

# Deep Generation of Face Images On basis of Sketches

Priyanka Jadhav, Harshal Pawar, Akanksha Kalbhor, Sonali Dongare

Nutan Maharashtra Institute of Engineering and Technology, Pune, India

**Abstract:** New deep image-to-image translation techniques or methods enable rapid generation of face images from incomplete or rough freehand sketches. However, existing solutions adapt too much to sketches and therefore require edge maps or professional sketches as input. To solve this problem, our main idea is to implicitly model the shape space of face images and synthesize it in this space to approximate the input sketch. We take a local to global approach. We study the insertion of elements into the main components of primary surfaces and transfer the corresponding parts of the input sketches towards the basic component varieties defined by the feature vectors of the surface component samples. Here is also another deep neural network that learns the mapping function from built-in component features to realistic images with various multi-channel feature maps as mediating results to improve information flow. Because our method basically uses input incomplete or rough freehand sketches as soft links, and thus is able to create realistic face images even from incomplete or rough sketches. Because our tool is very easy to use even for untrained artists, while still helping by providing fine control over shape details. Quantitative and qualitative analysis shows the high generation capacity of our system for existing and new other solutions. The fluency and practicality of our system is confirmed by a user study.

**Keywords:** Image-to-Image Translation, Feature Embedding, Sketch-based Generation, Face Synthesis.

## I. INTRODUCTION

Creating real images of human faces from scratch or incomplete freehand sketches benefits a variety of applications including crime investigation, character design, educational training, etc. Due to their clarity, brevity, and ease of use, sketches are often used to illustrate desired faces. Recently introduced image-to-image translation techniques based on deep learning enable the automatic generation of photographic images from sketches for various categories of objects, including human faces, and help produce spectacular results.

Most of these deep learning-based sketch-to-image systems often start with input sketches that are nearly fixed and then try to infer missing textures or shading information between strokes. These problems are written somewhat more like reconstruction problems with hard constraints on the input sketches. Due to their data-driven nature, they often train their networks from pairs of real photos and their corresponding edge maps, and therefore need test sketches with a quality comparable to that of the edge maps of real images to synthesize realistic face images. Creating such sketches is challenging, especially for people with less sketching experience. Our system consists of three main modules, namely FM (Feature Mapping), IS (Image Synthesis) and CE (Component Embedding). The CE module acquires an auto-encoder architecture and additionally learns five feature descriptors from the face sketch data, specifically for “right eye”, “left eye”, “nose”, “mouth” and “residues” for locally approximating component distributions. The other FM and IS modules form different deep learning subnets to generate conditional images, and the map components contain vectors for realistic or real images. Our main concept is to implicitly learn a space of plausible face sketches from real face data to solve this problem.

To deal with this problem, our primary idea is to study the space of plausible face sketches from real face sketch images and find the closest point in this space to estimate the input sketch. In this way, image synthesis can be guided by using rather soft sketch bindings. Thus, the credibility of the synthesized images can be increased even with rough and/or incomplete input sketches, without disregarding the characteristics shown in the sketches. Learning such a space globally (if it exists) is not very practical due to the small amount of training data against the expected high-dimensional feature space.

This inspires us to fully model the manifolds at the component level, which makes a better sense of the estimate that each component manifold is low-dimensional and locally linear [14]. This decision not only helps to locally split such splitters using a limited amount of face data, but also verifies a finer control of shape details.

Our system consists of three main modules, namely FM (Feature Mapping), IS (Image Synthesis) and CE (Component Embedding). The CE module obtains an auto-encoder architecture and partially learns five feature descriptors from the face sketch data, namely for “right eye”, “left eye”, “nose”, “mouth” and “remains”. for a local approach to component distributions. The other FM and IS modules form different deep learning subnets to generate conditional images, and the map components contain vectors for realistic or real images. It tries to infer missing textures or shading information between strokes.

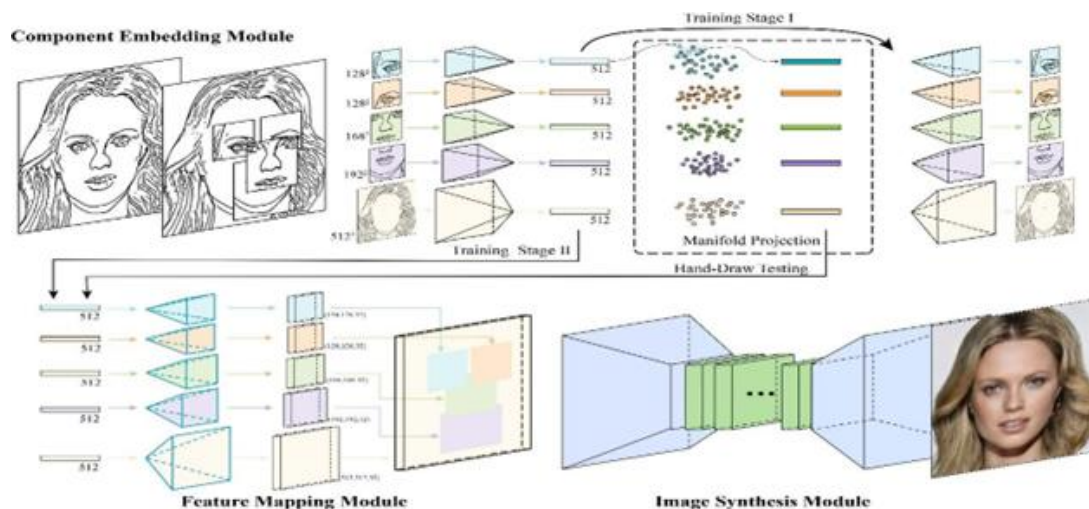


Fig1

## II.METHODOLOGY

The 3D shape space of human faces has been well thought out (see the classic morphable face model. A viable approach to synthesizing realistic faces from rough freehand sketches is to first predict an input sketch into such a 3D face space and also synthesize a face image from the generated 3D surface. Although such a global parametric model is not variable enough to align the rich details of the image or support local adjustments, which presents us with the advantages of a local-global structure for faithful local synthesis of details, our method focuses on modeling the shape spaces of facial components in the image processing domain.

We can do this by learning how to insert elements into face components. For each component type, the points corresponding to the component samples implicitly define the manifold. Although we don't explicitly study this manifold because we focus a lot on knowing the closest point in such a manifold given the new sketched face component to be refined. Noting that in the connotation of embedding spaces, similar components are close to each other, we consider that the underlying manifolds of the components are locally linear. We then follow the main idea of the classical locally linear embedding (LLE) algorithm to project the feature vector of the sketched face component into its component manifold.

We use the information in the feature space by guiding the conditional sketch to synthesize the image through feature embedding learning. Unlike traditional inter-sketch synthesis methods that study conditional GANs to transform sketches into images, our idea forces the synthesis process to traverse the component feature space and then map 1-channel feature vectors to 32-channel feature maps before use. conditional GAN. This greatly improves the flow of information and benefits the fusion of components. Below we first discuss our data preparation procedure. We can then represent our new pipeline for sketch-to-image synthesis and our approach to multiple projection.

### III. DATA PREPARATION

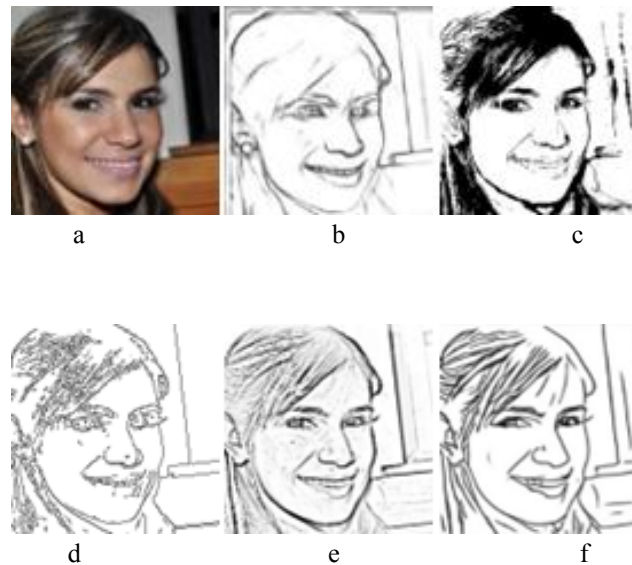


Fig.2. The comparisons of different edge extraction methods

A reasonably large dataset of face sketch-image pairs is required to train our network. There are several relevant datasets such as the CUHK face sketch database. Since a more abstract representation of surfaces using sparse lines is excluded, although sketches in this type of datasets include a shading effect. We thus contribute to a new dataset of pairs of face images and corresponding synthesized sketches. We build on the CelebAMask-HQ face image data, which contains high-resolution face images with semantic facial attribute masks. For simplicity, we now focus on the front faces, without any decorative accessories (eg glasses, face masks).

### IV. SKETCH TO IMAGE SYNTHESIS ARCHITECTURE

Module for inserting components. Since human faces share an understandable structure, we can decompose a face sketch into five components, labeled as  $S_c$ ,  $c \in \{1, 2, 3, 4, 5\}$  for "Left Eye", "Right Eye", "Nose", "Mouth" and "Remainder", separately. In order to process the details between the components, we need to easily define the first four components using four overlapping windows centered on the individual facial components (obtained from the pre-labeled segmentation mask in the datasets), as shown in Figure 3 (top left) The "Remainder" image correlated with the "Remainder" component is the same as the original sketch image, but the eyes, nose and mouth are gone. Here we deal with the "left eye" and the "right eye" separately, to explore the best flexibility of the generated faces (see two examples in Figure 3). To better control the details of individual components and each type face component, we study local feature embedding. We obtain feature descriptions of individual individual components using five auto-encoder networks, represented as  $\{E_c, D_c\}$ , where  $E_c$  is the encoder and  $D_c$  is the decoder for component  $C$ . autoencoder consists of 5 encoding layers as well as 5 decoding layers. Here we add a fully completed concatenated layer in the middle or middle to ensure that the latent descriptor has 512 dimensions for all five components.



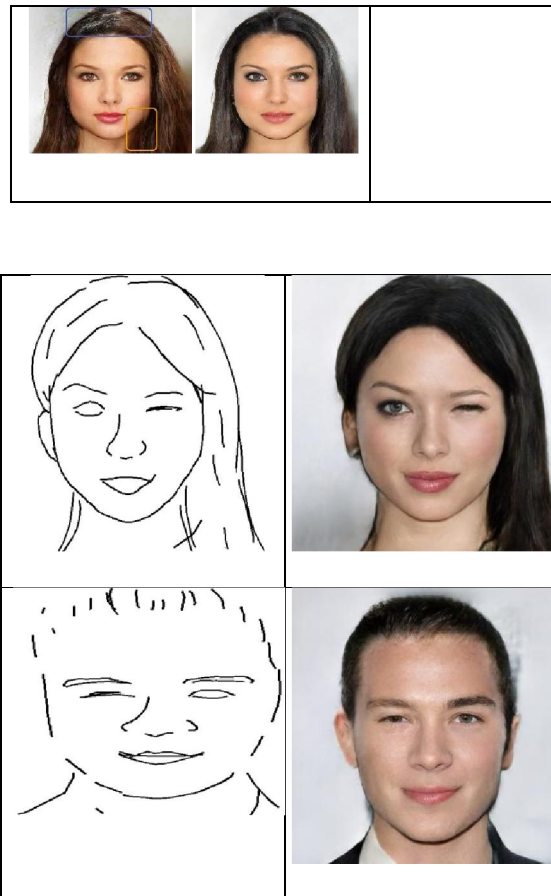


Fig. 3. examples of generation flexibility supported by using separate components for the left and right eyes.

### Feature Mapping Module:

Given the sketch, we can suggest its placement to the material pipeline to make sense. One solution the real image is first converted image vectors of projected pipe elements back to the object drawing using the {Dc} learned decoder, then process the art-image synthesis at the layer standard level and finally merge it. Brought pictures together. get a full face. However, a direct solution may lead to a conflict between the value of the local context and global model, because there is no system for managing the process of personal development.

Another solution is to first render the cut material into the finished sketch and then process the composite sketch to create the surface outline. It can be seen that this solution causes artifacts in the images as well, and these artifacts carry over to the composite images because deep learning by its very nature uses the image very hard for image-to-image synthesis. Limited as discussed previously.

Figure 4. Given the same conceptual graph (a), the conditional synthesis image on multiple projection eigenvectors yields more results than the conditional synthesis image on transition graphs (b) (c) (d). Look at the artifacts highlighted in the middle image (b) and composite result (c) of pix2pixHD. We have noticed that the above problem usually occurs in the area of overlapping clipping windows of different objects.

Since the image has only one channel, disparity between adjacent objects in the overlap areas are difficult to determine using a network of images. This motivates us to show the picture point vectors of various points for many special maps. This improves the flow of information, and combining elements instead of sketches helps resolve facial inconsistency. Because the descriptions of different objects have different meanings, we created an FM module with five independent decision models to transform image vectors into maps. Each decision model has an attached set and five layers.



It has 32 channels for each custom map and is the same size as the corresponding component in the diagram. The result of the map of "left eye", "right eye", "nose" and "mouth" is inserted into the feature map "residual" as the fact of the position of the face in the ace concept. sketch. image to save. original spatial relationships. As shown in the figure (bottom middle), we combine the maps using depth (for example "left/right" > "nose" > "mouth" > "rest").

Image synthesis module. The IS module, which compiles the maps, converts them into a real face. We use this model using a GAN structure model that uses a specific map as input to a discriminator-driven generator. Like the universal renderer in pix2pixHD, our renderer has an encoding part, a residual part, and a decision room. These units create a feedback map sequentially. Similar to [16], the separator focuses on identifying different kinds of patterns: we reduce the input to different scales and use different separators to distinguish different scale objects. We use this setting to display the expected high level of the site. Two-stage training. To train our network using our image pair data, we use a two-stage training strategy as shown in Figure 3. In the first stage, we train only the CE module by training five separate autoencoders to place features using component drawings.

Image Synthesis Module. The IS module, which assembles the maps, transforms them into a real face. We use this model using a GAN structure model that uses a specific map as input to a discriminator driven generator. Like the universal generator in pix2pixHD [16], our generator has a coding part, a residual part and a decision unit. The feedback map is made sequentially by these units.

As shown in Figure 1, we use a two-stage training strategy using our dataset of sketch image pairs to train our network (Section 3.1). In the first phase, we train only the CE module by training five separate autoencoders for feature placements using component drawings. The training is based on self-monitoring with the mean square error (MSE) loss between the drawn image and the reconstructed image. In Phase II, we adjust the parameters of the learning encoders and train all unknown networks end-to-end in FM and IS modules.

For the GAN in IS, in addition to the GAN loss, we also add the L1 loss to add to the renderer, thus increasing the pixel-level quality of the rendered image. We use the null hypothesis in discrimination to compare the differences between real and generated images. Due to the different characteristics of male and female portraits, we train the network using all methods, but for testing we limit the search space to male and female positions.

## V. APPLICATIONS

Our system can be adapted for numerous applications. In this section, we present two applications: face morphing and face modeling.

### A. Face Morphing

Conventional face morphing algorithms usually require key-level similarity between two face shapes to guide semantic interpolation. We present a simple but effective morphing approach by 1) dividing a pair of face sketches from the training database into five parts; 2) encode the component sketches into feature vectors in their respective characteristic spaces.

### B. Face Modelling

The traditional copying technique uses a continuous stitching technique on colored images. However, there will be many cases where local color doesn't matter. To solve this problem, we combine surface components to compose new surfaces that can maintain all color and light consistency. Separately, the sketch of the first coding surface components as feature vectors (perhaps from different subjects) can also be combined as a new surface using the FM module and the IS module. This can be used to replace or replace existing surface parts with components from other sources, or to combine components from multiple sources. Figure 20 shows some new faces synthesized by recombination of components such as eyes, nose, mouth and the rest of the four source sketches. Since our images are synthesized, the mesh can resolve conflicts between facial components from different sources in terms of image and flash.

### C. Criminal Investigation:

Finding a suspect based on a sketch image is a difficult task. The police use a lot of automatic biometric engineering to identify suspects in some crimes, the only information that ensures that the investigator will build a forensic sketch

using a skilled artist and sketch to convert the sketch into a digital image and this system or framework. it is easy to find suspicious.

## VI.CONCLUSION AND DISCUSSIONS

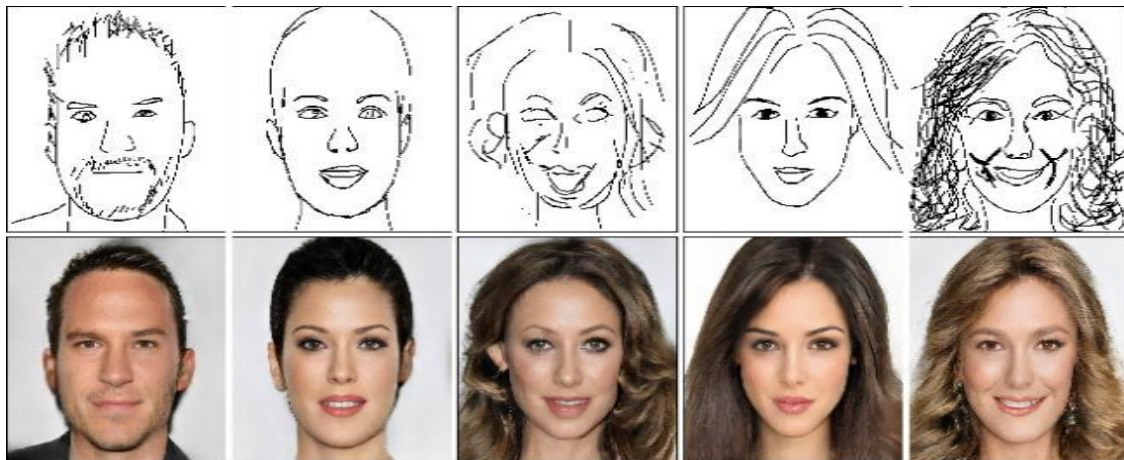
So, in this paper we have hand out a novel deep learning framework which synthesize realistic human face images from irregular or incomplete freehand sketch. We have taken both local-to-global approach by firstly de-composing a face sketch into separate components, clarifying its individual component by extruding them to components manifold defined from the existing components samples in the characteristics/feature spaces, plotting or mapping the filtered/refined characters / features vector to the features map for structural combination, and finally converting the combined features map to realistic images. As this approach normally support local editing and make the network involved easy to learn/ train from its training datasets which is not much in large scale.

So, our perspective outperforms the existing sketch to image synthesis approaches, often requires edge map or sketch with similar quality as inputs. Our user study has confirmed the versatility of our system. We have also adapted our system for the following applications as face morphing, face copy-paste and criminal investigation.

## VII.ACKNOWLEDGMENT

The authors would like to thank the publishers and researchers for making their resource available and also, we are thankful to our teachers for their guidance. We would like to express our gratitude to our guide Professor Sonali Dongare for his encouraging support and guidance in carrying out this work. We express our sincere thanks to Nutan Maharashtra Institute of Engineering and Technology Pune for permitting us to take our work this further.

## VIII.TESTED IMAGES



## REFERENCES

- [1] Volker Blanz and Thomas Vetter. 1999. A Morphable Model for the Synthesis of 3D Faces. In Proceedings of the 26<sup>th</sup> Annual Conference on Computer Graphics and Interactive Techniques. ACM, 187–194.
- [2] John Canny. 1986. A computational approach to edge detection. IEEE Transactions on Pattern Analysis and Machine Intelligence PAMI-8, 6 (1986), 679–698.
- [3] Tali Dekel, Chuhan Gan, Dilip Krishnan, Ce Liu, and William T Freeman. 2018. Sparse, smart contours to represent and edit images. In IEEE Conference on Computer Vision and Pattern Recognition (CVPR). 3511–3520.
- [4] Lin Gao, Jie Yang, Tong Wu, Yu-Jie Yuan, Hongbo Fu, Yu-Kun Lai, and Hao(Richard) Zhang. 2019. SDM-NET: Deep Generative Network for Structured Deformable Mesh. ACM Trans. Graph. 38, 6 (2019), 243:1–243:15.
- [5] Shiming Ge, Xin Jin, Qiting Ye, Zhao Luo, and Qiang Li. 2018. Image editing by object-aware optimal boundary searching and mixed-domain composition. Computational Visual Media 4 (01 2018). <https://doi.org/10.1007/s41095-017-0102-8>.

- [6] Xiaoguang Han, Chang Gao, and Yizhou Yu. 2017. DeepSketch2Face: a deep learning-based sketching system for 3D face and caricature modeling. *ACM Trans. Graph.* 36, 4, Article Article 126 (2017), 12 pages.
- [7] Takeo Igarashi, Satoshi Matsuoka, Sachiko Kawachiya, and Hidehiko Tanaka. 1997. Interactive Beautification: A Technique for Rapid Geometric Design. In *Proceedings of the 10th Annual ACM Symposium on User Interface Software and Technology (UIST '97)*. Association for Computing Machinery, 105–114.
- [8] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. 2017. Image-to-image translation with conditional adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 1125–1134.
- [9] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. 2016. Perceptual losses for real-time style transfer and super-resolution. In *European Conference on Computer Vision (ECCV)*. Springer-Verlag, 694–711.
- [10] Tero Karras, Samuli Laine, and Timo Aila. 2019. A style-based generator architecture for generative adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 4401–4410.
- [11] Diederik P. Kingma and Jimmy Ba. 2014. Adam: A Method for Stochastic Optimization. <http://arxiv.org/abs/1412.6980> Comment: Published as a conference paper at the 3rd International Conference for Learning Representations, San Diego, 2015.
- [12] Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. 2019. MaskGAN: Towards Diverse and Interactive Facial Image Manipulation. *arXiv preprint arXiv:1907.11922* (2019).
- [13] Yong Jae Lee, C Lawrence Zitnick, and Michael F Cohen. 2011. Shadowdraw: real-time user guidance for freehand drawing. *ACM Trans. Graph.* 30, 4, Article Article 27 (2011), 10 pages.
- [14] Yuhang Li, Xuejin Chen, Feng Wu, and Zheng-Jun Zha. 2019. LinesToFacePhoto: Face Photo Generation from Lines with Conditional Self-Attention Generative Adversarial Networks. In *Proceedings of the 27th ACM International Conference on Multimedia*. ACM, 2323–2331.
- [15] Sam T Roweis and Lawrence K Saul. 2000. Nonlinear dimensionality reduction by locally linear embedding. *Science* 290, 5500 (2000), 2323–2326.
- [16] Edgar Simo-Serra, Satoshi Iizuka, Kazuma Sasaki, and Hiroshi Ishikawa. 2016. Learning to Simplify: Fully Convolutional Networks for Rough Sketch Cleanup. *ACM Trans. Graph.* 35, 4, Article Article 121 (2016), 11 pages.
- [17] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. 2018. High-resolution image synthesis and semantic manipulation with conditional gans. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 8798–8807.
- [18] Xiaogang Wang and Xiaoou Tang. 2008. Face photo-sketch synthesis and recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 31, 11 (2008), 1955–1967.
- [19] Saining Xie and Zhuowen Tu. 2015. Holistically-Nested Edge Detection. In *IEEE International Conference on Computer Vision (ICCV)*. IEEE, 1395–1403.
- [20] Ran Yi, Yong-Jin Liu, Yu-Kun Lai, and Paul L Rosin. 2019. APDrawingGAN: Generating Artistic Portrait Drawings from Face Photos with Hierarchical GANs. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 10743–10752.
- [21] Wei Zhang, Xiaogang Wang, and Xiaoou Tang. 2011. Coupled information-theoretic encoding for face photo-sketch recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 513–520.
- [22] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. 2017. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *IEEE International Conference on Computer Vision (ICCV)*. 2223–2232.