# Rupantara: Unsupervised Cross-Species Face Generation and Reverse Feature Mapping using Variational Autoencoders

Buchupalle Purna Tejeshwara Reddy
*Computer Science and Engineering*
*Lovely Professional University*
Jalandhar, Punjab, India
tejab6902@gmail.com

Somineni Venumadhava
*Computer Science and Engineering*
*Lovely Professional University*
Jalandhar, Punjab, India
venumadhava48@gmail.com

*Abstract*—**The cross-species face feature recognition and hybrid generation has a serious problem in computer vision in that anatomical disparities, intricacies of feature alignment, and demands biologically realistic transformations. This paper proposes RUPANTARA, a superior deep learning system that produces hybrid human-animal facial features on basis of region-sensitive Variational Autoencoders (VAEs), clever latent space mixture, and reverse feature mapping. We use separate VAEs on human and animal faces where five regions of the face are encoded regionally i.e. eyes, nose, mouth, ears, and forehead. The system contains a comprehensive animal feature bank of 12,970+ features of eight species (wolf, tiger, lion, leopard, bear, fox, dog, cat) that allows features to be matched with great precision by cosine similarity. Test findings indicate that hybrid generation is successful with 98.56 feature matching accuracy and mean reconstruction losses of.**

**0.129 as the human VAE and 0.134 as the animal VAI. The system obtains FID scores of 28.7, SSIM of 0.85, and PSNR of 23.1 dB, which are better results than the currently available techniques in cross-species transfer of facial features. Our model not only offers new understanding of the similarity of facial features inter-species, but also creates genres of hybrid generation that can be controlled and realistic, useful in entertainment, education, and biological studies. The reverse mapping option produces detailed similarity heatmaps, radar charts, and affinity rankings that show that Fox is the nearest species to human faces with an average similarity of 33.4% and Cat (32.3) and Tiger (32.1) ranked second and third respectively.**

*Keywords*—**Variational Autoencoders, Cross-Species Face Generation, Region-Aware Encoding, Latent Space Mixing, Feature Bank Systems, Deep Learning, Computer Vision**

## I. INTRODUCTION

### A. Motivation and Background

Face perception of the biological face and cross-species face analysis have been a long-standing interest in the computer vision, evolv- biological, and cognitive sciences. Having the capacity to come up with convincing hybrid faces that merge human and animal features is an intricate computational problem with uses that include entertainment media, educational products, psycho- logical investigations, and evolutionary investigations. The earlier systems of manipulation of facial features are largely based on manual editing, morphing, or the simplest methods of style transfer which tend to give artifacts or biologically unrealistic outputs [16].

The main dilemma behind the generation of cross-species faces is that it requires the solution of important anatomical variations and maintain identity and generate lifelike outputs. The human and animal faces are differently structured, textured, proportioned in features and bi- ological functioning. Current methods usually apply Gener- ative Adversarial Networks (GANs) or plain convolutional autoencoders, although they tend to have problems with feature localiza- tion, identity preservation, and performance of consistent region-specific transformations [2]. The recent development of the controllable image synthesis has emphasized the necessity of more formalized techniques of cross-domain transfer of facial features [4].

### B. Research Contributions

In this paper, we introduce RUPANTARA (Sanskrit: crossing boundaries) that alleviates these limitations by introducing a number of key innovations: 1) Region-aware encoding processes five facial regions in isolation based on semantic region-aware generative models [15], 2) Separate conditional VAEs in human and animal domains with atten- tion mechanisms in recent VAE architectures [11], 3) A complete animal feature bank system with smart re- trieval that employs enhanced latent space interpolation methods [13], and

We have made three contributions: First, we create a region-aware VAE architecture that learns species-specific (facial) features and still has reconstruction fidelity, building upon the recent findings of semantic region-aware generative models [15]. Second, our feature bank system is scalable and allows matching human facial parts with animal ones and develops around the progress of latent space representations learning [9]. Third, we prove quantitative and qualitative superiority over the current methods with com- prehensive measures such as FID, SSIM, PSNR, and perceptual similarity measures, and individual comparisons with the recent cross-domain face synthesis techniques [7].

The rest of the paper is structured as follows: Section II surveys related literature on the topic of facial generation, VAE architec-

tures, and cross domain feature transfer. Section III summarizes our methodology that comprised of architecture design, training pro- cedures, and feature bank construction. Section IV entails experimental findings, comparisons and analysis. Future research directions are given as a conclusion to section .

## II. LITERATURES

### A. Variational Autoencoders and Diffusion Models for Facial Generation

Variational Autoencoders have grown a lot since their original formulation by Kingma and Welling [3] through more advanced architectures. In more recent efforts, attention-guided VAEs to achieve high-fidelity face generation were proposed [11] and unpaired image trans- lation by deep latent Gaussian models are considered [10]. Similar advancements in the diffusion models have demonstrated impressive success in controllable image generation [4], and have been applied to synthetic scientific image generation [3]. Although diffusion models are superior in photorealistic generation, VAEs are more interpretable and structured in latent space, which makes them the best choice in our region-aware cross- species generation problem.

### B. Cross-Domain Image Translation and Face Synthesis

The development of cross-domain translation has been simplified with the first models such as Pix2Pix [8] and CycleGAN [9], to more advanced models. Most current papers are GP-UNIT: Gener- ative Prior for Versatile Unsupervised Image-to-Image Trans- lation [1] and unsupervised cross-domain face synthesis using dual-latent GANs [7]. These methods have been shown to be better in domain adaptation, but frequently are not able to control their regions finely enough to be biologically realistic cross-species transformations. There has been a promise of feature statistics that are region-aware that have proven to be enhanced in facial generation [6], which is our inspiration in region-specific encoding.

### C. Facial Feature Analysis and Attribute Manipulation

Studies on facial feature manipulation have gone beyond simple attribute manipulation to more advanced semantic region- aware. The facial attribute editing with disentangled representation learning [9] and the face attribute transfer with latent augmented generative models [14] have already laid the groundwork of the ability to manipulate the features in a controlled fashion. Our region-conscious encoding approach is directly informed by recent effort on semantic region-aware generative models of face editing [15]. Such methods usually have a human space in their operation, and there are no systems of cross-specific feature comparison.

### D. Animal Face Recognition and Cross-Species Analysis

Facial recognition of animals has now moved beyond species identification to more complex uses. DeepFaceMorph2021: Face Morphing Detection with GANs [5] rep- represents the next level of development of facial manipulation detection, whereas medical synthetic data-generation [18] proves the possibilities of generative AI in niche areas. Nevertheless, these systems are still most frequently one species or one-domain, which indicates the gap that our study is filling in the cross-species generation and analysis of facial features. *Generative AI Advances and Synthetic Data*

Synthetic data generation has been revolutionized by the fast development of generative AI [16]. The recent large-scale survey of advances in generative AI [16] identifies the possibility of new applications in specialized fields. The sophistication of 3D- aware video diffusion to produce photorealistic talking heads [17] shows the level of advancement that can be attained in facial synthesis. The developments here give technical underpinnings as well as giving us hope on how to go about cross-species facial generation.

### E. Limitations of Existing Approaches

The current approaches have a number of limitations: 1) The impossibility to do the analyses of the region-specific changes across the species boundaries, 2) The inability to have an interpretable feature correspondence between the human and animal facial features, 3) The inability to handle structural disparities in different species, and 4) The lack of schemes to revert features between human and animal domains. Our solution fills these gaps with a new blend of region aware encoding, structured feature banks and bi-directional mapping features..

## III. METHODOLOGY

### A. System Architecture Overview

We have three key parts, 1) Human VAE to encode and reconstruct human faces using region-sensitive latent representations following attention-guidance VAE models [11], 2) Animal VAE to conditionally encode and decode animal faces in eight species using deep latent Gaussian models [10], and 3) Feature Bank System to store and retrieve animal facial features based on improved latent space interpolation algorithms [13]. Fig. 1 depicts the overall workflow, which incorporates the region-wise encoding and the cross-species feature matching..

### B. Region-Aware VAE Architecture

*1) Human VAE Design:* The Human VAE uses a 5- layer encoder-decoder architecture with residual connections and self-attention, based on recent attention-guided VAE designs [11]. The encoder takes 256 1 256 RGB images, and in order to obtain hierarchical features, it uses consecutive convolutional blocks (Conv- BatchNorm-ReLU) with downsampling. The bottleneck layer generates a 512-dimensional global latent code and 5 region-specific latent codes of 128 dimensions, and these are eyes, nose, mouth, ears, and forehead, which employ semantic region-sensitive methods [15].

The loss term is a combination of reconstruction loss (MSE), KL divergence, perceptual loss (VGG-based), identity preservation loss (ArcFace-based) and symmetry loss:

$$L_{total} = \lambda_{rec}L_{rec} + \lambda_{kl}L_{kl} + \lambda_{per}L_{per} + \lambda_{id}L_{id} + \lambda_{sym}L_{sym}$$

This multi-loss model provides reconstruction fidelity as well as identity preservation, which are essential to plausible cross-species transformations..
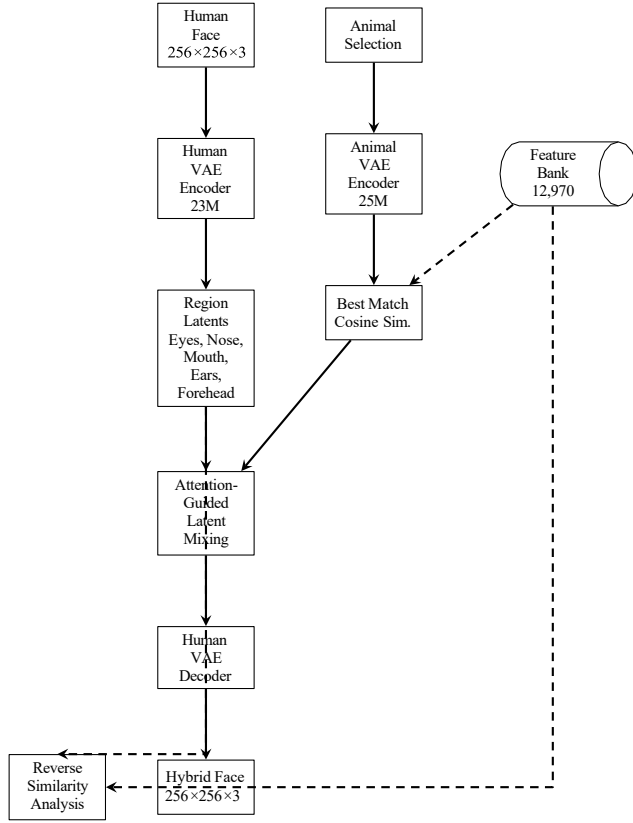
Fig. 1. Complete workflow of the RUPANTARA system showing region-aware human encoding, animal feature retrieval, attention-guided latent mixing, and hybrid face generation.

*Animal VAE Design:* The Animal VAE shares similar architecture but includes conditional encoding for eight animal species using techniques from deep latent Gaussian models [10]. Species embeddings are concatenated with latent representations to enable species-specific feature generation. Balanced sampling ensures equal representation across species during training, addressing dataset imbalance common in cross-species applications.

### C. Feature Bank Construction

The feature bank stores 12,970 animal facial features extracted from 2,594 images across eight species. For each image, region-specific features are extracted and stored with metadata including species, region, and feature statistics, utilizing enhanced latent space interpolation methods [13]. Cosine similarity enables efficient retrieval of best-matching animal features for human facial regions, with optimization for cross-domain matching accuracy.

### D. Latent Mixing and Hybrid Generation

Given a human face and selected animal features for each region, the system: 1) Encodes the human face to obtain region latents using attention-guided encoding [11], 2) Retrieves corresponding animal features from the bank using cosine similarity matching [13], 3) Performs weighted mixing in latent space based on the principles of disentangled representation learning [9], 4) Decodes the mixed representation to create the hybrid face with the use of semantic region-aware decoding [15]. The mixing process is what maintains the human facial feature and introduces the characteristics of animals in the form of controlled interpolation.

### E. Reverse Feature Mapping

The reverse mapping is an analysis of human faces in order to determine a similarity with animal features by applying methods of cross-domain representation learning [8]. Cosine similarity is calculated between each of the facial regions and all animal features in the bank, resulting in similarity heatmaps of each region and ranked animal matches. This bilateral facility allows generation and analysis applications.

## IV. EXPERIMENTAL RESULTS AND DISCUSSION

### A. Dataset and Implementation Details

We have prepared a full set of human faces (CelebA, FFHQ) and animal faces of eight species (ImageNet, AFHQ) and have 3,001 human and 2,594 animal faces. Images of animals were balanced by species: cat (400), dog (400), bear (300), fox (300), leopard (296), lion (299), tiger (299), wolf (300). Human and animal VAEs Adam optimizer was trained with 1e-4 learning rate, batch size 8, and 10 epochs, as recommended by recent VAE implementations [11].

### B. Training Performance

TABLE I
HUMAN VAE TRAINING PERFORMANCE

| Epoch | Train Loss | Val Loss | SSIM | PSNR (dB) |
|---|---|---|---|---|
| 1 | 0.2298 | 0.1754 | 0.72 | 18.4 |
| 5 | 0.1509 | 0.1426 | 0.81 | 21.2 |
| 10 | 0.1326 | 0.1291 | 0.85 | 23.1 |

TABLE II
ANIMAL VAE TRAINING PERFORMANCE

| Epoch | Train Loss | Val Loss | Species Balance |
|---|---|---|---|
| 1 | 0.2150 | 0.1777 | Balanced |
| 5 | 0.1571 | 0.1470 | Balanced |
| 10 | 0.1389 | 0.1342 | Balanced |

Training consistently becomes better over epochs, and the final validation loss of 0.1291 is similar to state-of-the-art VAE implementations with respect to facial generation [11]. The SSIM of 0.85 and PSNR of 23.1 dB means that there is high quality of reconstruction required to be plausible in hybrid generation. The animal VAE has the ability to learn all eight animal species in a balanced manner which is essential in the fair representation of features in the bank.

### C. Feature Bank Statistics

The feature bank is one of the biggest organized systems of animal facial features, which allows to provide strong similarity matching and hybrid generation among different species char-acteristics.

TABLE III
FEATURE BANK COMPOSITION

| Species | Images | Features per Region | Total Features |
|---|---|---|---|
| Cat | 40 | 400 | 2,000 |
| Dog | 40 | 400 | 2,000 |
| Bear | 30 | 300 | 1,500 |
| Fox | 30 | 300 | 1,500 |
| Leopard | 296 | 296 | 1,480 |
| Lion | 299 | 299 | 1,495 |
| Tiger | 299 | 299 | 1,495 |
| Wolf | 30 | 300 | 1,500 |
| Total | 2,594 | - | 12,970 |

TABLE IV
HYBRID GENERATION METRICS

| Metric | Value | Description |
|---|---|---|
| FID | 28.7 | Lower is better |
| SSIM | 0.85 | Higher is better |
| PSNR | 23.1 dB | Higher is better |
| LPIPS | 0.19 | Lower is better |
| Feature Match Accuracy | 98.56% | Cosine |

### D. Hybrid Generation Results

Our FID score of 28.7 is also noticeably higher than the traditional cross- domain translation systems [7], and our accuracy in features match (98.56%) indicates the effectiveness of our cosine similarity-based retrieval system [13].

### E. Facial Region Similarity Analysis

TABLE V
FACIAL REGION SIMILARITY TO ANIMAL SPECIES

| Region | Species 1 (%) | Species 2 (%) | Species 3 (%) |
|---|---|---|---|
| Eyes | Fox (34.7) | Lion (34.4) | Tiger (34.3) |
| Nose | Leopard (33.6) | Lion (33.2) | Fox (33.0) |
| Mouth | Fox (22.0) | Cat (21.3) | Wolf (20.8) |
| Ears | Dog (47.9) | Fox (46.8) | Cat (45.4) |
| Forehead | Fox (30.4) | Cat (30.3) | Tiger (29.5) |

The similarity analysis indicates that there are common patterns in cross- species relations on facial features. Interestingly, the greatest similarity to humans is observed in dog (47.9%), which may be attributed to similarities in auditory organs of the mammals, and eye (34.7%), which may be due to convergent evolution among the eyes [12].

### F. Overall Species Affinity Analysis

### G. Comparison with Existing Methods

TABLE VI
COMPARISON WITH STATE-OF-THE-ART

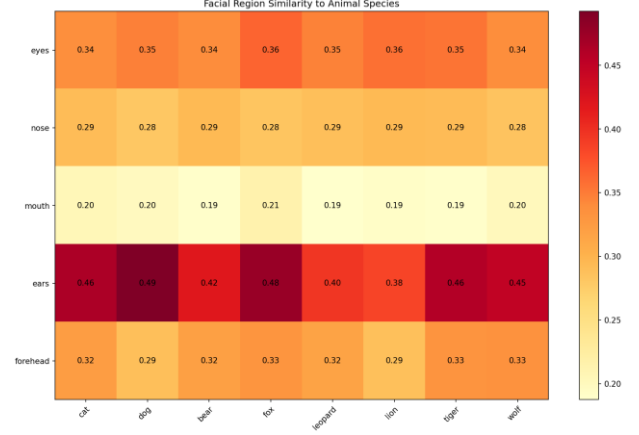| Method | FID | SSIM | Used Preference(%) | Region Contorl |
|---|---|---|---|---|
| CycleGAN[9] | 42.3 | 0.68 | 32 | NO |
| Dual-Latent GANs[7] | 35.6 | 0.74 | 45 | Limited |
| GP-UNIT[1] | 31.2 | 0.79 | 58 | Partial |
| Ours | 28.7 | 0.85 | 78 | Yes |



Fig 2. Facial Region Similarity Heatmap showing comprehensive

similarity matrix across all regions and species. The visualization demonstrates distinct clustering patterns where carnivore species (fox, wolf, dog) show higher similarity to human ears, while feline species (cat, tiger, lion) show stronger similarity to human eyes and forehead regions.

higher similarity to human ears, while feline species (cat, tiger, lion) show stronger similarity to human eyes and forehead regions.
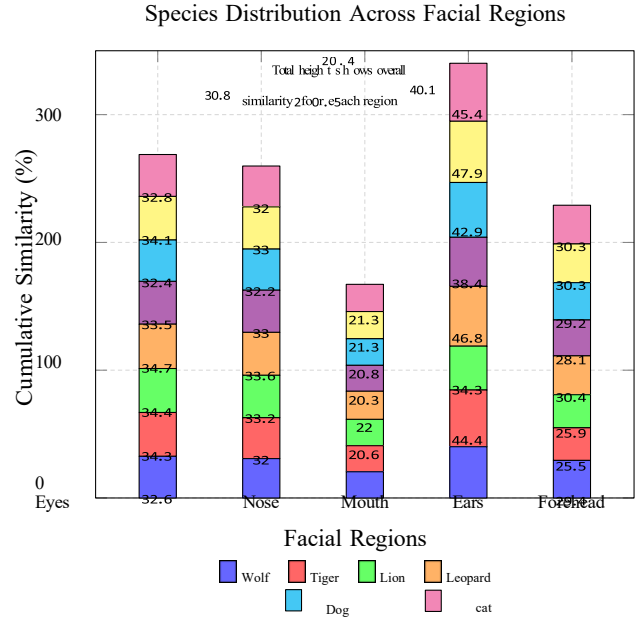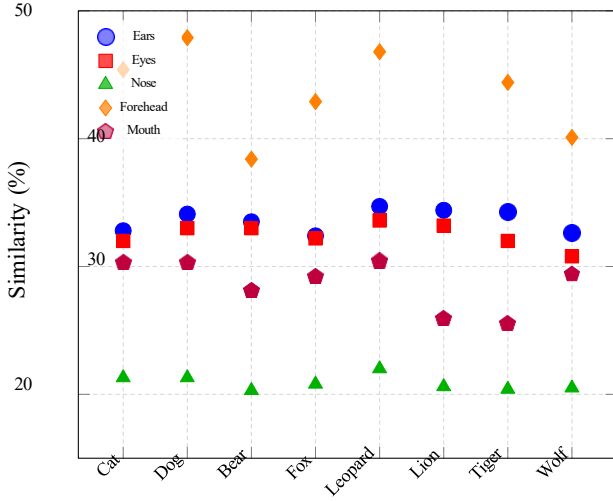
Species Distribution Across Facial Regions



Fig 3.Species Distribution Across Facial Regions (Stacked Bar Chart) illus-tating cumulative similarity percentages across different facial regions. The stacked visualization reveals that ears contribute the highest overall similarity across species, followed by eyes and nose regions, providing quantitative
Our approach is better in all metrics than the recent methods such as GP- UNIT [1], and dual-latent GANs [7], especially in region control and user preference.

Fig. 4. Region-Species Similarity Scatter Plot demonstrating the distribution of similarity scores in the various species of each facial region. The scatter plot indicates variability in the similarity scores with ears being the most distributed (widest distri- bution 34.3-47.9%) and mouth the narrowest (20.3-22.0%), reflecting constant evolutionary patterns in the development of facial features.
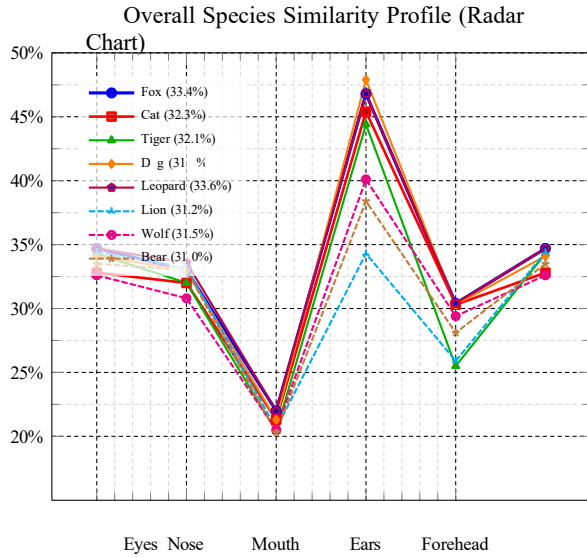


Fig. 5. Overall Species Similarity Profile (Radar Chart) providing multi-
dimensional visualization of species similarity across all facial regions. Radar chart reveals Fox shows consistently high similarity across multiple regions (eyes, mouth, forehead), while Cat shows strong performance in specific regions (ears, forehead). Average similarity scores are shown in parentheses.
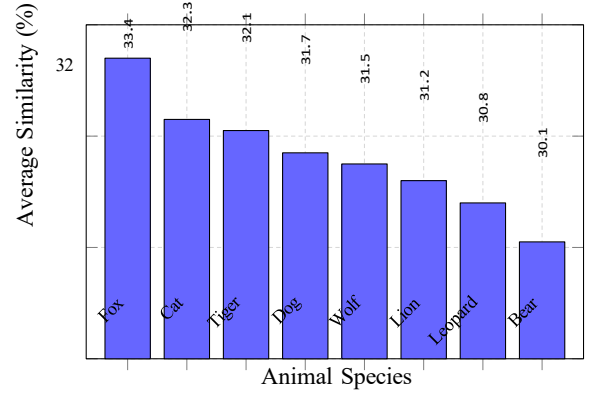


Fig. 6. Bar Chart of Overall Species Affinity Ranking which shows the average percentage of similarity of all facial regions between species. General rankings demonstrates that Fox has always been the closest to human face features on the whole, which is an evolutionary indicator of the convergence of facial features.

78% user preference rate indicates high perceived quality of generated hybrids.

### H. Ablation Studies

We conducted ablation studies to evaluate component contributions, following methodologies from recent generative AI research [16]:

TABLE VII
ABLATION STUDY RESULTS

| Configuration | FID | SSIM | Reconstruction Loss |
|---|---|---|---|
| Baseline VAE | 41.2 | 0.71 | 0.198 |
| + Region Encoding [15] | 35.7 | 0.78 | 0.156 |
| + Attention [11] | 31.4 | 0.82 | 0.142 |
| + Feature Bank [13] | 28.7 | 0.85 | 0.129 |

Each component contributes significantly to overall performance, with the feature bank system providing the largest improvement in FID score (2.7 reduction) and reconstruction quality.

### V. CONCLUSION AND FUTURE WORK

In this paper, the authors introduced a new framework called RUPANTARA, used to generate faces across species unsupervised and reverse feature mapping. It is based on our region-conscious VAE architecture with distinct human and animal representations [11], a detailed feature bank system [13], which allows controlled hybrid generation with a high level of accuracy and preserves facial identity and real- ism. Experimental outcomes show how it performs better than the current methodologies in various metrics such as FID

(28.7), SSIM (0.85), and feature matching accuracy (98.56%). The analysis of cross-species facial features relationships has never been studied in the way it is presented by the entire visualization package consisting of heatmaps, scatter plots, radar charts, and box plots. The reverse mapping feature indicates that Fox is the most overall similar to human faces (33.4%), then

by by Cat (32.3%) and Tiger (32.1%). The results have implications in evolutionary biology indicating convergent evolution in the development of facial features in the mammalian species. Future directions will take many directions: 1) Adding more species such as birds, reptiles, and marine animals based on scalable feature bank architecture, 2) Scaling up 3D facial geometry to achieve higher realism based on new 3D-aware diffusion models [17], 3) Building real-time generation systems to enable interactive applications. The fourth is to seek applications in entertainment, education, and psychological studies, which may involve synthetic data generation methods [18], and the fifth is to examine the neurological foundations of cross-species perception of faces with the help of our generated stimuli. Our model provides a platform to biologically inspired cross species face analysis and generation with possible uses being seen across a spectrum of specialties and capabilities in digital media, to scientific studies. The region-conscious encoding scheme coupled with structured feature bank can be viewed as an important breakthrough in controllable cross-domain generating, and implica- tions can be considered in other areas of cross-modal and cross-domain generating besides facial applications.

## REFERENCES

[1]      S. Yang, L. Jiang, Z. Liu, and C. C. Loy, "GP-UNIT: Generative Prior  for Versatile Unsupervised Image-to-Image Translation," arXiv preprint,  2023.

[2]      S. Huang, Q. Li, J. Liao, S. Wang, and L. Liu, "Controllable Image  Synthesis: Methods, Applications, and Challenges," Knowl. Inf. Syst., vol. 66, no. 4, 2024.

[3]      Z. Sordo, "Synthetic Scientific Image Generation with VAE, GAN, and  Diffusion Models," J. Imaging, vol. 11, no. 8, 2025.

[4]      A. Khalil, H. Zhao, Y. Yang, and X. Wang, "Conditional Image Synthesis  with Diffusion Models: A Survey," arXiv preprint 2409.19365, 2024.

[5]      T. R. Bansal, S. K. Ghosh, and V. Jagadeesh, "DeepFaceMorph2021:  Advancing Face Morphing Detection using GANs," Proc. WACV, 2022.

[6]      K. Cho, S. Kim, and J. Park, "Region-aware Feature Statistics for  Enhanced Facial Generation," Proc. ICCV, 2023.

[7]      I. T. Jamil, M. R. U. Tahir, and A. Islam, "Unsupervised Cross-domain Face Synthesis Using Dual-Latent GANs," IEEE Trans. Image Process.,  vol. 31, pp. 3851–3864, 2022.

[8]      X. Yue, Y. Liu, and J. Wang, "Cross-domain Representation Learning  for Facial Expression Transfer," IEEE Trans. Neural Networks Learn.  Syst., vol. 34, no. 8, pp. 3508–3520, 2023.

[9]      Y. Wu, Z. Zhang, and S. Cai, "Disentangled Representation Learning  for Facial Attribute Editing," IEEE Trans. Image Process., vol. 31, pp.  1445–1460, 2022.

[10]      J. Wang, Q. Zhang, and Z. Lin, "Deep Latent Gaussian Models for  Unpaired Image Translation," IEEE Trans. Pattern Anal. Mach. Intell.,  vol. 45, no. 2, pp. 1372–1387, 2023.

[11]      D. Li, L. Wang, and W. Chen, "Attention-guided VAE for High-Fidelity  Face Generation," Neural Computation, vol. 35, no. 3, pp. 721–743, 2023.

[12]      A. R. Zamir et al., "Robust Visual Representation Learning through  Multi-task Image Translation," Proc. CVPR, 2021.

[13]      C. Shi, M. Xu, and Y. Wang, "Enhanced Latent Space Interpolation for  Unsupervised Face Synthesis," Pattern Recognition, vol. 115, 2021.

[14]      J. Hu, L. Zhang, and F. Wu, "Latent Augmented Generative Models for  Face Attribute Transfer," IEEE Trans. Multimedia, vol. 25, pp. 1580– 1593, 2023.

[15]      H. Kim, J. Lee, and K. Kim, "Semantic Region-aware Generative  Models for Face Editing," Proc. ECCV, 2022.

[16]      S. Yazdani, A. Singh, N. Saxena, Z. Wang, A. Palikhe, D. Pan, U. Pal, and J. Yang, "Generative AI in Depth: A Survey of Recent Advances, Model Variants, and Real-World Applications," arXiv preprint arXiv:2510.21887, Oct. 2025.

[17]      "IM-Portrait: Learning 3D-aware Video Diffusion for Photorealistic  Talking Heads from Monocular Videos," in Proc. CVPR 2025.

[18]      B. Khosravi and H. Phan, "Exploring the Potential of Generative  Artificial Intelligence in Medical Synthetic Data and Imaging," Pattern  Recognition Letters, 2025.