.1. Depth Estimation Model

Following the work of Goodfellow *et al.* [12] and Mirza *et al.* [28], the proposed architecture is composed of a generative network G and a discriminative network D . G corresponds to an estimation function that predicts the depth map $I^{gen} = G(I^{gray})$ of a given face gray-level image I^{gray} , while D corresponds to a discriminative function that distinguishes between original (*i.e.* “real”) and generated (*i.e.* “fake”) depth maps.

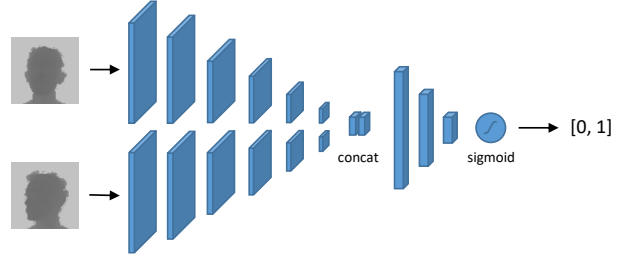


Figure 3. Architecture of the Siamese network: k3n64s2 - k3n128s2 - k3n256s2 - k3n256s2 - k3n256s2 - avgpoolk2 - fc128 - fc32 - fc1. k, n, s, fc correspond to kernel size, number of feature maps, stride, and fully-connected units. The output of the network is a continuous similarity score.

3.1.1 Network architecture

Generator. The generator network is based on the fully convolutional architecture depicted in Figure 2 that, following the paradigm of *conditional GANs*, takes a face intensity image as input and estimates the corresponding depth map.

The first part of the network acts as an encoder, mapping the input image into a 1024-dimensional embedding with a spatial size 16 times smaller than the input one. It is composed of four convolutional layers with kernel size 5, stride 2, and 128, 256, 512, and 1024 features maps, respectively. Each layer is followed by the *Leaky ReLU* activation function [26] with a negative slope of 0.2.

The second part of the network, acting as a decoder, generates a depth image by processing the face embedding produced by the encoder. It is composed of four transposed convolution layers (also known in the literature as fractionally-strided convolution layers) which increase the embedding resolution up to the original image size. The layers are applied with kernel size 5, stride 2, and 512, 256, 128, and 64 features maps, respectively, and they are followed by the *ReLU* activation function.

Then, a standard convolutional layer, followed by a hyperbolic tangent activation function (*tanh*), produces the final depth map estimation. Batch normalization is employed before each activation function, except the last one, for regularization purposes.

Discriminator. The discriminator network, depicted in Figure 2, takes as input a depth map and predicts the probability of the input to be a real or a generated depth map.

The first part of the discriminator shares the same architecture with the encoder part of the generator network.

Then, the 1024-dimensional embedding is flattened and a fully connected layer with one unit and a sigmoid activation function are applied obtaining a final score in the range [0, 1] where 0 corresponds to a “real” depth map and 1 to a “fake” one.

3.1.2 Adversarial training

During the training procedure, the discriminator network D , with parameters θ_d , is trained to predict whether a depth map is “real” or “fake” by maximizing the probability of assigning the correct label to each sample. Meanwhile, the generator network G , with parameters θ_g , is trained in order to generate realistic depth maps and fool the discriminator D . From a mathematical perspective, the training can be formalized as the optimization of the following min-max problem:

$$\min_{\theta_g} \max_{\theta_d} \mathbb{E}_{x \sim p_{dpt}(x)} [\log(D(x))] + \mathbb{E}_{y \sim p_{gray}(y)} [\log(1 - D(G(y)))] \quad (1)$$

where $D(x)$ is the probability of being a “real” depth image (consequently $1 - D(G(y))$ is the probability to be a “fake” depth image), p_{dpt} is the distribution of the real depth maps, and p_{gray} is the distribution of the intensity images. This approach leads to a generative model which is capable of generating “fake” images that are highly similar to the “real” ones, thus indistinguishable by the discriminator D .

To reach this goal, the following loss functions are employed during the training with the *Adam* optimizer [15] with initial learning rate of 0.0002 and betas 0.5 and 0.999. Regarding the discriminator network, a binary categorical cross entropy loss function, defined as

$$L_{adv}(\mathbf{y}, \mathbf{t}) = -\frac{1}{N} \sum_{i=1}^N [t_i \log y_i + (1 - t_i) \log(1 - y_i)] \quad (2)$$

where $y_i = D(I_i)$ is the discriminator prediction regarding the i -th input depth map and t_i is the corresponding ground truth, is applied to the discriminator output.

Regarding the generator network, we aim to generate images that are similar to the ground truth depth maps as well as capable of fooling the discriminator network (*i.e.* visually indistinguishable from real depth maps). For fulfilling the first goal, we apply the Mean Squared Error (MSE) loss function:

$$L_{MSE}(\mathbf{s}^g, \mathbf{s}^d) = \frac{1}{N} \sum_{i=1}^N \|G(s_i^g) - s_i^d\|_2^2 \quad (3)$$

where \mathbf{s}^g and \mathbf{s}^d are respectively the input gray-level images and the target depth maps.

To accomplish the second goal, we feed the discriminator with generated depth images and apply the adversarial loss on the discriminator prediction to evaluate if the generated images are capable of fooling the discriminator network. Then, we back-propagate the gradients up to the input of the generator network and update the generator parameters while keeping fixed the discriminator weights. There-

fore, we aim to solve the back-propagation problem minimizing:

$$\hat{\theta}_g = \arg \min_{\theta_g} L_G(\mathbf{s}^g, \mathbf{s}^d) \quad (4)$$

where L_G is a combination of two components, defined as the following weighted sum:

$$L_G(\mathbf{s}^g, \mathbf{s}^d) = \lambda \cdot L_{MSE}(\mathbf{s}^g, \mathbf{s}^d) + L_{adv}(G(\mathbf{s}^g), \mathbf{1}) \quad (5)$$

in which λ is a weighting parameter that controls the impact of the L_{MSE} loss with respect to the adversarial loss.

3.2. Dynamic face crop

The head detection task is out of the scope of this paper, therefore a trivial dynamic face crop algorithm is adopted in order to accurately extract face bounding boxes from the considered datasets including a small portion of background. In particular, given the head center position (x_H, y_H) in a depth map (we assume that the head center position is provided in the dataset annotations), a bounding box of width w_H and height h_H is extracted defining its width and height as

$$w_H = \frac{f_x \cdot R_x}{D} \quad h_H = \frac{f_y \cdot R_y}{D} \quad (6)$$

where R_x, R_y are the average width and height of a face (we consider $R_x = R_y = 320$), f_x, f_y are the horizontal and the vertical focal lengths in pixels of the acquisition device (in the considered case: $f_x = f_y = 365$ and $f_x = f_y = 370$ for *Pandora* and *Biwi* datasets, respectively), and D is the distance between the head center and the acquisition device which is estimated averaging the depth values around the head center. Computed values are used to crop the face in both depth maps and intensity images.

4. Experimental Results

In this section, experimental results, obtained through a *cross-subject evaluation*, are reported. In particular, we investigate the use of pixel-wise metrics in order to verify the generation capability of the proposed adversarial model (Section 4.2). Furthermore, we evaluate the quality of the estimated facial depth maps by means of a *Face Verification* task (Section 4.3). The code and the network models are publicly released¹.

4.1. Datasets

A description of the exploited datasets, namely *Pandora* and *Biwi Kinect Head Pose*, is provided. The explanation of how the *Pandora* dataset was split to take different head poses, occlusions, and garments into account is presented in the following as well.

¹<http://imagelab.ing.unimore.it/face-generation>

Metrics		Pandora [3]				Biwi [9]			
		<i>cGAN</i>	<i>AE</i>	[14]	[5]	<i>cGAN</i>	<i>AE</i>	[14]	[5]
L_1 Norm		11.792	16.185	18.172	19.046	10.503	10.444	47.191	16.507
L_2 Norm	↓	1,678.2	2,224.8	3,109.0	2,093.3	2,368.5	2,342.5	6,661.3	2,319.8
Absolute Diff		0.1019	0.1441	0.1512	0.1465	0.1838	0.1936	0.9062	0.2836
Squared Diff	↓	2.9974	5.3891	8.6444	3.9084	8.7122	9.0332	100.89	9.3032
RMSE _{lin}		18.677	25.213	33.526	22.599	24.865	24.699	72.084	24.521
RMSE _{log}	↓	0.1744	0.2752	1.0864	0.2105	0.2932	0.2970	1.2240	0.3390
RMSE _{scale-inv}		0.1345	0.2018	1.0774	0.1301	0.2687	0.2642	1.1759	0.2867
$\delta < 1.25$		0.8529	0.6854	0.7802	0.7556	0.7393	0.7230	0.4149	0.6395
$\delta < 1.25^2$	↑	0.9642	0.8728	0.8978	0.9554	0.9224	0.9064	0.5298	0.7943
$\delta < 1.25^3$		0.9915	0.9651	0.9638	0.9935	0.9609	0.9557	0.6360	0.9311
$\delta < 1.25^{\frac{1}{4}}$		0.5619	0.4763	0.5053	0.0740	0.5293	0.5532	0.3128	0.3580
$\delta < 1.25^{\frac{1}{3}}$	↑	0.6113	0.5047	0.5615	0.1036	0.5727	0.6034	0.3274	0.4549
$\delta < 1.25^{\frac{1}{2}}$		0.6954	0.5547	0.6496	0.5509	0.6348	0.6606	0.3513	0.5189
Face Verification	↑	0.7247	0.6570	0.5315	0.6442	0.6251	0.6043	0.5422	0.5966

Table 1. Accuracy comparison for the *pixel-wise* metrics and the *face verification* task, as a function of different generative approaches. *cGAN* and *AE* refer to the generative adversarial network and the encoder-decoder architecture proposed in this paper, respectively. Results are reported for both the *Pandora* and *Biwi* datasets and for the main competitors [5, 14]. Starting from the top, L_1 and L_2 distances are reported, then the absolute and the squared differences, the root mean squared error, and the percentage of pixels under a certain error threshold (δ -metrics). Further details about metrics are reported in [8] and in Section 4.4. Finally, the accuracy on the face verification task, detailed in Section 4.3, is reported. The arrows next to the metrics represent the positive changing direction: a better generation performance corresponds to a metric variation in the arrow direction.

4.1.1 Pandora Dataset

The Pandora dataset was introduced in [3] for the head pose estimation task in depth images. It consists of more than 250k paired face images, both in the RGB and the depth domain.

Depth maps are acquired with the *Microsoft Kinect One* device (also known as *Microsoft Kinect for Windows v2*), a *Time-of-Flight* sensor that assures great quality and high resolution for both the RGB (1920×1080 pixels) and the depth (512×424 pixels) data.

Even though the dataset was not created for the depth generation task, it can be successfully employed for that purpose as well, as it contains paired RGB-depth images.

Furthermore, it includes some challenging features, such as the presence of garments, numerous face occlusions created by objects (e.g. bottles, smartphones, and tablets) and arms, and extreme head poses (*roll*: $\pm 70^\circ$, *pitch*: $\pm 100^\circ$, *yaw*: $\pm 125^\circ$). As reported in the original paper, we use subjects number 10, 14, 16, and 20 as testing subset.

Each subject presents 5 different sequences S_i of frames. We split the sequences into two sets. The first one, referred as $\{S_1, S_2, S_3\}$, contains actions performed with constrained movements (yaw, pitch, and roll vary one at a time), for both the head and the shoulders. The second set, referred as $\{S_4, S_5\}$, consists of both complex and simple movements, as well as occlusions and challenging camou-

flage. Experiments are performed on both the subsets in order to investigate the effects of the mentioned differences. Moreover, we additionally split the dataset taking head pose angles into account. We create two mutually-exclusive head pose-based subsets, defined as

$$A_1 = \{s_{\rho\theta\sigma} \mid \forall \gamma \in \{\rho, \theta, \sigma\} : -10^\circ \leq \gamma \leq 10^\circ\} \quad (7a)$$

$$A_2 = \{s_{\rho\theta\sigma} \mid \exists \gamma \in \{\rho, \theta, \sigma\} : \gamma < -10^\circ \vee \gamma > 10^\circ\} \quad (7b)$$

where ρ , θ , and σ are the yaw, the pitch, and the roll angle, respectively, for each sample $s_{\rho\theta\sigma}$. In practice, A_1 consists of frontal face images, while non-frontal face images are included in A_2 .

When using the dataset for the face verification task, the problem of the high number of possible dataset image pairs arises. To overcome the issue, we created two fixed set of image pairs, a validation and a test set, in order to allow repeatable and comparable experiments.

We extract face images from dataset frames using the automatic face cropping technique presented in Section 3.2 then we resize them to the size of 96×96 pixels. We exclude from the dataset a very small subset of extreme head poses and occlusions, as well as frames in which the automatic cropping algorithm does not work properly, to avoid training instability.

	$\{S_i\}_{i=1,2,3}$		$\{S_i\}_{i=4,5}$		$\{S_i\}_{i=1,2,3,4,5}$	
	<i>original</i>	<i>generated</i>	<i>original</i>	<i>generated</i>	<i>original</i>	<i>generated</i>
A_1	0.8184	0.8614	0.7685	0.7155	0.7917	0.7950
A_2	0.7928	0.7499	0.7216	0.6586	0.7576	0.7007
$\{A_1, A_2\}$	0.8034	0.7851	0.7271	0.6696	0.7664	0.7247

Table 2. Confusion matrix for the Face Verification task on the *Pandora* dataset, as a function of different angle and sequence subsets on original and generated depth images. Subsets description is reported in Section 4.1.1. Tested generated images are estimated by the proposed network trained according to the adversarial approach described in Section 3.

4.1.2 Biwi Kinect Head Pose Database

The *Biwi Kinect Head Pose Database* was introduced by Fanelli *et al.* [9] in 2013. Differently from *Pandora*, it is acquired with the first version of the *Microsoft Kinect*, a *structured-light* infrared sensor. With respect to ToF sensors, this Microsoft Kinect version provides lower quality depth maps [30], in which it is common to find holes (missing depth values).

The dataset consists of about 15k frames, split in 24 sequences of 20 different subjects (four subjects are recorded twice). Both RGB and depth images have the same spatial resolution of 640×480 pixels. The head pose angles span about $\pm 50^\circ$ for *roll*, $\pm 60^\circ$ for *pitch*, and $\pm 75^\circ$ for *yaw*. We adopt the same procedure used for *Pandora* to crop the faces and obtain 96×96 pixel images.

4.2. Pixel-wise metrics

The overall quality of the generated facial depth maps is evaluated with the pixel-wise metrics proposed in [8].

In particular, the generation capability of the generator network, trained both as an encoder-decoder architecture and with the adversarial policy reported in Section 3.1.2, is compared with the recent *pix2pix* architecture [14] and the algorithm proposed in [5]. We test the models on both the *Pandora* and the *Biwi* dataset.

As reported in Table 1, our generator network, trained as a cGAN, performs better than the encoder-decoder architecture and the literature competitors. The δ -metrics, corresponding to the percentage of pixels under a certain error threshold, confirm that the proposed model achieves the best spatial accuracy, in particular on the *Pandora* dataset, which contains more fine details than the *Biwi* dataset. As highlighted in right part of the above-mentioned table, the limited size and variability of the *Biwi* dataset have a negative impact on the generative and the generalization capability of the tested architectures.

4.3. Face Verification test

Pixel-wise metrics allow for a mathematical evaluation of the generative performance of deep convolutional networks. Yet, they might not fully convey whether the original domain features are accurately preserved through the

generative process. Even when a human observer perceives no difference between “real” and “fake” images, the information content might still be represented in a slightly different fashion in terms of texture, colors, geometries, light intensity, and fine details.

In order to deeply investigate the quality of the generated images, the following Face Verification test, *i.e.* determining whether two given face images belong to the same subject, is employed. We exploit a deep convolutional Siamese network trained on original depth images, without adopting any kind of fine-tuning on the generated depth maps. The model compares two depth images and predicts their similarity as a value in the $[0, 1]$ range. Two input faces are considered as belonging to the same person if their similarity score is higher than 0.5.

The model architecture is depicted in Figure 3. It is composed by 5 convolutional layers with an increasing number of feature maps, an average pooling layer, and 3 fully connected layers.

Through this test, which needs both high and low level features in order to work properly, we can estimate how well the generated faces preserve the individual visual features of the subjects. As reported in the last line of Table 1, the proposed cGAN model allows for the highest test accuracy on the face verification task. Actually, the proposed architecture significantly outperforms the other tested architectures. In light of the reported results, we believe that our model is able to estimate both high level and fine details of facial depth maps and hence to obtain realistic and discriminative synthetic images.

It is worth to notice that our architecture overcomes the model proposed in [5], in which the deep model is specifically trained to preserve subject identities by minimizing a dedicated face recognition loss.

Table 2 presents how the Face Verification accuracy varies as a function of different head poses and image complexity (*i.e.* obstructions, garments, unconstrained movements) on the *Pandora* dataset, as reported in Section 4.1.1. Furthermore, we report the verification accuracy obtained by the same model on the original depth maps.

As expected, the Siamese network performs better on frontal faces (angle subset A_1) and constrained movements

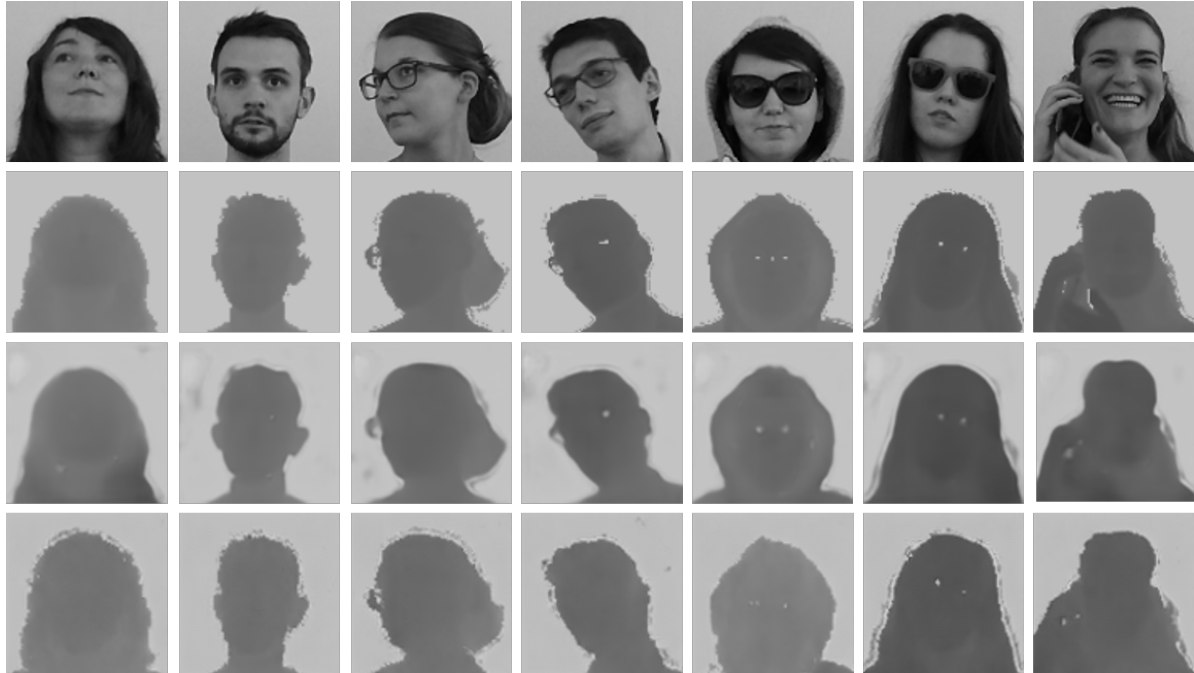


Figure 4. Sample outputs of the proposed method (last row). In the first and second rows, original gray-level and depth maps are reported, respectively. In the third row, depth maps generated by the encoder-decoder architecture are depicted.

(sequence subset $\{S_1, S_2, S_3\}$). Nevertheless, even the most complex task ($A_2, \{S_4, S_5\}$) reaches a 72.47% accuracy. Surprisingly, in the case of frontal faces and constrained movements, the generated depth images lead to better accuracy than the original ones. We hypothesize that this behavior is related to the generative process that tends to remove partial occlusions and to produce highly discriminative facial features.

4.4. Result discussion

Since this is one of the first work that tackles the facial depth maps generation through a Conditional GAN, a further evaluation of the proposed method and the obtained results is here reported. Indeed, we believe that the proposed method and this investigation could be a useful baseline for future work in this research field.

From Table 1, it can be noticed that the detail accuracy of the proposed model is quite good, compared with the tested competitors. However, we highlight that the δ -metrics reported in lines 8-10, often exploited in previous work, are effective to check the overall quality of depth maps generated from landscapes or wide-angle scenes [1, 40], but the threshold value is too high to take fine details into account. For this reason, we introduce a new set of δ -metrics with harder thresholds, reported in lines 11-13 of Table 1. We believe that this refined version of the δ -metrics could be useful for evaluation of future research.

Finally, we note that our approach is able to produce overall

accurate views of the generated facial depth maps, preserving the shape of the face and the garments (and also the characteristic signal noise of the *Microsoft Kinect One*), but it still struggles to produce accurate fine facial details. Detail preservation during the depth map generation process, after all, is still an open problem in the literature. Therefore, we suppose that the face verification network is more influenced by the shape of the face than by the fine facial details. In future work, we plan to investigate methods to produce more accurate and realistic depth maps, including better facial details, like the mouth, the ocular cavities, and the nose.

5. Conclusion

In this paper, we present an approach for the estimation of facial depth maps from intensity images.

In order to evaluate the quality of the generated images, we perform a face verification task employing a Siamese network pre-trained on the original depth maps.

By showing that the Siamese network accuracy does not degrade when tested on generated images, we demonstrate that the proposed framework produces high-quality depth maps, both in terms of visual appearance and discriminative information. We also demonstrate that the proposed architecture outperforms the encoder-decoder architecture and the literature competitors when trained with the adversarial policy. Thanks to the flexibility of our approach, we plan to extend our model by introducing task-specific losses and to apply it to different scenarios.

References

- [1] A. Atapour-Abarghouei and T. Breckon. Real-time monocular depth estimation using synthetic data with domain adaptation via image style transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, volume 18, 2018. [1](#), [2](#), [7](#)
- [2] M. H. Baig, V. Jagadeesh, R. Piramuthu, A. Bhardwaj, W. Di, and N. Sundaresan. Im2depth: Scalable exemplar based depth transfer. In *Applications of Computer Vision (WACV), 2014 IEEE Winter Conference on*, pages 145–152. IEEE, 2014. [2](#)
- [3] G. Borghi, M. Venturelli, R. Vezzani, and R. Cucchiara. Poseidon: Face-from-depth for driver pose estimation. In *IEEE International Conference on Computer Vision and Pattern Recognition*. IEEE, 2017. [2](#), [3](#), [5](#)
- [4] P. Cavestany, A. L. Rodríguez, H. Martínez-Barberá, and T. P. Breckon. Improved 3d sparse maps for high-performance sfm with low-cost omnidirectional robots. In *Image Processing (ICIP), 2015 IEEE International Conference on*, pages 4927–4931. IEEE, 2015. [1](#)
- [5] J. Cui, H. Zhang, H. Han, S. Shan, and X. Chen. Improving 2d face recognition via discriminative face depth estimation. *Proc. ICB*, pages 1–8, 2018. [2](#), [5](#), [6](#)
- [6] L. Ding and G. Sharma. Fusing structure from motion and lidar for dense accurate depth map estimation. In *Acoustics, Speech and Signal Processing (ICASSP), 2017 IEEE International Conference on*, pages 1283–1287. IEEE, 2017. [1](#)
- [7] D. Eigen and R. Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2650–2658, 2015. [2](#)
- [8] D. Eigen, C. Puhrsch, and R. Fergus. Depth map prediction from a single image using a multi-scale deep network. In *Advances in neural information processing systems*, pages 2366–2374, 2014. [1](#), [2](#), [5](#), [6](#)
- [9] G. Fanelli, M. Dantone, J. Gall, A. Fossati, and L. Van Gool. Random forests for real time 3d face analysis. *Int. J. Comput. Vision*, 101(3):437–458, February 2013. [2](#), [5](#), [6](#)
- [10] R. Garg, V. K. BG, G. Carneiro, and I. Reid. Unsupervised cnn for single view depth estimation: Geometry to the rescue. In *European Conference on Computer Vision*, pages 740–756. Springer, 2016. [1](#)
- [11] C. Godard, O. Mac Aodha, and G. J. Brostow. Unsupervised monocular depth estimation with left-right consistency. In *CVPR*, volume 2, page 7, 2017. [1](#), [2](#)
- [12] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014. [1](#), [3](#)
- [13] S. Iizuka, E. Simo-Serra, and H. Ishikawa. Let there be color!: joint end-to-end learning of global and local image priors for automatic image colorization with simultaneous classification. *ACM Transactions on Graphics (TOG)*, 35(4):110, 2016. [3](#)
- [14] P. Isola, J. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2017. [3](#), [5](#), [6](#)
- [15] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. [4](#)
- [16] K. Konda and R. Memisevic. Unsupervised learning of depth and motion. *arXiv preprint arXiv:1312.3429*, 2013. [1](#)
- [17] D. Kong, Y. Yang, Y.-X. Liu, M. Li, and H. Jia. Effective 3d face depth estimation from a single 2d face image. In *Communications and Information Technologies (ISCIT), 2016 16th International Symposium on*, pages 221–230. IEEE, 2016. [2](#)
- [18] L. Ladicky, J. Shi, and M. Pollefeys. Pulling things out of perspective. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 89–96, 2014. [2](#)
- [19] I. Laina, C. Rupprecht, V. Belagiannis, F. Tombari, and N. Navab. Deeper depth prediction with fully convolutional residual networks. In *3D Vision (3DV), 2016 Fourth International Conference on*, pages 239–248. IEEE, 2016. [3](#)
- [20] G. Larsson, M. Maire, and G. Shakhnarovich. Learning representations for automatic colorization. In *European Conference on Computer Vision*, pages 577–593. Springer, 2016. [3](#)
- [21] B. Li, C. Shen, Y. Dai, A. van den Hengel, and M. He. Depth and surface normal estimation from monocular images using regression on deep features and hierarchical crfs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1119–1127, 2015. [3](#)
- [22] J. Li, R. Klein, and A. Yao. A two-streamed network for estimating fine-scaled depth maps from single rgb images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3372–3380, 2017. [2](#)
- [23] F. Liu, C. Shen, G. Lin, and I. Reid. Learning depth from single monocular images using deep convolutional neural fields. *IEEE transactions on pattern analysis and machine intelligence*, 38(10):2024–2039, 2016. [2](#)
- [24] M.-Y. Liu and O. Tuzel. Coupled generative adversarial networks. In *Advances in neural information processing systems*, pages 469–477, 2016. [3](#)
- [25] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3431–3440, 2015. [3](#)
- [26] A. L. Maas, A. Y. Hannun, and A. Y. Ng. Rectifier nonlinearities improve neural network acoustic models. In *ICML Workshop on Deep Learning for Audio, Speech, and Language Processing*, 2013. [3](#)
- [27] J. Masci, U. Meier, D. Cireşan, and J. Schmidhuber. Stacked convolutional auto-encoders for hierarchical feature extraction. In *International Conference on Artificial Neural Networks*, pages 52–59. Springer, 2011. [3](#)
- [28] M. Mirza and S. Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014. [1](#), [3](#)
- [29] M. Reiter, R. Donner, G. Langs, and H. Bischof. Estimation of face depth maps from color textures using canonical correlation analysis. *na*, 2006. [2](#)
- [30] H. Sarbolandi, D. Lefloch, and A. Kolb. Kinect range sensing: Structured-light versus time-of-flight kinect. *Computer Vision and Image Understanding*, 2015. [6](#)

- [31] A. Saxena, S. H. Chung, and A. Y. Ng. Learning depth from single monocular images. In *Advances in neural information processing systems*, pages 1161–1168, 2006. 2
- [32] A. Saxena, M. Sun, and A. Y. Ng. Make3d: Learning 3d scene structure from a single still image. *IEEE transactions on pattern analysis and machine intelligence*, 31(5):824–840, 2009. 2
- [33] Z.-L. Sun and K.-M. Lam. Depth estimation of face images based on the constrained ica model. *IEEE Transactions on Information Forensics and Security*, 6(2):360–370, 2011. 2
- [34] Z.-L. Sun, K.-M. Lam, and Q.-W. Gao. Depth estimation of face images using the nonlinear least-squares model. *IEEE transactions on image processing*, 22(1):17–30, 2013. 2
- [35] M. W. Tao, P. P. Srinivasan, J. Malik, S. Rusinkiewicz, and R. Ramamoorthi. Depth from shading, defocus, and correspondence using light-field angular coherence. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1940–1948, 2015. 1
- [36] X. Wang and A. Gupta. Generative image modeling using style and structure adversarial networks. In *European Conference on Computer Vision*, pages 318–335. Springer, 2016. 3
- [37] R. J. Woodham. Photometric method for determining surface orientation from multiple images. *Optical engineering*, 19(1):191139, 1980. 1
- [38] S. Xie and Z. Tu. Holistically-nested edge detection. In *Proceedings of the IEEE international conference on computer vision*, pages 1395–1403, 2015. 3
- [39] D. Xu, E. Ricci, W. Ouyang, X. Wang, and N. Sebe. Multi-scale continuous crfs as sequential deep networks for monocular depth estimation. In *Proceedings of CVPR*, 2017. 2
- [40] D. Xu, W. Wang, H. Tang, H. Liu, N. Sebe, and E. Ricci. Structured attention guided convolutional neural fields for monocular depth estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3917–3925, 2018. 7
- [41] K. Yamaguchi, T. Hazan, D. McAllester, and R. Urtasun. Continuous markov random fields for robust stereo estimation. In *European Conference on Computer Vision*, pages 45–58. Springer, 2012. 1
- [42] J. Zhao, M. Mathieu, and Y. LeCun. Energy-based generative adversarial network. *arXiv preprint arXiv:1609.03126*, 2016. 3
- [43] Y. Zheng and Z. Wang. Robust depth estimation for efficient 3d face reconstruction. In *Image Processing, 2008. ICIP 2008. 15th IEEE International Conference on*, pages 1516–1519. IEEE, 2008. 2
- [44] T. Zhou, M. Brown, N. Snavely, and D. G. Lowe. Unsupervised learning of depth and ego-motion from video. In *CVPR*, volume 2, page 7, 2017. 1