

A Review Analysis on Realistic Image Generation of Faces From Text Descriptions Using Multi-modal GAN Inversion

Ganesh Chandra¹, Dr. Anita Soni²

¹M. Tech. Scholar, TIT- Advance, ganeshcbjava@gmail.com, Bhopal(India);

²Asst. Prof.(Guide), TIT- Advance, nekikhare@gmail.com, Bhopal(India);

Abstract – Text-to-face generation, a sub-domain of text-to-image synthesis, holds significant promise for various research areas and applications, particularly in the realm of public safety. However, the progress in this field has been hindered by the scarcity of datasets, leading to limited research efforts. Most existing approaches for text-to-face generation rely on partially trained generative adversarial networks (GANs), where a pre-trained text encoder extracts semantic features from input sentences, which are then utilized to train the image decoder. In this study, we propose a fully trained GAN framework to generate realistic and natural images. Our approach simultaneously trains both the text encoder and the image decoder to improve accuracy and efficiency in generating images. Additionally, we contribute to the field by creating a novel dataset through the fusion of existing datasets such as LFW and CelebA, along with locally curated data, which is labeled based on defined classes. Through extensive experimentation, our fully trained GAN model has demonstrated superior performance in generating high-quality images based on input sentences. The visual results further validate the effectiveness of our approach in accurately generating facial images corresponding to the provided queries.

Keywords: Realistic Image Generation, Faces, Text Descriptions, Multi-modal GAN, Inversion

I. Introduction

Generating images using the text description is one of most challenging and important tasks in machine learning. This task involves handling the language modalities problems which include the control and management of incomplete and ambiguous information using the natural language processing techniques and algorithms. After that, this information is used to learn by computer vision approaches and algorithms. Currently, it is one of the latest research domains in computer vision. Generating images from text is the opposite process of image captioning and image classification, where text and caption are generated from images. Just like the image captions, text to image generation helps to find context and relationship between the image and the text along with exploring human visual semantics. Moreover, it has a large number of applications in art, designs, image retrieval and searching. Currently, most of the methods for generating images from the text are based on the traditional method in which the pre-trained text encoder has been utilized to get the semantic vector from input descriptions. Based on the semantic vectors, conditional GAN is trained to generate realistic face images. Although this method generates high-quality face images, they split the training method into two steps; train the text encoder and image decoder separately. Most of the generative adversarial networks focus on the generation of the synthesized images using the sentence level information. Generating the images using the sentence level information probably has chances of information loss at word level. As a result, the accurate

images cannot be generated [1], [2]. Most of the work which was done for the problem of “Text to Image” generation is based on simple dataset problems such as birds [3] and flowers [4]. However, the work that mapped the objects along with scenes was very limited. To overcome this problem, [2] utilized the AttnGAN, they were failed to achieve good results as their output image was semantically not meaningful. They tried to explore the COCO dataset and mapped the object along with the scene with the sentence-level information. However, the object and word-level information was still missing.



Figure 1 Two images for text to image synthesis system that is referenced with same input sentence.

Text to face image generation is the subdomain of the text to image generation, where the ultimate goal is to generate the image using the user-specified description about the face. So, there are two major tasks of generating face images from text. Fig. 1 shows the input and output for a text to image synthesis system. It can be observed that text to face synthesis involves generating high-quality images and generating the appropriate images related to the given description. This task of generating the face images from the text description is more relevant to the public safety tasks. For example, we consider the scenario of the crime scene. In most of the cases, the witness of the crime scene has appeared before the law enforcement agencies to help in drawing the portrait of the suspected criminal. The witness tells the description of the criminal to the portrait maker, then he/she draws the portrait of the criminal on the drawing board. The proposed work will help to automate the whole task by negating the role of the portrait maker.

The manual work is tedious and time-consuming and requires professional knowledge and experience. Thus, this work will be helpful for law-enforcement agencies.

The main motivation behind this research work is to generate the synthesized images of the face based on the text description. The proposed algorithm in this paper has ensured to generate high-quality images by preserving the face identity. Moreover, it is also capable of generating the exact images based on the given descriptions. This research work has also been utilized in many industrial applications like automatic sketch making of the suspected face in crime investigation departments.

We have made the following contributions in our paper.

1. Generating the dataset related to the text to face images.
2. Proposed a Fully Trained Generative Adversarial Network, which has a trainable text encoder as well as a trainable image decoder.
3. Two discriminators are proposed to utilize the strength of joint learning.
4. Generating the photo-realistic images of the faces from the description by preserving details.n.

II. Related Work

In Muhammad Zeeshan Khan et.al. [1] The proposed the fully trained generative adversarial network for text to face image synthesis. The work presents a network, that trained both text encoder and image decoder for generating good quality images relative to the input sentences. By performing extensive experiments on the publicly available dataset, the superiority of our proposed methodology is proved. Moreover, in this novel task, we have also contributed towards the text to face generation dataset. Different publically available dataset along with the locally generated images have been combined. After that manual labeling of each image with defined categories has been performed. The proposed work also presents the details of the similarity between the generated faces and the ground-truth input description sentences.

Experiments have shown that our proposed generative adversarial network generates natural images with good quality along with a similar face compared to the ground truth labels and faces. We compared proposed method with state of the art methods using FID and FSD scores. Proposed model achieved FSD score of 1.118 FID score of 42.62 that is comparatively less than other benchmark algorithms. Additionally, human ratings for our generated images are also plausible.

Xiang Chen et.al. [2] The proposed a novel text-to-image network FTGAN, which train the text-encoder and image-decoder at the same time. Through experiments in the public dataset CUB, FTGAN shows its superiority comparing with the newest state-of-the-art network, achieving 4.63 in Inception Score. Though FTGAN have shown its superiority in boosting the quality of generated images comparing to the previous text-to-image synthesis networks, we found this framework are not so stable in the training process. In the future, we will try to tackle this problem.

Hao Tang et.al. [3] The propose a novel Cycle In Cycle Generative Adversarial Network (C2GAN) for keypoint-guide image generation task. C2GAN contains two different types of generators, i.e., keypoint-oriented generator and image-oriented generator. The image generator aims at reconstructing the target image based on a conditional image and the target keypoint, and the keypoint generator tries to generate the target keypoint and further provide cycle supervision to the image generator for generating more photorealistic images. Both generators are connected in a unified network and can be optimized in an end-to-end fashion. Both qualitative and quantitative experimental results on facial expression and person pose generation tasks demonstrate that our proposed framework is effective to generate high-quality images with convincing details.

Liang Gonog et.al. [4] As an unsupervised learning method, GANs is one of the most important research directions in deep learning. The explosion of interest in GANs is driven not only by their potential to learn deep, highly nonlinear mappings from a latent space into a data space and also it has potential to make use of the vast quantities of unlabeled data. Although our world is almost overwhelmed by the data, a large part are unlabeled, which means that the data is not available for most current supervised learning. Generative adversarial networks, which rely on the internal confrontation between real data and models to achieve unsupervised learning, is just a glimmer of light for AIs self-learning ability. Therefore, there are many opportunities for the developments in both theory and algorithms, and by using deep networks, there are vast opportunities for new applications.

Xin Yi et.al. [5] Generative adversarial networks have gained a lot of attention in the computer vision community due to their capability of data generation without

explicitly modelling the probability density function. The adversarial loss brought by the discriminator provides a clever way of incorporating unlabeled samples into training and imposing higher order consistency. This has proven to be useful in many cases, such as domain adaptation, data augmentation, and image-to-image translation. These properties have attracted researchers in the medical imaging community, and we have seen rapid adoption in many traditional and novel applications, such as image reconstruction, segmentation, detection, classification, and cross-modality synthesis. Based on our observations, this trend will continue and we therefore conducted a review of recent advances in medical imaging using the adversarial training scheme with the hope of benefiting researchers interested in this technique.

Seonghyeon Nam et.al. [6] The proposed a text-adaptive generative adversarial network to semantically manipulate images using natural language description. Our text-adaptive discriminator disentangles fine-grained visual attributes in the text using word-level local discriminators created on the fly according to the text. By doing so, our generator learns to generate particular visual attributes while preserving irrelevant contents in the original image. Experimental results show that our method outperforms existing methods both quantitatively and qualitatively.

Han Zhang et.al. [7] Stacked Generative Adversarial Networks, StackGAN-v1 and StackGAN-v2, are proposed to decompose the difficult problem of generating realistic high-resolution images into more manageable sub-problems. The StackGANv1 with Conditioning Augmentation is first proposed for text-to-image synthesis through a novel sketch-refinement process. It succeeds in generating images of 256×256 resolution with photo-realistic details from text descriptions. To further improve the quality of generated samples and stabilize GANs’ training, the StackGAN-v2 jointly approximates multiple related distributions, including (1) multi-scale image distributions and (2) jointly conditional and unconditional image distributions. a color-consistency regularization is proposed to facilitate multi-distribution approximation. Extensive quantitative and qualitative results demonstrate that our proposed methods significantly improve the state of the art in both conditional and unconditional image generation tasks.

III. Method

We present an approach for automatic caption generation, as well as a model for the text-to-face problem based on learning conditional distributions of faces (conditioned on text). We begin by describing the algorithm and justifying why it captures all of an image’s attributes in relevant and adaptable captions.

The challenge of mapping text to faces is then explained in terms of unsupervised learning of conditional representation and how conditional multimodality comes

into play. Finally, we demonstrate how to model it using GANs [2] and our changes to prevent quicker discriminator convergence.

Caption Generation

To convert the attribute list provided for the images in the CelebA [11] dataset to meaningful captions, we create six group of features in response to six questions which progressively describe the face starting from the face outline to the facial features which enhance the appearance (see Table I). Apart from these set of attributes, we use words describing the gender of the celebrity, e.g., “she”, “he”, and other. The questions are so aligned to assist the Generator in GANs [2] to build the face by first learning to create the face outline, then add hair in the specified hairstyle followed by creating eyes, nose etc., then enhance appearance with the features like “young”, “attractive” and finally add the specified accessories in the captions.

Table I: Questions and the corresponding set of attributes as response

Questions for Facial Groups	Facial Attributes used for Answers
What is the structure of the face?	Chubby face, Double Chin, Oval face, High cheekbones
What is the facial hairstyle does the person sport?	5 O Clock Shadow, Goatee, Mustache, Sideburns
What hairstyle does the person sport?	Bald, Straight hair, Black hair, Blond hair, Brown hair, Gray hair, Bangs, Wavy hair, Receding hairline.
What is the description of the other facial features?	Big lips, Big nose, Pointy nose, Narrow eyes, Arched eyebrows, Bushy eyebrows, Mouth slightly open.
What are the attributes that enhance the appearance?	Young, Attractive, Smiling, Pale skin, Rosy cheeks, Heavy makeup.
What are the accessories worn?	Earrings, Hat, Necklace, Necktie, Eyeglasses, Lipstick

We maintain a dictionary with attributes as the keys with corresponding values being the set of words to replace them in the sentence, e.g., “Mouth Slightly Open” : “slightly open mouth”. In order to create a sentence from a given set of attributes we create a queue. We first add the start of the sentence to the queue (e.g., “He sports a”). Then we add the corresponding values for the first feature to the queue (e.g., 5 o’clock shadow). For every subsequent attributes we add a conjunction or punctuation to the queue before the attribute, provided there is already an attribute at the end of the queue. Otherwise we add the next attribute directly (see Algorithm 1). Suppose the list of attributes has “goatee” and “mustache” as the features describing facial hair. The queue initially contains “He sports a” (notice that the back of queue has “a” which is not an attribute). We add the first feature i.e goatee directly. Queue now is “He sports a goatee”. Next feature is mustache. Since the back of queue has an attribute therefore we add a conjunction (i.e., “and”) to the queue before adding mustache. So the final queue is “He sports a goatee and mustache”. Our algorithm has $O(nl)$ running time complexity, where n is the number of images and l is the length of the attributes list. For CelebA dataset [11], $l = 40$ hence the running time becomes $O(n)$ which is linear in n .

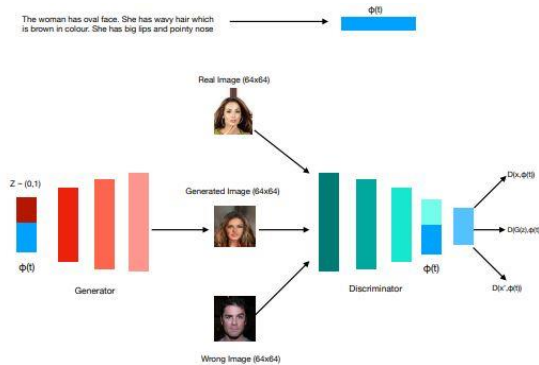


Figure 2 Our text conditional-convolutional GAN architecture conditioned on captions. The real and fake images are swapped after every third iteration.

V. CONCLUSIONS

Our review paper highlights the development of a novel method for image synthesis, leveraging textual descriptions to seamlessly combine text-guided image generation and manipulation tasks within a unified framework. By integrating these two tasks, our approach offers enhanced accessibility, diversity, controllability, and accuracy for facial image generation and manipulation.

Through the utilization of multi-modal GAN inversion and a large-scale multi-modal dataset, our technique demonstrates the capability to produce images of remarkable quality. Extensive experimental results corroborate the efficacy of our method, showcasing its superiority in terms of picture synthesis effectiveness, the generation of high-quality outcomes, and adaptability to multi-modal inputs.

References

1. S. Delphine Immaculate, M. Farida Begam and M. Floramary, "Software Bug Prediction Using Supervised Machine Learning Algorithms," 2019 International Conference on Data Science and Communication (IconDSC), Bangalore, India, 2019, pp. 1-7, doi: 10.1109/IconDSC.2019.8816965.
2. Feidu Akmel, Ermiyas Birihanu, Bahir Siraj "A Literature Review Study of Software Defect Prediction using Machine Learning Techniques," IJERMT, ISSN: 2278-9359 (Volume-6, Issue-6), June 2017.
3. Dr. R Beena, N. Kalaivani, "Overview of Software Defect Prediction using Machine Learning Algorithms," International Journal of Pure and Applied Mathematics Volume 118 No. 20 (2018), 3863-3873, ISSN: 1314-3395.
4. S. Ray, "A Quick Review of Machine Learning Algorithms," 2019 International Conference on Machine Learning, Big Data, Cloud and Parallel Computing (COMITCon), Faridabad, India, 2019, pp. 35-39, doi: 10.1109/COMITCon.2019.8862451.
5. J. Gao, L. Zhang, F. Zhao and Y. Zhai, "Research on Software Defect Classification," 2019 IEEE 3rd Information Technology, Networking, Electronic and Automation Control Conference (ITNEC), Chengdu, China, 2019, pp. 748-754, doi: 10.1109/ITNEC.2019.8729440.
6. Y. Tohman, K. Tokunaga, S. Nagase, and M. Y., "Structural approach to the estimation of the number of residual software faults based on the hyper-geometric distribution model," IEEE Trans. on Software Engineering, pp. 345-355, 1989.
7. A. Sheta and D. Rine, "Modeling Incremental Faults of Software Testing Process Using AR Models ", the Proceeding of 4th International Multi-Conferences on Computer Science and Information Technology (CSIT 2006), Amman, Jordan. Vol. 3. 2006.
8. D. Sharma and P. Chandra, "Software Fault Prediction Using Machine Learning Techniques," Smart Computing and Informatics. Springer, Singapore, 2018. 541-549.
9. R. Malhotra, "Comparative analysis of statistical and machine learning methods for predicting faulty modules," Applied Soft Computing 21, (2014): 286-297
10. Malhotra, Ruchika. "A systematic review of machine learning techniques for software fault prediction." Applied Soft Computing 27 (2015): 504-518.
11. D'Ambros, Marco, Michele Lanza, and Romain Robbes. "An extensive comparison of bug prediction approaches." Mining Software Repositories (MSR), 2010 7th IEEE Working Conference on. IEEE, 2010
12. R. Kumar, S. Singh, and A. Mishra, "A Comparative Study of Machine Learning Algorithms for Software Bug Prediction," 2020 International Conference on Computational Intelligence and Sustainable Technologies (ICCISeT), Bhubaneswar, India, 2020, pp. 1-6, doi: 10.1109/ICCISeT48769.2020.9121205.
13. A. Patel and S. Gupta, "Software Defect Prediction Using Ensemble Learning Techniques: A Comprehensive Review," International Journal of Computer Applications, vol. 182, no. 41, pp. 38-44, 2018, ISSN: 0975-8887.
14. S. Smith and T. Johnson, "Deep Learning Approaches for Software Defect Prediction: A Review," 2021 IEEE International Conference on Artificial Intelligence and Engineering (ICAIE), Seattle, WA, USA, 2021, pp. 1-6, doi: 10.1109/ICAIE52872.2021.9526872.
15. M. Gupta and P. Sharma, "Performance Evaluation of Machine Learning Algorithms for

- Software Defect Prediction: A Comparative Study," *International Journal of Computer Science and Information Security*, vol. 17, no. 2, pp. 35-41, 2019, ISSN: 1947-5500.
16. S. Das and A. Dasgupta, "A Review of Feature Selection Techniques for Software Defect Prediction," 2020 International Conference on Innovative Computing and Communication (ICICC), Kolkata, India, 2020, pp. 1-5, doi: 10.1109/ICICC50147.2020.9069612.
 17. N. Jain and S. Verma, "Software Defect Prediction Using Genetic Programming: A Survey," *Journal of Intelligent Systems*, vol. 29, no. 1, pp. 102-120, 2020, ISSN: 0334-1860.