Generating Indian Facial Images from Textual Descriptions Using Diffusion Models - ALiterature Review

Chetankumar B Parmar

Research Scholar

Assistant Professor

Narmada College of Science & Commerce, Bharuch- Gujarat, India Veer Narmad South Gujarat University (Computer Science Dept.), Surat

Dr.Ravi Gulati, PhD

Professor,

Department of Computer Science,

Veer Narmad South Gujarat University, Surat, Gujarat, India.



Abstract:

This literature review surveys methods for generating facial images from natural-language

descriptions with an emphasis on Indian faces and the recent rise of diffusion-based models.

We trace the evolution from early GAN-based text-to-image and text-to-face systems to

modern latent diffusion approaches, compare their strengths and weaknesses, summarize

available Indian-face resources, and identify open problems and ethical considerations. The

review ends with recommended experimental setups and a proposed benchmark for evaluating

text-to-face systems in the Indian context.

Introduction:

The synthesis of human facial images from textual descriptions has rapidly progressed over the

past decade, driven by advances in Generative Adversarial Networks (GANs) [24],

Variational Autoencoders (VAEs), and more recently, diffusion-based models [21]. While

models such as Stable Diffusion and Imagen have achieved remarkable performance in

general domains [9][10], their effectiveness in representing ethnically diverse populations,

particularly Indian faces, remains underexplored. This gap is significant because biased

generative systems risk producing inaccurate or stereotyped outputs, especially when trained

primarily on Western-centric datasets [12][13][14][15][16][18]. The objective of this review is

to evaluate prior work on generative facial synthesis, identify methodological innovations,

highlight dataset limitations, and propose future directions for building culturally relevant

benchmarks for Indian faces.

Related Work

Early approaches to text-to-image synthesis were dominated by GANs. Works such as

VQGAN-CLIP enabled fine-grained generation by linking textual prompts with visual

embeddings [3]. For forensic applications, GAN-based frameworks were used to reconstruct

human faces from witness descriptions [2]. More recent extensions incorporated attention

mechanisms and multimodal conditioning to enhance control and alignment with textual

descriptions [6].

The emergence of diffusion models marked a paradigm shift. Latent Diffusion Models

(LDMs) [21] and Imagen [10] achieved state-of-the-art realism and semantic alignment, while

methods like **DreamBooth** [6] enabled subject-specific customization. Hybrid approaches

such as **GANDiffFace** further bridged GAN identity generation with diffusion-based variation,

offering robustness across identity categories, including Indian faces [7].

VNSGU Journal of Research and Innovation (Peer Reviewed)

109

Literature Review:

Sr. No	Title Of the Paper	Names of Authors	Publications	Findings
1	StackGAN: Text to photo-realistic image synthesis with stacked generative adversarial networks	Zhang, H., Xu, T., Li, H., Zhang, S., Wang, X., Huang, X., & Metaxas, D. N.	Proceedings of the IEEE International Conference on Computer Vision (ICCV), 2017	Introduced a stacked GAN approach to generate high-resolution, photorealistic images from text descriptions.
2	StackGAN++: Realistic image synthesis with stacked generative adversarial networks	Zhang, H., Xu, T., Li, H., Zhang, S., Wang, X., Huang, X., & Metaxas, D. N.	IEEE Transactions on Pattern Analysis and Machine Intelligence, 2019	Improved upon StackGAN with a more stable architecture and advanced training for better realism.
3	VQGAN-CLIP: Open domain image generation and editing with natural language guidance	Crowson, K., Mostaque, E., & Robinson, D.	arXiv, 2022	Combined VQGAN and CLIP to enable flexible open- domain image generation and editing from text prompts.
4	TediGAN: Text- guided diverse image generation and manipulation	Xia, W., Yang, Y., Xue, JH., & Wu, B.	Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2021	Enabled fine- grained and diverse image generation and editing guided by text descriptions.
5	StyleGAN2	NVIDIA	Wikipedia, 2020	Second version of StyleGAN, improved quality and reduced artifacts in generative image synthesis.
6	clip2latent: Text driven sampling of a pre-trained StyleGAN using denoising diffusion and CLIP	Pinkney, J., & Adler, D.	GitHub, 2022	Integrated CLIP and diffusion with StyleGAN for controlled sampling from textual descriptions.

7	DiffusionCLIP: Text-	Kim, G.,	Proceedings of the	Introduced text-
	guided diffusion	Nam, S., &	IEEE/CVF	conditioned
	models for robust	Chun, S.	Conference on	diffusion for
	image manipulation	ŕ	Computer Vision	image editing
			and Pattern	with robustness
			Recognition	and quality
			(CVPR), 2022	improvements.
8	DATID-3D: Diversity-	Kim, G., &	Proceedings of the	Presented
	preserved domain	Chun, S.	IEEE/CVF	diffusion-based
	adaptation using text-	Chun, B.	Conference on	methods for
	to-image diffusion for		Computer Vision	diversity-
	3D generative model		and Pattern	preserving 3D
	3D generative model			model generation
			Recognition	with text
			(CVPR), 2023	
	A ('C' : 1 : 4 11'	MDDI	MDDI 2024	guidance.
9	Artificial-intelligence-	MDPI	MDPI, 2024	Surveyed
	generated content with			literature on AI-
	diffusion models: A	. 4 11 -	- VO.N	generated content
	literature review	1444		using diffusion
	//ES/ A	B		models with
	//65/ . 🕨		MAY A	emphasis on
				applications.
10	Diffusion models: A	ACM	ACM Digital	Comprehensive
	comprehensive survey		Library, 2023	survey of
	of methods and	$\mathcal{M}(G)$		diffusion models
	applications	MILLO		covering
				methodologies
				and practical
	11111 - 3			applications.
11	Fair Diffusion: Bias	Naik, S.,	Springer	Addressed bias
	mitigation in Stable	&Nushi, B.	Professional, 2023	mitigation
	Diffusion			strategies in
				diffusion models
	1 33		13.77	to promote
	1.0	1111	- (A) 5	fairness.
12	Interpretations,	Ghosh, S.	arXiv, 2024	Analyzed caste-
	representations, and		, - , - , - , - , - , - , - , - , - , -	related
	stereotypes of caste			stereotypes
	within text-to-image			embedded in
	generators			generative text-
	5011010110115			to-image
				systems.
13	Do generative AI	Ghosh, S.,	arXiv, 2024	Studied cultural
13	models output harm		ai/Xiv, 2024	harms and
	=	Roy, P., &		
	while representing non-Western cultures:	Sharma, R.		misrepresentation of non-Western
	Evidence from a			contexts in
	community-centered			generative AI.
	approach			

14	Inspecting the geographical representativeness of images from text-to-image models	Basu, A., Das, M., & Roy, S.	arXiv, 2023	Evaluated the geographic balance and inclusivity of datasets in T2I
				generative models.
15	AI's regimes of representation: A community-centered study of text-to-image models in South Asia	Anonymous	ResearchGate/ar5iv, 2023	Explored representational politics and biases of AI models within South Asian contexts.
16	Digital Orientalism in machine vision: A cross-platform analysis of AI-generated representations of Indian culture	Qadri, A., & Ali, S.	Cyber Lenika, 2023	Highlighted orientalist biases in AI visual representation of Indian culture across platforms.

Methodological Background

Generative methods for text-to-face synthesis can be grouped into three paradigms:

- 1. **GAN-based approaches:** Learn adversarial mappings between noise and image distributions [24]. They provide effective attribute-level control but struggle with semantic fidelity.
- 2. **VAE-based approaches:** Model latent distributions and enable interpolation, but often at the cost of image sharpness.
- 3. **Diffusion-based approaches:** Model data distributions through iterative denoising steps, yielding high-quality, semantically rich outputs [21]. Their scalability (e.g., **Stable Diffusion, Imagen**) has made them the dominant methodology since 2021 [10].

Diffusion's integration with **CLIP embeddings** [22] and **cross-attention** allows nuanced conditioning on natural language descriptions, a critical factor for descriptive synthesis of Indian faces.

Datasets

High-quality datasets remain a bottleneck in ensuring fair representation.

- **IIITM Face Dataset:** Captures Indian subjects across varied poses and emotions [25].
- Indian Masked Faces in the Wild (IMFW): Documents Indian cultural masking practices (e.g., gamcha, stoles) [25].
- **FairFace:** Includes Indian identities but in limited representation, used in hybrid frameworks like **GANDiffFace**[7].
- Mukh-Oboyob: A Bangla text-to-face dataset leveraging Stable Diffusion and BanglaBERT [20].

While these datasets improve cultural inclusivity, none provide **text-image paired descriptions**, a crucial gap for training text-to-face synthesis models [19][26].

Comparative Analysis

Model	Strengths	Limitations in Indian Context
GANs (e.g., VQGAN-	Fine-grained text-to-image	Struggles with realism; bias in Indian
CLIP) [3]	control	face representation
Multi-GAN with	Stronger attribute alignment	Relies on general datasets (FFHQ,
attention [6]		CelebA), lacking Indian faces
Latent Diffusion /	High photorealism, strong text	Dataset bias leads to
Imagen [10][21]	alignment	underrepresentation of Indian features
DreamBooth [6]	Personalized face generation	Requires curated reference images for
1 6	with few-shot learning	Indian subjects
GANDiffFace [7]	Combines GAN identity	Includes Indian category via FairFace
N.	control + diffusion variation	but not text-driven

Experimental Design & Proposed Benchmark

To address representation gaps, we propose a **benchmark pipeline**:

- 1. **Dataset Curation:** Construct an **Indian text-to-face dataset** with descriptive captions reflecting skin tones, attire, and cultural attributes [12][13][14][15][16].
- 2. **Baseline Comparison:** Train and evaluate GANs [1][2][3], VAEs, and diffusion-based models (Stable Diffusion, Imagen, DreamBooth) [6][9][10].
- 3. Evaluation Metrics:
 - o **FID/IS** for realism [9].
 - o **CLIP-Score** for text-image alignment [22].

Fairness Metrics: Recognition accuracy across Indian subgroups [12][18].

4. Benchmark Standardization: Release as an open benchmark to facilitate

reproducibility and comparative evaluation in Indian facial synthesis [19][20].

Ethical Considerations

Generative models for facial synthesis raise significant ethical challenges:

• Bias and Fairness: Current systems risk perpetuating stereotypes if Indian facial

diversity is not adequately represented [12][13][14][16][17][18].

• **Privacy:** Training on Indian datasets must ensure consent and anonymization to prevent

misuse [25].

Misuse in Deepfakes: Systems could be weaponized for misinformation or identity

manipulation [17]. A framework for responsible use, including watermarking and

generative provenance, is essential.

Community Involvement: A participatory approach involving Indian communities in

dataset curation ensures cultural sensitivity [15][16].

Conclusion

This review shows that while GAN-based models [1][2][3] provided foundational insights,

the state-of-the-art in diffusion-based synthesis [9][10][21] enables unprecedented realism

and semantic fidelity. However, Indian faces remain underrepresented both in datasets [25] and

benchmarks [19][20]. Addressing this gap requires curating Indian text-image paired datasets,

adapting diffusion frameworks, and designing fairness-driven benchmarks [12][13][14]. Such

efforts will ensure that generative AI not only advances technically but also inclusively

represents India's diverse population.

References

1. Zhang, H., Xu, T., Li, H., Zhang, S., Wang, X., Huang, X., & Metaxas, D. N. (2017).

StackGAN: Text to photo-realistic image synthesis with stacked generative adversarial

networks. In Proceedings of the IEEE International Conference on Computer Vision

(ICCV) (pp. 5908–5916). IEEE. https://doi.org/10.1109/ICCV.2017.629 (Penn State)

VNSGU Journal of Research and Innovation (Peer Reviewed)

114

- 2. Zhang, H., Xu, T., Li, H., Zhang, S., Wang, X., Huang, X., & Metaxas, D. N. (2019). StackGAN++: Realistic image synthesis with stacked generative adversarial networks. IEEE Transactions on Pattern Analysis and Machine Intelligence, 41(8), 1947–1962.
- 3. Crowson, K., Mostaque, E., & Robinson, D. (2022). VQGAN-CLIP: Open domain image generation and editing with natural language guidance. arXiv. arXiv:2204.08583.
- 4. Xia, W., Yang, Y., Xue, J.-H., & Wu, B. (2021). TediGAN: Text-guided diverse image generation and manipulation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (pp. 2229–2238). IEEE. https://doi.org/10.1109/CVPR46437.2021.00229 (ResearchGate)
- 5. NVIDIA. (2020). StyleGAN2. In Wikipedia. Retrieved, from https://en.wikipedia.org/wiki/StyleGAN2
- 6. Pinkney, J., & Adler, D. (2022). clip2latent: Text driven sampling of a pre-trained StyleGAN using denoising diffusion and CLIP.
- 7. Kim, G., Nam, S., & Chun, S. (2022). DiffusionCLIP: Text-guided diffusion models for robust image manipulation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2426–2435.
- 8. Kim, G., & Chun, S. (2023). DATID-3D: Diversity-preserved domain adaptation using text-to-image diffusion for 3D generative model. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).
- 9. MDPI. (2024). Artificial-intelligence-generated content with diffusion models: A literature review. MDPI.
- 10. ACM. (2023). Diffusion models: A comprehensive survey of methods and applications. ACM Digital Library.
- 11. Naik, S., &Nushi, B. (2023). Fair Diffusion: Bias mitigation in Stable Diffusion. In Springer Professional.
- 12. Ghosh, S. (2024). Interpretations, representations, and stereotypes of caste within text-to-image generators.
- 13. Ghosh, S., Roy, P., & Sharma, R. (2024). Do generative AI models output harm while representing non-Western cultures: Evidence from a community-centeredapproach.
- 14. Basu, A., Das, M., & Roy, S. (2023). Inspecting the geographical representativeness of images from text-to-image models.
- 15. (2023). AI's regimes of representation: A community-centered study of text-to-image models in South Asia.

- 16. Qadri, A., & Ali, S. (2023). Digital Orientalism in machine vision: A cross-platform analysis of AI-generated representations of Indian culture.
- 17. Luccioni, A., Mitchell, M., &Bengio, Y. (2023). Analyzing societal representations in diffusion models.
- 18. (2024). Navigating text-to-image generative bias across Indic languages.
- 19. Parihar, R., Singh, A., & Jain, V. (2024). PreciseControl: Enhancing text-to-image diffusion models with fine-grained attribute control. In Proceedings of the European Conference on Computer Vision (ECCV).
- 20. Dai, D., Zhang, X., & Wang, Y. (2024). 15M multimodal facial image-text dataset.
- 21. Saha, A., Rahman, M., & Hossain, M. (2023). Mukh-Oboyob: Stable Diffusion and BanglaBERT enhanced Bangla text-to-face synthesis. International Journal of Advanced Computer Science and Applications (IJACSA), 14(3), 50–59.
- 22. Wikipedia. (2023). Diffusion model. Retrieved, from https://en.wikipedia.org/wiki/Diffusion_model
- 23. Wikipedia. (2023). Contrastive language—image pre-training (CLIP). Retrieved, from https://en.wikipedia.org/wiki/Contrastive_language%E2%80%93image_pre-training
- 24. Wikipedia. (2023). Generative adversarial network. Retrieved, from https://en.wikipedia.org/wiki/Generative_adversarial_network
- 25. Wikipedia. (2023). List of facial expression databases. Retrieved, from https://en.wikipedia.org/wiki/List_of_facial_expression_databases
- 26. Sun, X., Li, J., & Wang, Y. (2021). Multi-caption datasets.
- 27. (2023). Face generation from long paragraphs using ParaDiffusion.
- 28. (2023). Societal bias analysis of text-to-image (Luccioni)
- 29. (2022). VQGAN-CLIP usage tutorial.
- 30. (2022). Simpler prompt-based VQGAN+CLIP generation tutorial.