

CONDITIONAL DIFFUSION INPAINTING FOR SKETCH-TO-FACE SYNTHESIS

Sanhita Pathak^{*1}, Vinay Kaushik², and Brejesh Lall¹

¹Indian Institute of Technology Delhi

²Indian Institute Of Information Technology Sonapat

Sanhita.Pathak@dbst.iitd.ac.in, vkaushik@iiitsonepat.ac.in, brejesh@ee.iitd.ac.in

Abstract

Generating realistic faces from monochromatic sketches is challenging due to missing details like expressions and skin tones. GANs struggle with stability and structure, while diffusion models face issues with monochrome inputs and high costs. DS-Face, a latent diffusion-based conditional inpainting framework, addresses these challenges using a frozen Paint-by-Example diffusion model with ControlNet conditioning and DINO-V2 embeddings from a GAN-generated coarse image. Trained on the CUFS dataset, DSFace achieves superior realism, perceptual quality, and structural alignment.

1 Introduction

Generating realistic human faces from monochrome sketches is a key challenge in image-to-image translation [7, 8], benefiting domains such as character design and forensic analysis. However, the narrow distribution of single-channel sketch data and the lack of semantic cues hinder robust feature extraction and generalization. Traditional GAN-based methods [9–11] leverage semantic masks or paired supervision but often lose fine details such as wrinkles, expressions, and accessories. Diffusion models (DM) [12] and CLIP [13] have shown promise in text-to-image synthesis, yet their direct adaptation to sketch-to-image translation remains limited due to weak structural correspondence. Conditioning-guided diffusion models like ILVR [14] and SDEdit [8] use RGB references for control but still fail to preserve accurate facial geometry.

We propose a latent diffusion-based framework trained on a sketch–face dataset that reformulates sketch-to-face synthesis as a conditional inpainting problem. By leveraging the conditioning flexibility of diffusion models, we employ coarse image guidance, A PBE based baseline and DINO-V2 embeddings to provide global structural consistency. Furthermore, ControlNet and parsed facial maps enable fine-grained control over local facial attributes. The proposed model demonstrates superior realism, structural alignment, and quantitative performance

on the CUFS dataset. Our contributions are as follows, 1) Reframe sketch-to-image generation as conditional inpainting using latent diffusion. 2) Employ DINO-V2 embeddings and GAN-generated coarse images for structural alignment. 3) Integrate ControlNet and semantic face parsing for fine-grained control.

2 Methodology

We propose solving sketch-to-face generation as a conditional diffusion inpainting task utilizing ControlNet encoder for structural control, and conditioning Dino-v2 embeddings from a coarse input generated by a GAN module.

The proposed methodology consists of two stages. In **Stage 1**, a Generative Adversarial Network (GAN) is employed to generate a coarse face image I_{coarse} from the input facial sketch I_s . The resulting coarse image I_{coarse} serves as a global structural prior, providing essential conditioning for the subsequent stage. The **Stage 2** consists of diffusion based pipeline, which consists of a frozen Diffusion Denoising Unet. The inputs to this module are face segmentation I_{seg} , face agnostic binary mask I_{mask} , face agnostic rgb image I_{ag} , input image with added noise I_t . The diffusion pipeline is guided by the structural conditioning from a ControlNet module for enabling the structural coherence with realistic facial features, the features guiding the controlNet are extracted from input sketch image I_s . An additional conditioning input that guides the ControlNet is the Dino-v2 embeddings computed from the provided coarse input image from stage 1, I_{coarse} . Dino-v2 module enables the pipeline to capture more detailed information about facial attributes, thereby enhancing control over the generation of facial regions. The ControlNet is trained to predict total noise, with the final clean generated face denoted as I_{out} .

3 Experiments

3.1 Implementation Details

We conduct end-to-end training using one NVIDIA A100 GPU with image resolution of 512×384 . Training runs for 150 epochs using the AdamW optimizer

^{*}Corresponding Author.

Metrics	DualGAN [1]	Pix2Pix [2]	UGATIT [3]	NICE-GAN [4]	FRAN [5]	DiSS [6]	DCNP [7]	Ours
↑SR-SIM	0.8687	0.8681	0.8758	0.8800	0.8822	0.8655	0.8787	0.8941
↑MSSIM	0.7938	0.7956	0.8090	0.8206	0.8381	0.8041	0.8181	0.8618
↓FID	126.53	85.55	117.34	84.96	73.21	86.22	65.49	62.58
↑FSIM	0.7893	0.7913	0.7955	0.8123	0.8187	0.7763	0.8121	0.8352
↑VIF	0.1128	0.1098	0.1367	0.1427	0.1798	0.1253	0.1475	0.2321
↑IS	1.32	1.35	1.42	1.41	1.35	1.34	1.40	1.49

Table 1. Comparison of various models using different metrics on CUSK datasets. Bold values indicate the best performance.

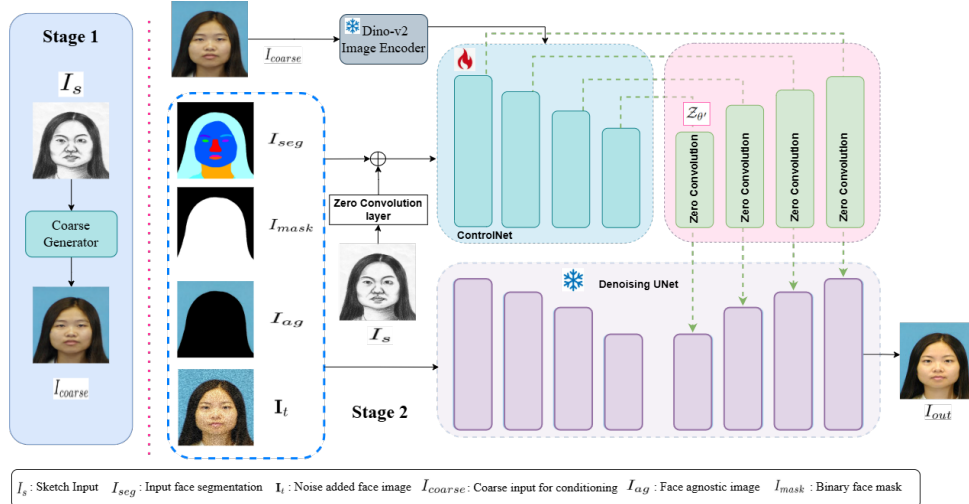


Figure 1. Architecture diagram of our proposed approach for sketch to face image synthesis



Figure 2. Qualitative comparison of our proposed approach with state-of-the-art approaches for sketch to face image synthesis.



Figure 3. Qualitative Ablation of effect of various proposed modules on quality of generated faces.

085 with a learning rate of 2×10^{-5} . Experiments are
 086 performed on the CUFS dataset [15], which includes
 087 606 identities from CUHK, AR, and XM2VTS, each
 088 containing one FPS pair, split in an 80:20 ratio for
 089 training and testing.

090 3.2 Results

091 We evaluate the proposed model for sketch-to-face
 092 synthesis through qualitative and quantitative compar-
 093 isons with state-of-the-art methods, including
 094 DualGAN [1], Pix2Pix [2], UGATIT [3], NICE-
 095 GAN [4], FRAN [5], DiSS [6], and DCNP [7]. As
 096 shown in Figure 2, earlier methods often yield blurry
 097 results or artifacts, while our approach produces
 098 sharper, more realistic images with well-preserved

facial features and textures, benefitting from DINO-
 v2 and ControlNet conditioning.

Quantitatively, Table 1 highlights our approach’s
 superior performance, achieving the highest SR-
 SIM (0.8941), MSSIM (0.8618), FSIM (0.8352), VIF
 (0.2321), and IS (1.49), along with the lowest FID
 (62.58). These results confirm that our approach del-
 ivers enhanced realism, structural consistency, and
 perceptual fidelity compared to existing approaches.

3.3 Ablation Study

We evaluate our proposed method through quanti-
 tative and qualitative ablations (Figure 3). The
 baseline PBE model, conditioned only on sketch
 input, fails to preserve structural and perceptual
 fidelity, yielding poor FID, MSSSIM, and IS scores.
 Adding a ControlNet encoder enhances facial geom-
 etry consistency, while incorporating face segmen-
 tation improves perceptual quality through semantic
 guidance, though texture and color remain inconsis-
 tent. Finally, conditioning ControlNet with DINO-
 v2 embeddings from the coarse generated image
 provides refined structural and appearance details.

4 Conclusion

We propose a latent diffusion model that gener-
 ates realistic faces from sketches using conditional
 inpainting with semantic guidance. It achieves state-
 of-the-art realism and fidelity on the CUFS dataset,
 advancing sketch-based synthesis applications.

127 **References**

- 128 [1] Z. Yi, H. Zhang, P. Tan, and M. Gong. “Dua-
129 lgan: Unsupervised dual learning for image-
130 to-image translation”. In: *Proceedings of the*
131 *IEEE international conference on computer*
132 *vision*. 2017, pp. 2849–2857.
- 133 [2] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros.
134 “Unpaired image-to-image translation using
135 cycle-consistent adversarial networks”. In: *Pro-*
136 *ceedings of the IEEE international conference*
137 *on computer vision*. 2017, pp. 2223–2232.
- 138 [3] J. Kim. “U-gat-it: unsupervised genera-
139 tive attentional networks with adaptive
140 layer-instance normalization for image-
141 to-image translation”. In: *arXiv preprint*
142 *arXiv:1907.10830* (2019).
- 143 [4] R. Chen, W. Huang, B. Huang, F. Sun, and
144 B. Fang. “Reusing discriminators for encoding:
145 Towards unsupervised image-to-image transla-
146 tion”. In: *Proceedings of the IEEE/CVF con-*
147 *ference on computer vision and pattern recog-*
148 *nition*. 2020, pp. 8168–8177.
- 149 [5] W. Wan, Y. Yang, S. Huang, and L. Gan.
150 “FRAN: feature-filtered residual attention net-
151 work for realistic face sketch-to-photo transfor-
152 mation”. In: *Applied Intelligence* 53.12 (2023),
153 pp. 15946–15956.
- 154 [6] S.-I. Cheng, Y.-J. Chen, W.-C. Chiu, H.-Y.
155 Tseng, and H.-Y. Lee. “Adaptively-realistic
156 image generation from stroke and sketch
157 with diffusion model”. In: *Proceedings of the*
158 *IEEE/CVF winter conference on applications*
159 *of computer vision*. 2023, pp. 4054–4062.
- 160 [7] M. Zhu, Z. Wu, N. Wang, H. Yang, and X.
161 Gao. “Dual conditional normalization pyramid
162 network for face photo-sketch synthesis”. In:
163 *IEEE Transactions on Circuits and Systems*
164 *for Video Technology* 33.9 (2023), pp. 5200–
165 5211.
- 166 [8] C. Meng, Y. Song, J. Song, J. Wu, J.-Y. Zhu,
167 and S. Ermon. “Sdedit: Image synthesis and
168 editing with stochastic differential equations”.
169 In: *arXiv preprint arXiv:2108.01073* (2021).
- 170 [9] H. Caesar, J. Uijlings, and V. Ferrari. “Coco-
171 stuff: Thing and stuff classes in context”. In:
172 *Proceedings of the IEEE conference on com-*
173 *puter vision and pattern recognition*. 2018,
174 pp. 1209–1218.
- 175 [10] S.-Y. Chen, W. Su, L. Gao, S. Xia, and H. Fu.
176 “Deep generation of face images from sketches”.
177 In: *arXiv preprint arXiv:2006.01047* (2020).
- [11] C.-H. Lee, Z. Liu, L. Wu, and P. Luo. 178
“Maskgan: Towards diverse and interactive fa- 179
cial image manipulation”. In: *Proceedings of* 180
the IEEE/CVF conference on computer vision 181
and pattern recognition. 2020, pp. 5549–5558. 182
- [12] J. Ho, A. Jain, and P. Abbeel. “Denoising 183
diffusion probabilistic models”. In: *Advances* 184
in neural information processing systems 33 185
(2020), pp. 6840–6851. 186
- [13] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, 187
G. Goh, S. Agarwal, G. Sastry, A. Askell, P. 188
Mishkin, J. Clark, et al. “Learning transferable 189
visual models from natural language supervi- 190
sion”. In: *International conference on machine* 191
learning. PMLR. 2021, pp. 8748–8763. 192
- [14] J. Choi, S. Kim, Y. Jeong, Y. Gwon, and 193
S. Yoon. “Ilvr: Conditioning method for de- 194
noising diffusion probabilistic models. In 2021 195
IEEE”. In: *CVF international conference on* 196
computer vision (ICCV). Vol. 1. 2021, p. 2. 197
- [15] X. Wang and X. Tang. “Face photo-sketch 198
synthesis and recognition”. In: *IEEE transac-* 199
tions on pattern analysis and machine intelli- 200
gence 31.11 (2008), pp. 1955–1967. 201