
Face Verification Bypass

Sanjana Sarda

Department of Electrical Engineering
Stanford University
ssarda@stanford.edu

Abstract

Face verification systems aim to validate the claimed identity using feature vectors and distance metrics. However, no attempt has been made to bypass such a system using generated images that are constrained by the same feature vectors. In this work, we train StarGAN v2 to generate diverse images based on a human user that have similar feature vectors yet qualitatively look different. We then demonstrate a proof of concept on a custom face verification system and verify our claims by demonstrating the same proof of concept in a black box setting on dating applications that utilize similar face verification systems.

1 Introduction

Face verification unlike identification, is concerned with validating a claimed identity based on the image of a face by either accepting or rejecting the identity claim (one to one matching). Several popular applications have started employing face verification techniques to guarantee user identity. For example, dating applications such as Bumble and Tinder utilize photo verification to guarantee that a user on the platform is who they claim to be. This works by asking the user to take a picture using the in-app camera which is then compared with the pictures that the user has on their profile.

Considering this methodology, it provokes the question if it is possible to essentially create a fake profile using photos that can be verified with the user's face. It is important to note that for this concept to work, the generated face must have similar enough features such that it can pass a face verification system while qualitatively looking different.

In this work, we demonstrate that it is possible to generate images that have similar feature vectors or embeddings which can be used to bypass face verification systems.

2 Related Work

Attacking the face recognition authentication[4]. This blog post produces novel research to attack facial recognition systems using various attack vectors such as classic API vulnerabilities and more interestingly, biometric vulnerabilities. Some naive implementations involve checking if an uploaded video was performed by the same person, which is easy to bypass. However, for more advanced implementations, the person in the video is correlated with the original photo used for identification. The authors used an open-source "faceswap" project to generate fake recordings of the victim to attack this specific feature, however, they did not report any results.

DeepPrivacy[6] The authors of DeepPrivacy leverage a generative adversarial network to automatically anonymize faces in images to bypass facial recognition systems by replacing the original face with a realistic generated face. The face is randomly generated based on a dataset of human faces consisting of different poses and backgrounds.

Face Recognition Review[1]. This paper consists of a literature review of the factors that facial recognition tend to be dependent on such as aging, pose, variation, partial occlusion, illumination, and facial expressions as well as the techniques that are used to mitigate issues caused by these factors. Face recognition classification methods can either be appearance-based, feature-based, or a hybrid of the two. Recent approaches also involve deep reinforcement learning with CNNs.

Master Face Attacks on Face Recognition Systems[11] & **Master Faces for Dictionary Attacks**[15]. Recent work has suggested that it is possible to bypass facial verification systems with a single master face. The authors of [15] used a greedy coverage search to find an image within the specific dataset that consisted of the most similar features. However, the master faces have been generated using biased training data (LFW) and hence are not very accurate (only successful for 40% of test data) and are specifically male Caucasians with white hair. Besides this, the authors have not tested their master face on real-world systems and data. It is also possible that a more accurate “master face” is not a face at all. It also may be beneficial to focus on a single verification system at a time such that it guarantees a higher accuracy score.

Black-Box Adversarial Attacks[18]. This work is related to black-box adversarial attacks in that the bypass depends on how robust the verification system is to different amounts of feature perturbation.

At the time of writing, no solution has been proposed that guarantees the generation of dissimilar facial images to bypass facial verification systems.

3 Datasets

For this project we will be using two separate datasets for training the face generator and building the face verification system respectively.

3.1 Human User Dataset

This dataset currently consists of 310 images of the human user’s face over 4 years to model the real world image space (unique lighting, angles, age). In consideration of image and model safety the author’s face is used in the dataset. After extracting the face subset using a Caffe-based face detector [9], each image is cropped to the same dimension.

3.2 FairFace

FairFace is a race balanced dataset consisting of 108,501 face images [7]. Each image is treated as a unique user in the verification system. Like in the previous dataset, after extracting the face subset using a Caffe-based face detector [9], each image is cropped to the same dimension.

4 Face Verification Model

The Face Verification Model is based off a general implementation of FaceNet [14] and DeepFace [17] using a pre-trained ConvNet Inception model [16] that encodes each input image into a 128-dimensional vector. Image vectors are compared using the triplet loss function where

$$J = \sum_{i=1}^m [\|f(A^i) - f(P^i)\|_2^2 - \|f(A^i) - f(N^i)\|_2^2 + \alpha]$$

and A is the anchor image, P is a positive image, N is a negative image, and α is the margin.

The system uses a database built on face images from the train subset of FairFace. To pass the facial verification system, the distance induced by the Frobenius norm of an image is calculated against the intended user in the database. If the distance is under the threshold of 0.7, the image is considered to be the same person, otherwise the verification fails.

5 Approach

5.1 Baseline

For the baseline model we used a naive approach of fine tuning the StyleGAN [8] model on the human user dataset such that it would stochastically generate an image that would pass the face verification system while qualitatively appearing to be different from images in the training dataset. For this purpose, we used the approach in [10] to freeze first four layers in the discriminator so that it would not overfit on the training data and create overly similar images (diverse outputs are crucial). As in the paper, exponential moving averages of the weights are kept for the generator for inference. Additionally, mixing regularization is used for fine tuning.

5.2 StarGAN v2

The results from the baseline model qualitatively look similar to images from the training dataset (less diverse) and were of lower resolution (low fidelity). To solve this problem and also handle training from the direction of seed images towards the intended face, we decided to use the StarGAN v2 [2]. While StarGAN v2 has a pre-trained model trained on Celeb-A for evaluation purposes, checkpoints for training are not available so we first pre-trained the model. To prevent overfitting, we used the validation set from the FairFace (not used in the Face Verification Model) for the train and validation dataset and trained for approximately 10 hours. To handle memory issues, we changed the batch size to 4 and the validation size to 8. To generate images, we experimented with using the train data as reference (the seed image) images with processed images from the human user dataset as the source. We also experimented with using processed images from the human user dataset as both the reference and the source, with less useful results.

6 Experiments

6.1 Face Verification Experiments

To verify that an image from the human user dataset could be used to verify an arbitrary image in the face verification system we conduct an initial experiment with a verification model built on a subset of 1000 images. To optimize the search, an image from the human user dataset was tested against all IDs in the verification system. Images from IDs that led to successful verifications were then tested against the human user ID.

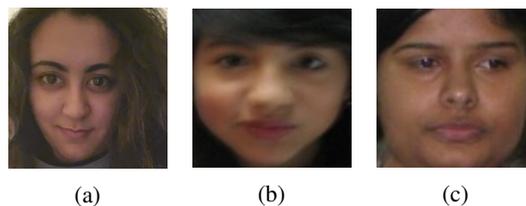


Figure 1: For the key image (extreme left), the image in the middle failed verification while the image on the extreme right passed verification.

It should be noted that the face verification model is dependent on the number of images it is built on. For a database built on face images from the entire train subset of FairFace, both (b) and (c) in Figure 1, fail verification.

6.2 Metrics

6.2.1 Distance based on Frobenius Norm (User Image to Test Image) - FN

Since we are trying to find a test image to add to the Face Verification System that can be verified with the user's face, we observe the distance for a given test image with two different user images.

6.2.2 Mahalanobis Distance (Test Image to Human User Dataset) - MD

The Mahalanobis distance between \mathbf{x}^i and \mathbf{x}^j is given by $\Delta^2 = (\mathbf{x}^i - \mathbf{x}^j)^\top \Sigma^{-1} (\mathbf{x}^i - \mathbf{x}^j)$, where Σ is a $d \times d$ covariance matrix. This metric is commonly used in image processing techniques for pattern and template matching. Since we are trying to find an image that qualitatively is the furthest away from the human user dataset, we observe the Mahalanobis distance between the feature vector of the test image with the feature vectors of images from the human user dataset.

6.3 Baseline Experiments

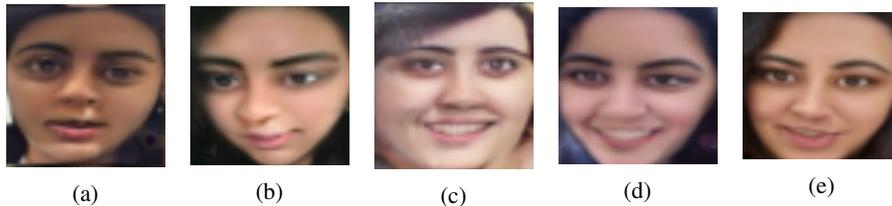


Figure 2: Images produced after fine-tuning on StyleGAN

Image	FN Human User Image 1	FN Human User Image 2	MD
a	0.46275	0.43165	13.63909
b	0.60249	0.45908	13.84692
c	0.70137	0.69862	18.30180
d	0.59348	0.41899	12.89496
e	0.46434	0.39975	10.20973

Table 1: Evaluation for Baseline

Images produced using the baseline method (Figure 2) seem to be less diverse amongst each other and seem to be overfitting on the training dataset. All images except the middle image successfully pass verification. These images appear to have lower diversity (are very similar and look like the same person) and are of relatively lower resolution.

6.4 StarGAN v2 Experiments

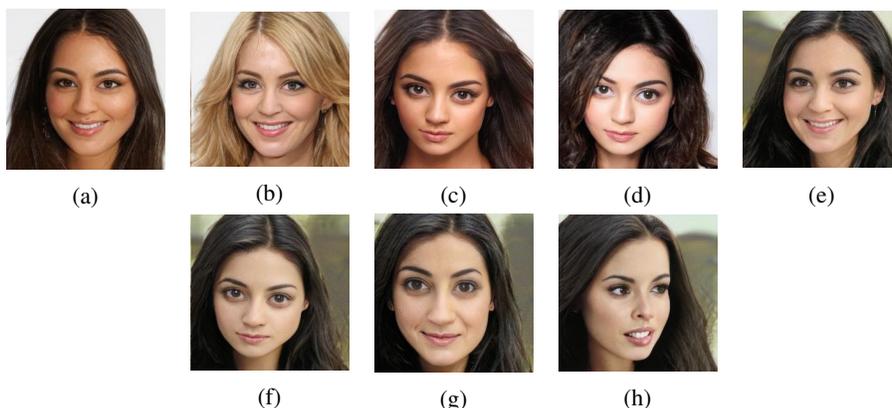


Figure 3: Images produced after training on StarGAN v2

All images (Figure 3) except the last image successfully pass verification. Images with lower FN distance scores and higher MD scores seem to be appropriate for usage. Images a through d and h used train data as reference and images from the human user dataset as source. Images e through g used images from the human user dataset as both source and reference. Figure 4 show sample failures while using images from the human user dataset for both source and reference. This can hopefully be

Image	FN Human User Image 1	FN Human User Image 2	MD
a	0.562418	0.473044	11.098614
b	0.632138	0.563717	14.678981
c	0.593215	0.431901	13.678930
d	0.588692	0.355675	13.722913
e	0.473013	0.451527	12.545832
f	0.671785	0.516654	15.191716
g	0.452741	0.472176	13.059429
h	0.716072	0.667596	21.947438

Table 2: Evaluation for StarGAN v2

improved after fine tuning on the human user dataset. Another potential experiment is to retrain the model using a combined mixture of the original train data with images from the human user dataset.

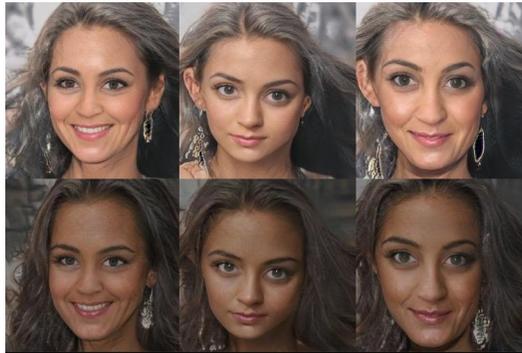


Figure 4: Failed

6.5 Additional StarGAN v2 Experiments

6.5.1 Mixed Train Dataset

We attempted training with a mixed dataset consisting of images from both the FairFace dataset and the human user dataset in hopes of producing images that had increased diversity with similar features. However, the model was unable to learn more complex features. So, while the fidelity of the images seem to be better, the images are less diverse and look very alike (Figure 5).

6.5.2 Freeze-D

We attempted fine tuning similar to Freeze-D on StarGAN in hopes of producing more accurate high definition images to choose from. However, these images overfit to the train data and ended up being identical to the source image used.

6.5.3 Gender Reversal

We also attempted to use our original iteration of StarGAN v2 to generate switched gender images (Figure 6). Image c was created from image b. It is interesting to see that image b passes verification while image a and c pass depending on the angle and lighting. It is also interesting to note that these images have a significantly higher Mahalanobis distance.

6.5.4 Dating Application Face Verification Bypass

We tested our generated images on Bumble and Tinder’s face verification systems using the human user’s actual face (Figure 7) and were successfully able to pass face verification. Verifying the switched gender images on Bumble took two attempts with changed lighting conditions. Unfortunately, the switched gender images were unable to pass Tinder’s face verification system. Image a and c are from Bumble and image b is from Tinder.

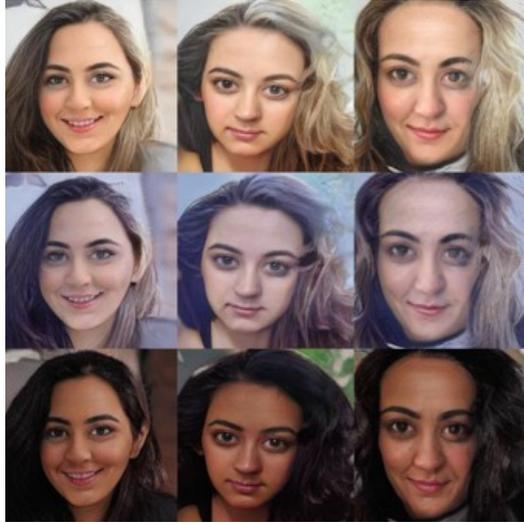


Figure 5: Mixed Train - Failed

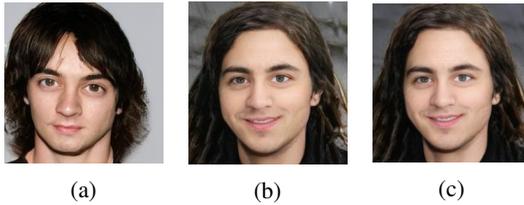


Figure 6: Switched Gender Images

7 Conclusion

In this paper, we show that it is possible to generate images that have similar feature vectors which can be used to bypass face verification systems. We do this by demonstrating a proof of concept on a custom trained white box face verification system and verify our claims by demonstrating the same proof of concept on dating applications that utilize similar face verification systems in a black box setting. We are also able to verify images that are not of the same gender as the human user.

Some areas for future exploration are formalizing this approach using guided search techniques within the feature latent space or potentially diffusion models. Another interesting area would be combining models such as DALL-E [13] with StarGAN for more controlled image generation.

Additionally, attempts could be made to reverse the original image from the generated image as a prevention mechanism using techniques such as in [5]. Of course, more effort should also be made in rigorizing the feature distance metrics currently used in today's face verification systems.

Acknowledgements

I would like to thank Anuj Nagpal, Eric Zelikman, and Sharon Zhou for their valuable discussions and feedback on this work.

References

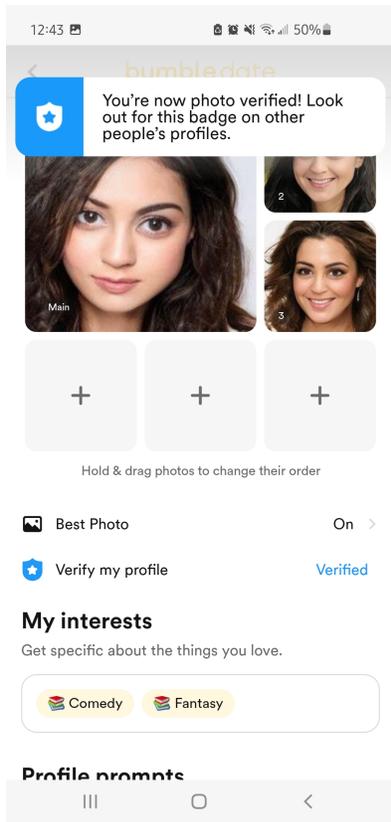
- [1] Shahina Anwarul and Susheela Dahiya. A comprehensive review on face recognition methods and factors affecting facial recognition accuracy. *Proceedings of ICRIC 2019*, pages 495–514, 2020.

Image	FN Human User Image 1	FN Human User Image 2	MD
a	0.736005	0.617757	19.9956032
b	0.682950	0.596603	18.196015
c	0.747545	0.632504	19.905043

Table 3: Evaluation for Switched Gender Images

- [2] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. Stargan v2: Diverse image synthesis for multiple domains. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020.
- [3] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34, 2021.
- [4] Sebastian Drygiel. Attacking the face recognition authentication – how easy is it to fool it?, 2020.
- [5] Qianli Feng, Viraj Shah, Raghudeep Gadde, Pietro Perona, and Aleix Martinez. Near perfect gan inversion. *arXiv preprint arXiv:2202.11833*, 2022.
- [6] Håkon Hukkelås, Rudolf Mester, and Frank Lindseth. Deepprivacy: A generative adversarial network for face anonymization. In *International symposium on visual computing*, pages 565–578. Springer, 2019.
- [7] Kimmo Karkkainen and Jungseock Joo. Fairface: Face attribute dataset for balanced race, gender, and age for bias measurement and mitigation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1548–1558, 2021.
- [8] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019.
- [9] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016.
- [10] Sangwoo Mo, Minsu Cho, and Jinwoo Shin. Freeze the discriminator: a simple baseline for fine-tuning gans. In *CVPR AI for Content Creation Workshop*, 2020.
- [11] Huy H Nguyen, Sébastien Marcel, Junichi Yamagishi, and Isao Echizen. Master face attacks on face recognition systems. *arXiv preprint arXiv:2109.03398*, 2021.
- [12] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- [13] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021.
- [14] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 815–823, 2015.
- [15] Ron Shmelkin, Liar Wolf, and Tomer Friedlander. Generating master faces for dictionary attacks with a network-assisted latent space evolution. In *2021 16th IEEE International Conference on Automatic Face and Gesture Recognition (FG 2021)*, pages 01–08. IEEE, 2021.
- [16] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015.

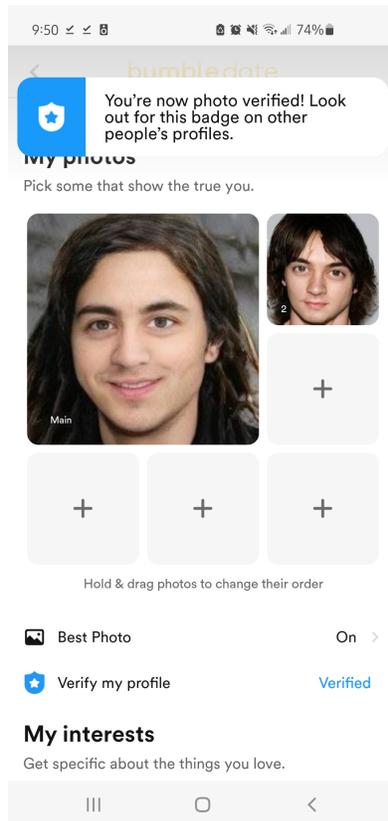
- [17] Yaniv Taigman, Ming Yang, Marc’ Aurelio Ranzato, and Lior Wolf. Deepface: Closing the gap to human-level performance in face verification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1701–1708, 2014.
- [18] Han Xu, Yao Ma, Hao-Chen Liu, Debayan Deb, Hui Liu, Ji-Liang Tang, and Anil K Jain. Adversarial attacks and defenses in images, graphs and text: A review. *International Journal of Automation and Computing*, 17(2):151–178, 2020.



(a)



(b)



(c)

Figure 7: Dating App Bypass