Exposing GAN-Generated Profile Photos from Compact Embeddings Shivansh Mundra¹, Gonzalo J. Aniano Porcile¹, Smit Marvaniya¹, James R. Verbus¹, Hany Farid^{1,2} LinkedIn¹ and University of California, Berkeley²

Overview



Average of 400 GAN (left) and 400 LinkedIn profile images (right) revealing a highly regular GAN facial structure. We learn simple/compact embeddings to capture these GAN regularities.

Dataset



- . StyleGAN1 (10k) 4. generated.photos (10k) 2. StyleGAN2 (10k) 5. stable diffusion (1.5k)
- 3. StyleGAN3 (10k) 6. profile photos not shown (100k)

Methods

Each image is:

- converted to grayscale
- 2. resized to 128×128 pixels
- 3. scaled into the intensity range [0,1]

PCA

- learn a PCA basis
- 2. project onto 128-D basis
- 3. compute reconstruction error as l2-norm (RE)

Autoencoder (AE)

- train an autencoder
- 2. one hidden layer (n=128)
- 3. ReLU activation, Adam optimization, *l*2 regularization, constant learning rate
- compute reconstruction error as l2-norm (RE)

Fourier

- represent using a fixed basis
- 2. projected onto 12x12 low-pass
- compute reconstruction error as l2-norm (RE)

Logistic Regression (LR)

- train LR classifier on PCA representation
- 2. train LR classifier on AE representation
- 80/20 training/testing split (all results on right are for testing)

CNN [*]

- 1. ResNet-50 pre-trained on ImageNet
- refined to classify an image as real or fake





odel	classifier	training	testing	TPR	
CA	RE	StyleGAN1	StyleGAN1	71.7%	
CA	RE	StyleGAN2	StyleGAN2	82.9%	
CA	RE	StyleGAN3	StyleGAN3	79.1%	
CA	RE	StyleGAN(123)	StyleGAN(123)	70.7%	
CA	LR	StyleGAN(123)	StyleGAN(123)	99.6%	
AE	RE	StyleGAN1	StyleGAN1	79.7%	
٩E	RE	StyleGAN2	StyleGAN2	92.0%	
٩E	RE	StyleGAN3	StyleGAN3	86.0%	
٩E	RE	StyleGAN(123)	StyleGAN(123)	79.0%	
AE	LR	StyleGAN(123)	StyleGAN(123)	99.5%	
urier	RE	StyleGAN3	StyleGAN3	3.5%	1
AE	RE	StyleGAN3	generated.photos	68.2%	
\E	RE	StyleGAN3	Stable Diffusion	0.9%	Г
NN	CNN	[*]	StyleGAN1	60.1%	
NN	CNN	[*]	StyleGAN2	46.8%	
NN	CNN	 [*]	StyleGAN3	9.6%	

generic Fourier basis fails

generalizes (mostly) to other GAN engines, but doesn't generalize to diffusion synthesis

LR on 128-D outperforms CNN

Robustness to attack

retrain PCA/AE on scaled/cropped images

TPR for PCA + RE = 22.7%

TPR for AE + RE = 38.8%

- TPR for PCA + LR = 77.9% (down from 99.6%)
- TPR for AE + LR = 78.8% (down from 99.5%)

[*] Sheng-Yu Wang, et al. CNN-generated images are surprisingly easy to spot... for now. *ICCV*, 2020.