
Comparison of DCGAN and Deep CNN-based VAE for Generating Fake Faces

Jordan Hanson

UFID: 9392-2594

Abstract

Generative models have become increasingly popular in recent years with LLMs such as GPT. Presented in this paper is a comprehensive comparative analysis of two state-of-the-art generative models: Generative Adversarial Networks (GANs) and Variational Autoencoders (VAEs). This paper analyzes the results of two specific architectures trained on the CelebA data-set to generate fake images. These architectures are deep CNN-based Variational Autoencoder and deep convolutional generative adversarial network (DCGAN).

1. Introduction

1.1. Variational Autoencoders

A Variational Autoencoder (VAE) is a type of generative model that learns the latent space of a data set to generate images. First, the model encodes an image into a latent space z where $z = \text{Encoder}(x) \approx q(z|x)$. This is done using an encoder network defined later in the paper. Then the model decodes this latent space z into a reconstructed image $x' = \text{Decoder}(z) \approx p(x|z)$. Two different losses are used to make the model generate images that resemble the input and the learned latent space resembles a Gaussian distribution. The VAE is designed to learn a latent distribution from input images and reconstruct the images using a series of encoding, sampling, and decoding steps.

1.2. Generative Adversarial Networks

A Generative Adversarial Network (GAN) is a type of generative model that learns to generate images that resemble the input data. There are two main components of a GAN, the generator and the discriminator. The generator is a network that generates images from a random noise value called the z value. The generator tries to create images that resemble the input data. The discriminator is a network that takes an image as an input and tries to determine if the image is real or fake. The training process involves pitting these two networks against each other. The generator is trained

to fool the discriminator and the discriminator is trained to determine which images are from the CelebA dataset and which are generated by the Generator. This is repeated until the generator is able to produce realistic images.

2. Data

2.1. CelebA

The data-set used to train these models is the CelebA data-set. This data-set contains over 200,000 face images with each image having a resolution of 178×218. There are over 10,000 people represented in the data-set. The data-set comes with a csv file that contains binary attributes for each image such as smiling, eyeglasses, and bald. To work with this data-set: the images were converted to 64×64, normalized, and centered. The converted images are in the shape of 3x64x64 where 3 is the number of color channels and 64 is the height and width of each channel.

3. Architecture and Training

3.1. VAE Network Architecture

Presented in this paper is a deep CNN-based Variational Autoencoder to generate fake faces trained on the CelebA dataset. The architecture is based off figure 2 from the research paper "Deep Feature Consistent Variational Autoencoder" (Hou et al., 2016). The architecture consists of an encoder and a decoder network, which are designed to compress and reconstruct the input images.

Encoder: The Encoder is a neural network that takes an input image (represented as a tensor) and outputs the mean and log-variance of the latent distribution associated with the input. As shown in the left column of figure 2, the Encoder consists of four convolutional layers (Conv2d), followed by batch normalization (BatchNorm2d) and leaky ReLU activation functions (LeakyReLU). The output of the encoder is then flattened and passed through two fully connected layers (Linear) to produce the mean and log variance.

Decoder: The Decoder is a neural network that takes a latent space (z) reconstructs the original image from it. As shown on the right column of figure 2, the Decoder be-

055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107
108
109

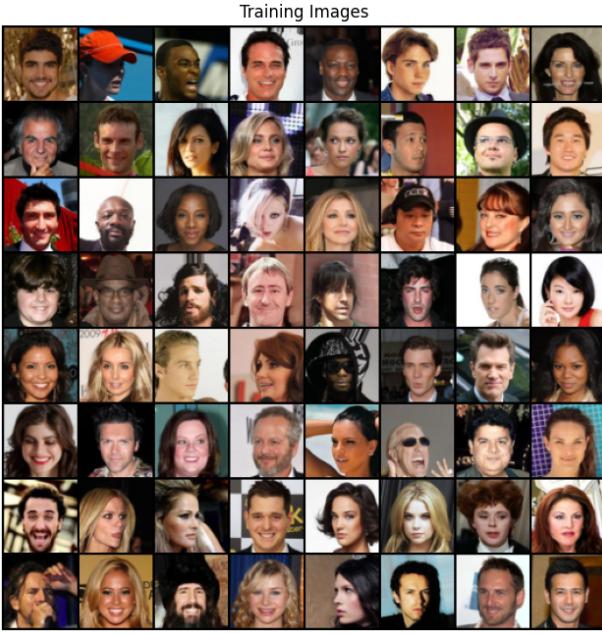


Figure 1. CelebA example images after transforming and denormalizing images.

gins with a fully connected layer (Linear) and an unflatten operation, followed by a series of upsample (Upsample) and convolutional layers (Conv2d) with batch normalization (BatchNorm2d) and leaky ReLU activation functions. The final layer is a convolutional layer that produces the image.

After encoding the image we sample from the mean and log variance to generate a reparameterized z value. This done during training to allow backpropogation to function correctly. This is done in my sample function where $z = \mu + \sigma \cdot \epsilon$ where μ is the mean found by the encoder, σ is the log variance found by the encoder, and ϵ is a random noise value from a normal distribution.

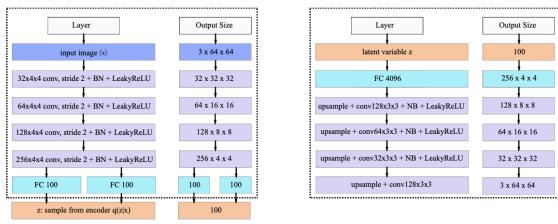


Figure 2. Architecture for deep CNN-based Variational Autoencoder. (Hou et al., 2016)

3.2. GAN Network Architecture

Presented in this paper implemented a Deep Convolutional Generative Adversarial Network (DCGAN) to generate fake faces trained on the CelebA dataset. The architecure used is shown in figure 3 and outlined in the research paper "Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks" (Radford et al., 2016). This GAN architecture consists of two main components: the Discriminator and the Generator. Both components are implemented as neural networks using the PyTorch library. Both the Generator and Discriminator are based off the architecture shown in figure 3.

Discriminator: The Discriminator is a neural network that takes an input image (represented as a tensor) and outputs the probability that the image is real (from CelebA dataset). The Discriminator consists of four convolutional layers (Conv2d), followed by batch normalization (BatchNorm2d) and leaky ReLU activation functions (LeakyReLU). The final layer is a convolution layer followed by a sigmoid activation function, which outputs the probability that the value is real (1 being completely real and 0 being fake).

Generator: The Generator is a neural network that takes a random noise vector (z) as input and generates an image. The noise vector is first reshaped into a 3D tensor, which is then passed through a series of transposed convolutional layers (ConvTranspose2d). Each transposed convolutional layer is followed by batch normalization (BatchNorm2d) and ReLU activation functions. The final layer is a transposed convolutional layer followed by a Tanh activation function, which produces the image.

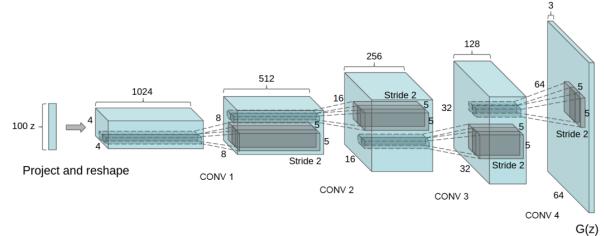


Figure 3. Architecture for Deep Convolutional Generative Adversarial Network (DCGAN). (Radford et al., 2016)

3.3. GAN Training

The training process involves making these two networks compete against each other. To train the discriminator, real images and fake images from the Generator (using random noise as input) are fed into the discriminator. The output is then fed into a Binary Cross Entropy loss (BCE) function with a label of 0 for the fake output and a label of 1 for the real output. The discriminator weights are then optimized

110 based on the loss function using an ADAM Optimizer. The
 111 Generator is trained by generating fake images from random
 112 noise and feeding it into the discriminator. The output of
 113 the discriminator is then fed into the BCE loss function with
 114 a label of 1 and the generator weights are optimized using
 115 ADAM. This is done for about 60 epochs, with a batch size
 116 of 128, and a latent space (z) of dimension 100.
 117

118 3.4. VAE Training

119 The VAE was trained for about 50 epochs with a batch size
 120 of 256. For the training process, images are fed into the
 121 encoder to get the μ (mean) and σ (log variance) for that
 122 image. This is then sampled from using a sample function
 123 to generate a z value.
 124

125 where $z = \mu + \sigma \cdot \epsilon$ where ϵ is random noise from a normal
 126 distribution with mean of 0 and standard deviation of 1.
 127

128 This z value is then used to generated a reconstructed image.
 129 The loss functions are then applied to the ouputs. I use two
 130 different loss functions to update the weights of the network
 131 using gradient descent. These two loss functions are recon-
 132 struction loss and KL divergence. The reconstruction loss
 133 using a mean squared error function (using a sum reduction)
 134 measures the difference between the input image and the
 135 reconstructed image. The KL divergence measures how far
 136 the learned latent space differs from a normal distribution.
 137 $KL\ divergence = -0.5 * \sum_{batch} (1 + \sigma - \mu - e^{\sigma})$ where σ is
 138 the log variance and μ is the mean, both found using the
 139 encoder.
 140

141 The KL divergence ensures that the model generates smooth
 142 and continuous data while the reconstruction loss ensures
 143 the reconstructed image resembles the input image. For
 144 each batch, these losses are added together and the weights
 145 are updated using an ADAM optimizer.

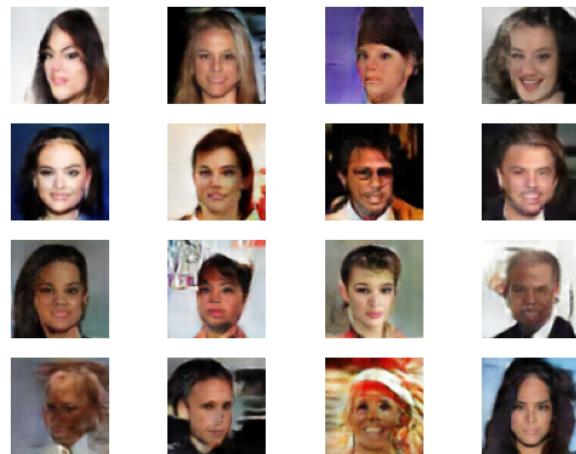
146 4. Results

147 4.1. GAN Results

148 After about 50 epochs, the GAN was able to generate face
 149 like images. The faces as shown in figure 4 are very good
 150 for only 50 epochs but definitely show some morphing be-
 151 tween different images. The images were created by feeding
 152 random data drawn from a Gaussian distribution into the
 153 generator. Some images clearly resemble faces but some
 154 images are too noisy to resemble a face.
 155

156 4.2. VAE Results

157 As shown in figure 5, the VAE was able to generate face
 158 like images from random noise. The images clearly show
 159 a general face in the middle of each image but the outside
 160 seems to morph between different attributes. Some images
 161 clearly resemble a known person with added attributes such
 162
 163
 164



165 *Figure 4. Fake Faces Generated by a DCGAN from Inputting Fake
 166 Noise.*

167 as one of the images shows a clearly male actor with blurred
 168 long hair. In this case, the VAE has learned to encode im-
 169 portant features of human faces, such as facial structure,
 170 eyes, nose, mouth, and hair, into its latent space. When gen-
 171 erating new images from random noise, the VAE samples
 172 from this latent space and combines the encoded features
 173 in various ways, which can result in images that resemble
 174 known individuals but with added or altered attributes.
 175

176 The morphing effect observed in the generated images can
 177 be attributed to the fact that the VAE's latent space is a
 178 continuous space. So it represents certain images as points
 179 in the latent space and when random noise is inputted into
 180 the model the noise often lies in between known points and
 181 so the model combines features resulting in this blurring of
 182 attributes.
 183

184 4.3. Comparison of Models

185 While both models were able to generate images that repre-
 186 sented, both models clearly have advantages and disadvan-
 187 tages. In terms of training speed, the VAE clearly had an
 188 advantage here. After only a couple epochs, the VAE was
 189 able to generate faces even if they resembled the original
 190 images and the GAN took a couple epochs to show faces.
 191 The GAN also took significantly longer to run for 50 epochs,
 192 the GAN took around 6.5 hours while the VAE took about
 193 4.5 hours for 50 epochs. The results clearly show that the
 194 VAE is able to generate better looking faces but the images
 195 look like they just take images it has already seen and add at-
 196 tributes to it. The GAN creates more original images but the
 197 images look less like faces with some being really bad and
 198 some being pretty good. The VAE also has the advantage
 199

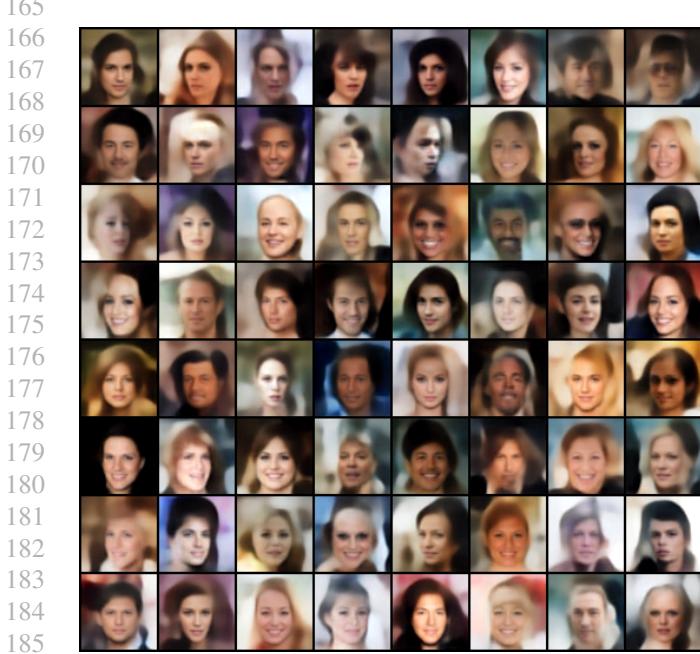


Figure 5. Fake Faces Generated by VAE from Inputting Fake Noise.

of being able to isolate features and manipulate the latent space with more control.

The VAE and GAN have different advantages and disadvantages and depends on the use case for which model is best. The GAN is able to generate more original images but takes more training to get good results. This is likely due to the fact that two models are trained from scratch to generate and determine which image is correct. The VAE shows more realistic results but seems to generate more blurred images where attributes are mixed together.

4.4. Exploring VAE Latent Space

The nature of the VAE allows more finetuned choice of latent space (z) for image reconstruction. By feeding certain images into the encoder, we are able to determine the latent space of different attributes. This is done by feeding all images into the encode and sample functions to find all the z values. We then can separate all the z values into two groups those with the attribute and those without the attribute. We then can find the mean of the z values for each group. By finding the difference of these two groups we can determine the latent space for a specific attribute ($z_{att} = zMean_{attribute} - zMean_{nonAttribute}$). The attributes are defined in a csv file included with the CelebA dataset.

In figure 6, the reconstruction of the latent space of a smile is shown. Clearly the smile is the most prominent property of the image but it seems that long hair as well as a couple

other attributes are included as well. This is likely due to the data-set having more people with long hair and those attributes smiling. In figure 7, the smile z value is added to the z value of other non smiling images. This is done using an intensity factor of 2.5 being applied to the smile z value. $z = z_{nonattimage} + intensity \cdot z_{smile}$. This shows the power of the VAE where attributes can be added to images.

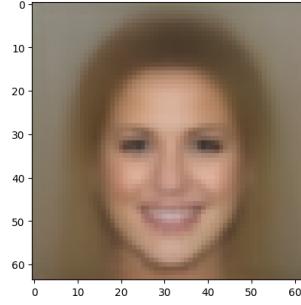


Figure 6. Reconstruction of Latent space of a smile.

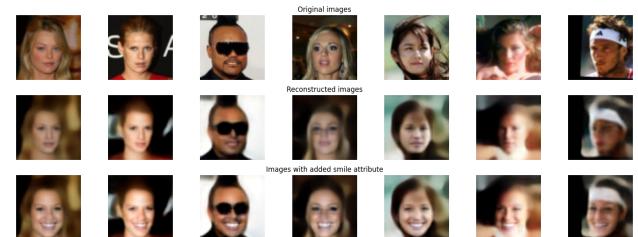


Figure 7. Applying Smile Attribute to Non Smiling Faces.

Figure 8, shows the same process this time with eyeglasses. The same process of determining the z value of eyeglasses are added to the z value of non eyeglasses pictures. Both figure 7 and 8, show good results where non attributes images are altered to display the attribute.

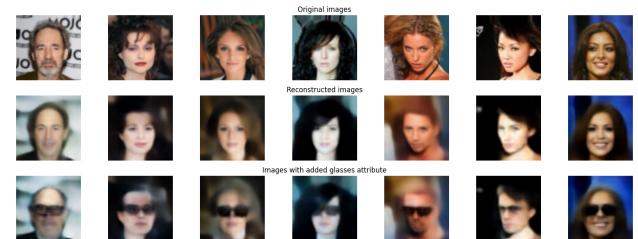


Figure 8. Applying Eyeglasses Attribute to Faces without Eye-glasses.

220 **5. Future Work**

221 This experimentation could be expanded by increasing the
222 epochs and including other face data-sets to improve generated
223 image quality. These architectures could also be applied to other types of images to see what they are able to uncover and generate. Further experimentation could try to solely generate faces by combining different attributes together using their z values. This could be interesting to see how much control researchers can apply to the VAE by picking the z value that is used to reconstruct images. More experimentation into the GAN's latent space and choice of noise going into the generator could expand its capabilities. Researchers may be able to apply a similar process to the GAN that was done in this paper with attribute isolation in the VAE. With more control over z values the GAN may be able to generate better images with less training time.

237 **6. Conclusion**

238 In this paper, two different generative models are compared
239 based on the quality of their results. These two models are deep CNN-based Variational Autoencoder and deep
240 convolutional generative adversarial network (DCGAN).
241 Both of these models have advantages and disadvantages
242 that may make them better or worse for different use cases.
243 The deep CNN-based VAE gives more control over the latent
244 space of an image while the DCGAN allows potentially
245 better results and more original images.

246 **7. Citations and References**247 **References**

- 248 Hou, X., Shen, L., Sun, K., and Qiu, G. Deep feature
249 consistent variational autoencoder, 2016.
250 Radford, A., Metz, L., and Chintala, S. Unsupervised rep-
251 resentation learning with deep convolutional generative
252 adversarial networks, 2016.

253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274