**Imperial College**
**London**

IMPERIAL COLLEGE LONDON

INDEPENDENT STUDY OPTION

# Image Generation and Recognition (Emotions)

*Author:*
Hanne Carlsson

*Supervisor:*
Dimitrios Kollias

May 2019

# Abstract

*Generative Adversarial Networks (GANs) were proposed in 2014 by Goodfellow et al. [16], and have since been extended into multiple computer vision applications. This report provides a thorough survey of recent GAN research, outlining the various architectures and applications, as well as methods for training GANs and dealing with latent space. This is followed by a discussion of potential areas for future GAN research, including: evaluating GANs, better understanding GANs, and techniques for training GANs. The second part of this report outlines the compilation of a dataset of images 'in the wild' representing each of the 7 basic human emotions, and analyses experiments done when training a StarGAN [11] on this dataset combined with the FER2013 [15] dataset.*

# Contents

# Chapter 1

# Introduction

Generative Adversarial Networks (GANs) are a fairly recent addition to deep learning research, and have advanced generative image modeling dramatically [6] since they were proposed by Goodfellow et al. [16] in 2014. GANs are applicable for both semi-supervised and unsupervised learning tasks [12], and they have achieved impressive results in various image generation tasks, such as: image-to-image synthesis, text-to-image synthesis, and image super-resolution [69].

There are many different types of GANs (see Chapter 2), but they all consist of two simultaneously trained models: a generator and a discriminator. The generator captures the distribution of the data and tries to generate a plausible image from that distribution, while the discriminator estimates the probability of the generated image being from the training data rather than made by the generator [16]. The two are in competition with each other, with the generator learning to fool the discriminator and the discriminator learning to tell real and fake images apart. Formally, the objective of GANs is to find the Nash equilibrium to the following two player minimax problem with value function $V(D, G)$ [16]:

$$\min_G \max_D V(D, G) = \mathbb{E}_{\boldsymbol{x} \sim p_{data}(\boldsymbol{x})}[\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})}[\log(1 - D(G(\boldsymbol{z})))] \qquad (1.1)$$

The generator defines a probability distribution $p_g$ of samples $G(\boldsymbol{z})$ obtained from a noise distribution $\boldsymbol{z} \sim p_{\boldsymbol{z}}$. Goodfellow et al. [16] proved that this minimax game has a global optimum for $p_g = p_{data}$, when the discriminator is unable to differentiate between the training data and generator distributions. The generator never sees the real training images, and instead learns the distribution of them through its interaction with the discriminator [12].

## 1.1. Project Aim

This project aimed to explore current GAN research and understand the trends and advancements made since their proposal in 2014. The first part of this project involved conducting an extensive literature review of GAN research to date (see Chapters 2-5), and then discussing the apparent research trends in the context [35] of open areas of research for future work (see Chapter 6).

The second part of this project involved implementing and applying GANs. First, a dataset of images 'in the wild' displaying the 7 basic emotions was compiled (see Chapter 7), and then these were combined with the FER2013[15] dataset to train a StarGAN [11] for image-to-image translation of emotions (see Chapter 8).

# Chapter 2

# GAN Architectures

Since GANs were introduced in 2014, many variants have been proposed for different applications. This chapter aims to walk through some of the key differences in various GAN architectures.

## 2.1. Conditional GANs

Conditional GANs were introduced by Mirza et al. as a way to condition the model and be able to direct what it generates [48]. This is done by feeding the data to condition on to both the generator and discriminator as an additional input layer [12] (see Figure 2.1). An advantage of using conditional GANs is that it enables better representations for one-to-many mappings, meaning conditioning on a single class (eg. Dog) can generate multiple types of dogs with various colours and features.

The Information Maximising GAN (InfoGAN) proposed by Chen et al. [9] in 2016 is a completely unsupervised adaptation of the Conditional GAN where the discriminator estimates both if an image is real or fake, and what its class label is [12] (see Figure 2.1). By adding an information regularisation term to the GAN minimax equation (Equation 1.1), the InfoGAN enforces high mutual information[1] between the condition and the generator distribution [9]. Despite being able to successfully learn interpretable representations on difficult datasets (eg. CelebA [44]), the Info-GAN barely adds any computational cost to the original GAN.

In 2018 a projection based approach to condition information in the discriminator was proposed by Miyato et al. [50]. Instead of concatenating the conditioning label with the training data [48], their paper instead proposes taking the inner product between the embedded condition vector and the feature vector [50]. An advantage of using this projection method in the discriminator is that it does not suffer from mode collapse[2], like previous conditional methods do, and it helps increase the diversity within each class.

---

[1]Mutual information measures the amount of information that is learned from one random variable $Y$ about another random variable $X$: $\mathbb{I}(X;Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$ [9].

[2]Mode collapse is when the discriminator stops causing the generator outputs to be dissimilar, resulting in the generator only outputting very similar looking images, with no diversity [59].
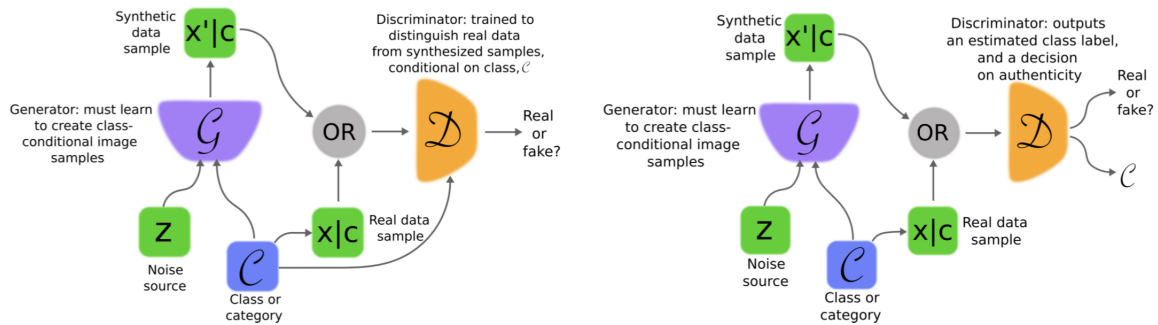
**Figure 2.1:** The image on the left is the conditional GAN proposed by Mirza et al. [48], and the image on the right is the InfoGAN proposed by Chen et al. [9]. The image is from the GAN survey paper by Creswell et al. [12].

Multiple GAN variants make use of conditional GANs as part of their overall architecture, some of these include: Text to Image Synthesis [57], GAWWN [56], pix2pix [22], StackGAN [70], and StarGAN [11].

## 2.2. Convolutional GANs

Convolutional neural networks (CNNs) are very well suited to solving image related problems, requiring fewer parameters than a normal feed-forward neural network [27] and therefore being much faster to train [42]. Convolutional GANs are able to make use of the advantages of CNNs in GANs.

In 2015 Radford et al. [54] proposed the Deep Convolutional GAN (DCGAN), where both the discriminator and generator are deep CNNs. Their approach uses the all convolutional net [62] in both the generator and discriminator which replaces pooling with strided convolutions, allowing down-sampling and up-sampling to be learned during training [12]. They also apply batch normalisation [21] everywhere (except for the generator output layer and discriminator input layer) which stabilises training.

The 3D-GAN proposed by Wu et al. [67] in 2016 also makes use of an all convolutional net [62] in both its generator and discriminator, synthesising 3D images of objects using volumetric convolutions [12]. The generator consists of 5 volumetric fully convolutional layers, and given a 200 dimensional noise vector $z$ it outputs a $64 \times 64 \times 64$ dimensional image of an object. Their implementation did not make use of conditional GANs and instead they trained one 3D-GAN for each class of objects they were interested in (eg. chairs, or tables).

A few months later in 2016 an approach for image-to-image translation, called pix2pix, was proposed by Isola et al. [22]. The pix2pix model is able to generate photo-realistic images from label maps, reconstruct objects from edge maps, and add colour to black and white images [22]. This model uses a U-net[58] architecture for the generator, in order to avoid the bottleneck caused by using an encoder-decoder
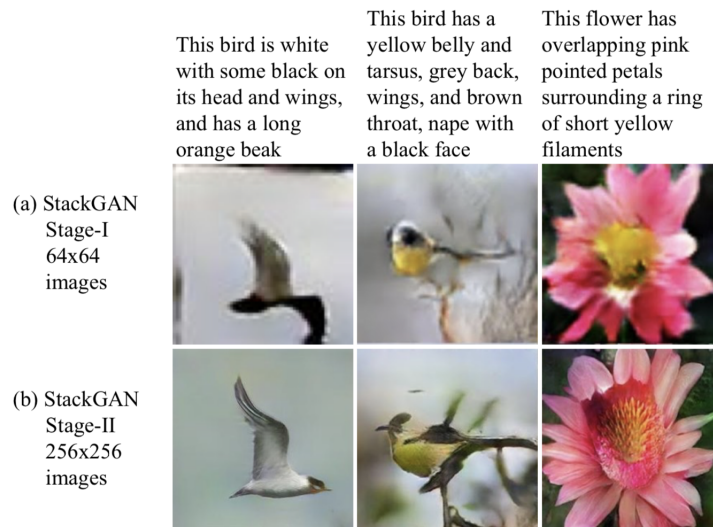
**Figure 2.2:** The $64 \times 64$ images produced by the Stage-I GAN, and the final $256 \times 256$ images produced by the Stage-II GAN (and therefore full StackGAN) given the various text descriptions on the top. This image is from the StackGAN paper [70].

network, and a convolutional PatchGAN[43] classifier for the discriminator, discriminating on a patch-level rather than on an image-level [22]. They use a $70 \times 70$ PatchGAN to force sharp photo-realistic outputs, after finding that increasing the patch size beyond this lowers the overall output quality. Their framework is not application specific, but instead makes use of a conditional GAN architecture to be able to deal with numerous image-to-image translation problems.

Numerous other GAN variants make use of convolutional networks as part of their architecture, some of these include: InfoGAN [9], GAWWN [56], StackGAN [70], CycleGAN [72], Progressive Growing of GANs [23], StarGAN [11], BicycleGAN [73], SAGAN [69], BigGAN [6], and StyleGAN [24].

## 2.3. StackGAN

The StackGAN, proposed by Zhang et al. in 2016 [70], is a two-stage GAN that generates photo-realistic images from text descriptions (see Figure 2.2). By decomposing this difficult problem into two easier 'stages'[3], they were able to significantly improve on previous state-of-the-art methods. The first stage (Stage-I GAN) sketches the general outline of what is described in the text, and adds the main background and object colours. The second stage (Stage-II GAN) then takes as input both the low-resolution Stage-I output and the text description, focusing on generating the finer details in the image, making it more photo-realistic and true to the text description.
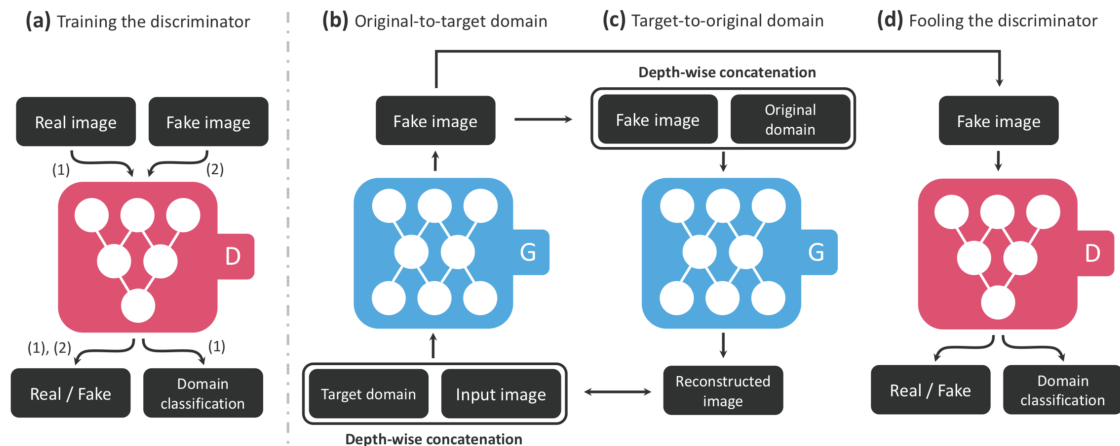
**Figure 2.3:** (a) Shows how the StarGAN discriminator is trained, and (b)-(d) show how the StarGAN generator is trained. The image is from the StarGAN paper [11].

## 2.4. Cyclical GANs

Cyclical GANs make use of a cycle consistency loss which enables translating from one domain to another, and then back again to the starting domain without loss of information. Although cycle consistency loss has been used before, the CycleGAN proposed by Zhu et al. [72] was the first time this was applied to GANs. The Cycle-GAN enables translating between two different image domains by training on two unordered image collections (one for each domain). This is possible because the CycleGAN assumes an underlying relationship between the two domains, making it possible to train without any paired training data.

A few months after the CycleGAN, the BicycleGAN was proposed by Zhu et al. [73]. The BicycleGAN model maps input to an output distribution, rather than a deterministic output (like conditional GANs do in previous image-to-image translation models). The BicycleGAN is a combination of a Conditional Variational Autoencoder GAN (cVAE-GAN), which encodes the original input image into latent space, and a Conditional Latent Regressor GAN (cLR-GAN), which gives the generator a randomly drawn latent vector which looks realistic (but not necessarily like the input image) and then uses an encoder to try and recover the latent vector [73]. This combination of conditional GANs enables the BicycleGAN to produce outputs that are both photo-realistic and diverse in nature.

The StarGAN, proposed by Choi et al. [11], is unique in that it enables training image-to-image translations between multiple domains with only a single generator and discriminator — unlike previous image-to-image translation models which can only learn relationships between two domains at a time [11]. It only needs one generator because it uses it twice: first to translate the input image to the target domain, then using a cycle consistency loss it reconstructs the input image from the

---

[3]During training they iteratively train Stage-I GAN for 600 epochs while fixing the Stage-II GAN, and then they fix the Stage-I GAN to train the Stage-II GAN for another 600 epochs [70]

translated image [11] (see Figure 2.3). StarGAN can train on multiple datasets with different label sets, because it uses a mask vector[4] which enables the model to ignore unspecified labels.

## 2.5. Self-Attention GAN

The Self-Attention GAN (SAGAN) proposed by Zhang et al. [69] in 2018 uses self-attention and long-range dependency modelling to generate images where the objects and scenarios within it are related in a way consistent with realistic images. Self-attention [10, 52] (or intra-attention) finds the relationship between different parts of a sequence in order to represent and understand each part of the sequence. The SAGAN is the first time self-attention has been applied to GANs. Most GAN models use convolutional layers as part of their architecture (see Section 2.2) since they are very good at modelling local dependencies, however CNNs cannot efficiently model long-range dependencies on their own. By applying self-attention to both the generator and the discriminator the SAGAN is able to effectively model both local and global dependencies in an image, ensuring that highly detailed features in various parts of the generated image are consistent [69].

---

[4]A mask vector is a $n$-dimensional one-hot vector, where $n$ is the number of datasets.

# Chapter 3

# Training GANs

In the paper where GANs were first introduced [16], the generator network used both Rectified Linear Units (ReLU) [14, pg.189] and Sigmoid [51] as activation functions, and the discriminator used maxout [17] activations. Dropout [63] was also applied when training the discriminator [16]. The goal of Goodfellow et al. was never to find the ideal training parameters, however, but rather to introduce GANs and show the theory behind the optimal solution.

Despite an optimal solution theoretically existing, finding it in practice can be very difficult. The main reason for this is because for the original GAN the optimum is a saddle point, rather than a minimum like most other machine learning models. Due to these difficulties, some of the problems which often occur when training GANs are: struggling to have both generator and discriminator converge, mode collapse, and the discriminator converging to zero which gives the generator no useful gradient updates [12]. This chapter will discuss some of the findings in later papers that focused on overcoming these issues and improving the GAN training process.

## 3.1. Training Tricks

The first paper to suggest techniques to improve the training process was the DC-GAN paper [54], which outlined guidelines for training and constructing both the generator and discriminator. The specific architecture of the DCGAN is outlined in Section 2.2, but the addition most related to the training of GANs was the application of batch normalisation throughout the model which helped stabilise learning in deeper networks [12]. Radford et al. also showed that using Leaky ReLU [45] in the discriminator, and using tanh [14, pg.191] for the output layer of the generator improved performance [54]. Using tanh allowed the model to learn to saturate quicker and cover the full colour space of the training distribution [54].

The DCGAN was trained with mini-batch stochastic gradient descent (SGD) with a mini-batch size of 128, and its weights were initialised from a Gaussian with $\mu = 0$, $\sigma = 0.02$ [54]. They also used the Adam [25] optimiser, which became common practice in GAN training after the release of this paper.

A few months later, Salimans et al. [59] released a paper proposing further guide-

lines on stabilising GAN training [12]. They observed that gradient descent was not always sufficient when searching for the Nash equilibrium, and introduced five techniques to encourage convergence: [59]:

1. **Feature matching**: making the objective of the generator be to match the expected value of intermediate discriminator layers. This is effective when the GAN is unstable during training, but there is no guarantee of it working in practice.

2. **Minibatch discrimination**: deals with mode collapse by having the discriminator look at a mini-batch of examples rather than a single example, enabling it to tell if the generator is producing the same outputs. This method was found to be superior to feature matching.

3. **Historical averaging**: keeps a historical average of parameters in order to penalise parameters that drastically differ from the average.

4. **One-sided label smoothing**: changes the discriminator target from 1 to $0.9$ to stop it from becoming overly confident and providing weak gradients.

5. **Virtual batch normalisation**: uses the advantages of batch normalisation found by Radford et al. [54], but normalises each example based on a reference batch fixed at the start of training. This is very computationally expensive, and was therefore only applied to the generator network.

Isola et al. [22] suggested an alternative method to optimising training by elaborating on the original method by Goodfellow et al. [16] of alternating between updating the generator and discriminator after each iteration, but when optimising the discriminator they divide the objective by 2 which slows down the rate at which the discriminator learns compared to the generator [22]. They also, quite unusually, apply dropout to the generator network at test time, and they apply batch normalisation with the statistics of their test batch rather than a training batch.

## 3.2. Variations in GAN Training Techniques

The Wasserstein GAN (WGAN) [2], proposed by Arjovsky et al., uses a different cost function to previous GANs, namely one derived from an approximation of the Wasserstein distance. The advantages of using this cost function are that it provides more meaningful gradients for updating the generator, the WGAN does not suffer from mode collapse, and it makes the GAN easier to train by avoiding the vanishing gradient problem that other GANs suffer from [12]. Another main advantage of the WGAN compared to the original GAN is that the optimum is a minimum rather than a saddle point, like previous GANs. Arjovsky et al. found that using the Adam optimiser did not work for this cost function as it made training very unstable, so instead they used the RMSProp [66] optimiser, and decreased the initial learning rate, which helped stabilise training [2]. Training WGANs was improved further by Gulrajani et al. [19] with the proposal of gradient penalty (WGAN-GP) which helped alleviate some of the issues caused by weight clipping in the original WGAN.
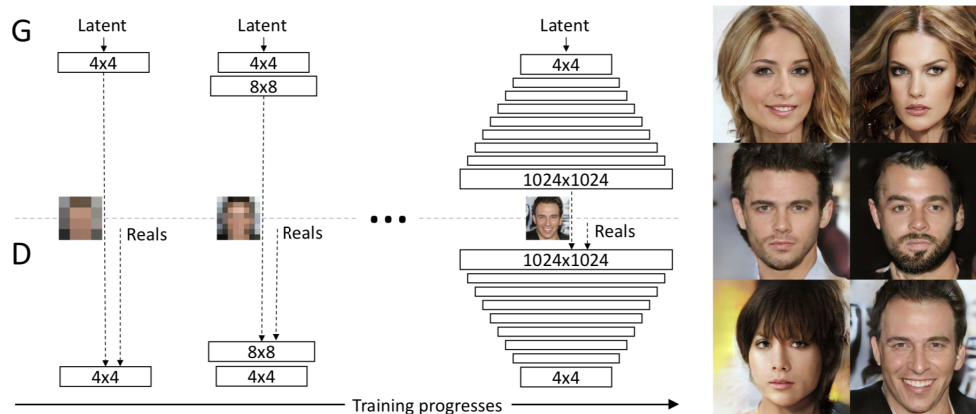
**Figure 3.1:** During training the progressively growing GAN begins with the generator and discriminator both having a spatial resolution of $4 \times 4$. As training advances, they keep adding layers to both networks to increase the resolution. This image is from the paper where this method was introduced [23].

Karras et al. [23] also introduced a new way of training GANs. They progressively increased the resolution of images by adding layers to the generator and discriminator during training (see Figure 3.1), which resulted in more stable training since the majority of training time was spent on smaller low-resolution images that are easier to train on [23]. They also adapt the suggestion by Salimans et al. [59] to use minibatch discrimination, resulting in a much simpler approach [23]:

1. Compute standard deviation for every feature (in each spatial location) over the minibatch.

2. Average estimates over all features and locations.

3. Replicate and concatenate the value to all spatial locations over the minibatch, resulting in an additional feature map. This layer can be inserted anywhere in the discriminator but they achieved the best results when inserting it towards the end of the discriminator network.

Another aspect of training that differs for the progressively growing GANs is that instead of using batch normalisation they initialise their weights from $\mathcal{N}(0, 1)$ and then scale them at runtime, and they normalise each pixel's feature vector to unit length in the generator after each convolutional layer [23].

A year later (2018) Zhang et al. proposed two techniques to help stabilise training GANs on challenging datasets, which they used for their SAGAN [69]. Firstly, they applied spectral normalisation [49] to both the generator and discriminator, which allowed for fewer discriminator updates per generator update and reduced the computational training cost [69]. Secondly, they used different learning rates for the generator and discriminator, enabling them to get better results in the same wall-clock time[1] [69].

---

[1]Wall-clock time is the amount of time it takes to complete a job, which in this case refers to the time it takes to train the GAN.

Using the SAGAN model as a baseline, the BigGAN proposed by Brock et al. [6] increased the batch size, width, and depth of the network to train GANs at the largest scale attempted to date [6]. They found that increasing batch size and width improved performance, presumably because it allowed each batch to cover more modes and thereby provide better gradients, and because increasing the network width increases the number of parameters and thereby the model capacity [6]. Increasing the depth of the network, however, did not improve performance until they adapted their model to a BigGAN-deep variant which used a slightly different residual block structure than BigGAN [6]. Scaling up by this amount significantly decreased training stability, which they dealt with in the following ways:

- Use a shared embedding (for class embeddings used in conditional batch normalisation layers of the generator) rather than a layer for each embedding. This decreases overall training time and computation.

- Add skip-$z$ connections (see Section 4.2).

- Truncation trick (see Section 4.1).

- Orthogonal regularisation (see Section 4.1).

- Apply early stopping to deal with training collapse.

Brock et al. found that training stability is not dependent on just the generator or the discriminator, but on their interaction throughout training. They also find that increasing training stability by adding regularisation to the discriminator comes at the cost of decreasing the overall training performance, and results in the discriminator overfitting to training data [6]. The solution which offered the best trade-off was allowing for training collapse later on in training and implementing early stopping when that occurs.

# Chapter 4

# Structure of Latent Space

The generator in all (vanilla) GAN models take a randomly sampled vector from the latent space as input and map it to the domain which it tries to model [12]. The latent space has fewer dimensions than the domain space but represents the semantic [26] structure of the space, in a similar way to word2vec [47] which represents the semantics of words in $n$[1]-dimensional space.

A way of exploring this semantic space [30] is to show arithmetic in it. In the word2vec space it is possible to 'add' and 'subtract' meanings from words in a way that makes sense, eg. king$-$man$+$woman$=$queen. Radford et al. [54] applied vector arithmetic on face samples, adding/subtracting things like glasses and emotions [39], but they needed to use the average of three latent vectors in order to obtain stable results. Wu et al. [67] showed greater progress with latent space arithmetic by doing shape arithmetic for their 3D-GAN. This allowed them to add or subtract 'arm' or 'leg' vectors from 'table' or 'chair' outputs (see Figure 4.1), thus demonstrating the learned representations of the 3D-GAN output [67].

The way the latent space has been modelled and fed to the generator varies between different GAN implementations, depending on the overall goal of the GAN model. In the AlignDRAW GAN proposed by Mansimov et al. [46] (which generates images from captions) the latent variables are randomly sampled from a normal distribution where the mean and variance are dependent on the previous hidden states of the generative LSTM [20], rather than drawn from a $\mathcal{N}(0, I)$ distribution [46]. For the InfoGAN model, Chen et al. [9] chose to decompose the noise vector into two separate parts: one treated as the noise, and one which targets the semantic features of the distribution [32]. In the pix2pix model, Isola et al. [22] do not provide any noise to the generator, and instead apply dropout to several layers throughout the generator network at both training and test time to act as noise in the network.

This chapter will cover some of the methods attempted when dealing with latent space, and the ways latent code is added to generator networks.

---

[1]The dimensionality of word2vec vectors vary depending on how they were trained and what they will be used for, the most common dimensions used are 50, 100, 300, or 500.

**Figure 4.1:** Examples of the kinds of shape arithmetic that can be performed on the latent space of the 3D-GAN. The image is from the 3D-GAN paper [67].

## 4.1. Adapting the Latent Space

Zhang et al. [70] proposed using *Conditioning Augmentation* for their text-to-image StackGAN. This was an attempt to mitigate previous discontinuity issues in the latent space caused by having limited amounts of data when transforming text embeddings into the latent space [70]. Their proposed method randomly samples additional conditioning variables from a Gaussian distribution with a mean and covariance matrix which are both functions of the text embedding. This mechanism helped produce more training pairs than before from the same size dataset.

Zhu et al. [73] investigated how dimensionality of the latent space can impact the overall model. They found that having too low of a dimensional space may limit the amount of diversity the latent space can represent, whereas having a very high-dimensional space can make sampling difficult since too much information will be encoded in the space. The optimal dimensionality will be dependent on the intended application of the model and the individual dataset used [73].

The *Truncation Trick* was proposed by Brock et al. [6] for the BigGAN, and is also used by Karras et al. [24] for the StyleGAN (although only for low resolutions to avoid affecting high resolution details). They truncate their $z$ vector by a threshold, and then every value with a magnitude above the chosen threshold is re-sampled. As the threshold value decreases, the Inception Score (IS[2]) increases, and the Fréchet Inception Distance (FID[3]) initially improves, but then as the threshold approaches zero it worsens to penalise the lack of variety (see Figure 4.2). Interestingly they obtain their best overall results when they train their model with a $\mathcal{N}(0, I)$ latent distribution and then sample with the Truncation Trick.

Another 'trick' applied by Brock et al. [6] is the use of *Orthogonal Regularisation* [7] to allow the full latent space to map to 'good' samples. This aims to minimise the cosine similarity between filters and results in adding smoothness to the models. They find that $60\%$ of their models are responsive to the Truncation Trick when they have been trained with Orthogonal Regularisation [6].

Karras et al. [24] applied *Mixing Regularisation* to their StyleGAN implementation,

---

[2]IS is a metric for evaluating GANs that does not penalise a lack of variety but instead rewards precision [6]. (The higher the score the better)

[3]FID is another commonly used metric for evaluating GANs. It penalises a lack of variety in generated images, but also rewards precision [6]. (The lower the distance the better)

**Figure 4.2:** These images show the effects of decreasing the threshold of the Truncation Trick on the BigGAN. From left to right the threshold is 2, 1, 0.5, 0.04. This image is from the BigGAN paper [6].



**Figure 4.3:** Examples of images generated by applying mixing regularisation and mixing two latent codes at coarse spatial resolutions during training. Aspects such as pose, hair style, and face shape are copied from source B, while the colours and finer facial features are from source A. This image is from the StyleGAN paper [24].

forcing some images to be generated using two random latent codes instead of just one during training (see Figure 4.3). If two latent codes are used, then the generator will randomly switch between the two throughout the network. This stops the network from assuming that adjacent styles are correlated.

## 4.2. Injecting Latent Code

Zhu et al. [73] explored two ways of injecting latent code into the generator for their BicycleGAN model: injecting it into the generator's input layer, and injecting it into

every intermediate layer of the network [73]. They obtained similar performance from both methods, showing that for their model it did not make much difference where the latent code was added.

Brock et al. [6] also inject latent code into multiple intermediate layers of the generator (skip-$z$ connections). They found that for BigGAN and BigGAN-deep this only boosted the overall accuracy by around $4\%$, but that it improved the training speed by a further $18\%$ [6]. They suggest that this is because skip-$z$ connections allow the latent space to influence features at different resolutions throughout the generator network [6].

# Chapter 5

# Applications of GANs

Along with the general advancements of GAN research, discovering new applications for GANs is another active research field [12]. The application and goal of a GAN model, as mentioned numerous times throughout this report, has a major impact on its general architecture and training set up. This chapter will highlight some of the broad categories GAN applications fall into and the type of GAN models that fall into each category.

## 5.1. Image Classification

Image classification is a key part of computer vision research [33, 18, 60, 8, 13, 64, 65], and being able to apply GANs to this is very useful. One way GANs can be used for image classification is as a way to quantitatively assess features extracted from unsupervised learning. Radford et al. [54] reuse the outputs of the convolutional layers in the discriminator (after training) as a feature extractor and assess their quality by applying a regularised L2-SVM classifier to each feature vector [12]. Due to the difficulty in classifying an image as real or fake (generated by a GAN) as GANs improve, using image classification is an important tool in judging the performance of a GAN [12].

## 5.2. Image Synthesis

One of the most important qualities of a GAN is its ability to generate (often photo-realistic) images [28, 29]. This becomes especially useful when the generated image can be conditioned on certain constraints relating to its features [12]. Conditioning was first introduced with the Conditional GAN, where the generated image was conditioned on a certain class [48]. This was expanded into 'text-to-image' synthesis, where the generated image was conditioned on a short text description.

The AlignDRAW [46] GAN attempted to solve this problem, but was unable to do so by training an end-to-end model. In 2016 Reed et al. [57] proposed the first end-to-end differentiable architecture which conditioned on text descriptions rather than class labels. Their results were not particularly photo-realistic and detailed, but they showed the possibilities of conditioning on natural language, and in the three years since their proposal text-to-image synthesis has made impressive improvements.

**Figure 5.1:** The positioning of the birds are determined by the keypoint coordinates, which enables shrinking, stretching, and translation of the bird in the image. Caption is the text input, GT is the Ground Truth, and the remainder of the images are the generated output based on the keypoint coordinates. This image is from the GAWWN paper [56].

Later that same year, Reed et al. [56] proposed the Generative Adversarial What-Where Network (GAWWN) which was able to condition on where specific parts of the text descriptions should be located in the image, using bounding boxes. They show that they are able to perform shrinking, stretching, and translation of the specific objects by altering the bounding box coordinates, without changing the text or noise variables the model is conditioned on (see Figure 5.1).

The StackGAN [70] also performed text-to-image synthesis, but split the problem into two 'stages', one producing a general low-resolution outline of the text description, and the next improving the resolution of the first image by adding high-resolution details (more details on the specific architecture are given in Section 2.3).

The quality of images generated by GANs has improved greatly since the GAN was introduced [16]. Very recently the BigGAN [6] was proposed, which produced incredibly high resolution and detailed photo-realistic images at a larger scale than previously attempted. This showed the potential possibilities of photo-realistic images being generated from GANs by scaling up.

## 5.3. Image-to-Image Translation

Another common application for GANs is image-to-image translation. This is when an image is given as input and the GAN translates it to a different domain. The first image-to-image translation model with a generic enough approach to be applied to multiple different problems was proposed by Isola et al. [22] and nicknamed pix2pix. Their model uses a conditional GAN and they show it works for numerous applications such as adding colour to a black and white image, translating an edge map to a photo-realistic image, translating a image between night and day, and creating a map from aerial photographs. These problems had previously required separate models trained for each specific domain translation [12].

This was later extended by the CycleGAN [72] which introduced a cycle-consistency loss enabling reverse translations without a loss of information. This cyclical ap-
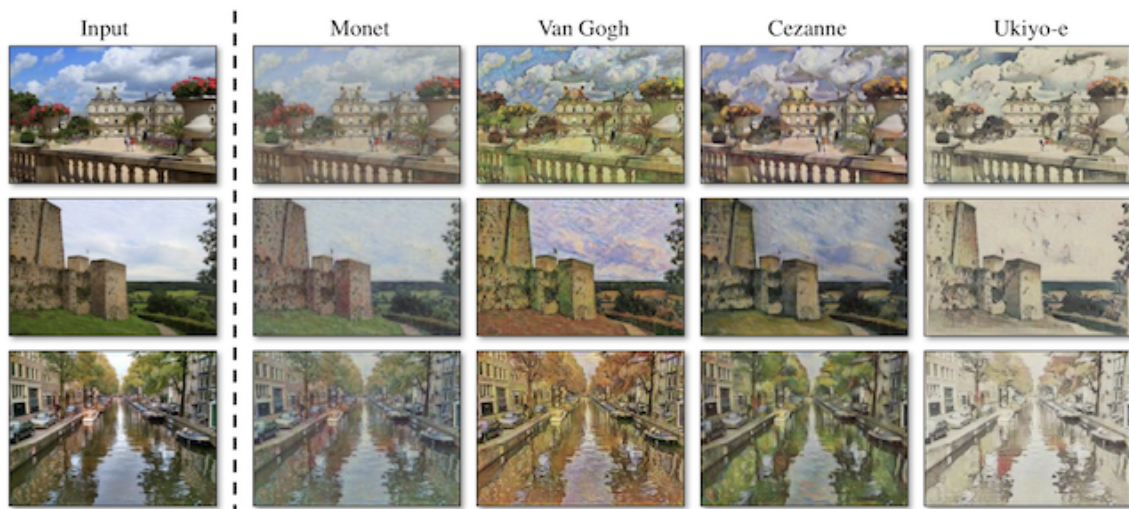
**Figure 5.2:** The CycleGAN is able to learn the characteristics of various domains, without training on any specific pairs between them, and use these characteristics to translate an input image into a different domain. Here an input photograph has been translated into the style of various famous painters. This image is taken from the CycleGAN paper [72].

proach meant that matching image pairs were no longer needed for training, and instead it was sufficient to train on two individual domains for the GAN to learn the characteristics of each well enough to translate between the two (see Figure 5.2). Other models which made use of this cycle-consistency loss for image-to-image translation were the BicycleGAN and the StarGAN, and the details of their architecture are discussed in Section 2.4.

# Chapter 6

# Discussion on GANs

This chapter will discuss the methods commonly used to evaluate GANs, and the advantages and (more importantly) disadvantages of using these, as well as the areas of future work within GAN research where there are still unanswered questions and room for improvement.

## 6.1. Evaluating GANs

As GANs output images, an obvious way of evaluating GANs is to look at the outputs and judge how realistic and reasonable they seem. The issue with this is it is hard to compare various GAN models, or even alterations of the same model, without having a way to quantify how 'good' a given model is for a certain task. To date there is no overall metric used to compare all GANs, but there are a few main ones that have been used to compare models with the same applications.

Image-to-Image models are often quantitatively evaluated using the FCN score. FCN stands for Fully-Convolutional Network, and the FCN score uses a classifier to perform semantic segmentation on the output images, and compares this to a ground truth labelled image. The intuition behind this metric is that if the GAN has successfully generated a realistic image, then a classifier should be able to label its semantic components successfully [72]. For Image-to-Image models that also perform photo-to-label translations (eg. CycleGAN [72]) metrics such as per-pixel accuracy, per-class accuracy, and mean class Intersection-Over-Union (Class IOU) are commonly used.

For conditional image synthesis GANs two commonly used metrics are Fréchet Inception Distance (FID) and Inception Score (IS). These metrics were briefly mentioned in Section 4.1, but essentially they are ways to quantitatively measure the sample quality of a model. They are not perfect measurements, but as they have been used to measure the effectiveness of multiple models, they are useful metrics in that they enable quantitative comparison between models [6].

In terms of human evaluation of GAN models, the most common way this is done is by using Amazon Mechanical Turk (AMT). This way the opinions of multiple people can be quantitatively used, which — although incredibly subjective — is a valid way

of getting a general feel for how good or bad a model is. The most common way AMT is used is that either a single picture is shown at a time and the person has to say whether they believe the image is real or fake (essentially acting as a human discriminator) or multiple pictures are shown and they have to select the real/fake one. These human evaluations are incredibly useful, but they are very hard to use when comparing two models (unless a human evaluation study has been completed to compare the two models).

## 6.2. Areas for Future Work

GAN research has made great strides since they were first introduced in 2014 by Goodfellow et al. [16], but they are still far from perfect and there are multiple areas where future work is needed. The main areas that will be discussed in this section are evaluating GANs, understanding the inner workings of GANs, and training instability.

Regarding the methods for evaluating GANs, as mentioned in Section 6.1, although multiple metrics exist all of them are flawed in some way. There is still no one good metric to successfully measure both how well a GAN is performing for its assigned task, and how good the quality of its output is, while also being able to compare multiple models with each other. This is an area of research that is still very open to future work, especially since the current method of using multiple different metrics easily leads to conflicting conclusions regarding the quality of models [12].

Another area for future work is regarding how and why GANs actually work. Part of this is better understanding the latent space and what information it entails. Some of the current methods for manipulating the latent space and using it to control GAN output were discussed in Chapter 4, but there are still a lot of unknowns regarding how the latent space is actually encoding the information[61]. Work was done by Bau et al. [4] to visualise and understand GANs and their inner workings, enabling them to compare the internal representations across various models and datasets as well as be able to manipulate objects within a generated scene. There are, however, still many unknown questions regarding GANs, such as: what is the relationship between layers of a GAN? how can we predict which relationships a GAN will and will not be able to learn?

Finally, as discussed in Chapter 3, training GANs is still a challenge, and they are very prone to instability and mode collapse. Training GANs requires searching for a Nash equilibrium of the minimax game (Equation 1.1), but sometimes no equilibrium exists and the GAN is unable to converge to the optimal solution. A lot of work has been done with regards to finding optimal conditions for training, but since they are application dependent, it is still very much an open problem to find ideal training conditions. Brock et al. [6] trained GANs at the largest scale yet, and found both that their model was very prone to instability during training, and that it achieved incredibly impressive photo-realistic results. This brings about some interesting questions: how far can they be scaled up while still showing a performance improvement? how much scaling up is too much? how much are we able to scale up given our current

computing resources?

Being able to better understand and control GANs during training is an area of GAN research that could drastically improve GAN performance as a whole.

# Chapter 7

# Emotion Dataset

As part of this project a database of images representing each of the 7 basic emotions (angry, disgust, fear, happy, neutral, sad, surprise) 'in the wild' was compiled. An image is referred to as 'in the wild' if it is not posed, and instead taken in a natural environment [68, 31, 34]. This type of database is quite invaluable in the field of Computer Vision[36, 38, 41], because most databases consist of perfectly posed pictures with a natural background, which makes it hard for models trained on those databases to generalise to the real world. By using images 'in the wild', or a mixture of posed and 'in the wild' images, models are able to better generalise to real world scenarios.

The initial aim was to gather enough images to be able to further categorise each emotion with its intensity, however the size of the final dataset was not large enough for this to seem reasonable. This chapter will detail the process of finding the images used for this database, as well as the pre-processing steps that were done before using the dataset in the Implementation (see Chapter 8).

## 7.1. Finding Images

The images for this dataset were creative commons licensed images collected from the following websites:

- https://pixabay.com/

- https://www.pexels.com/

- https://www.flickr.com/

- https://www.istockphoto.com/

- https://www.freeimages.com/

- https://burst.shopify.com/

The images were downloaded using the Google Chrome extension "Download All Images" offered by https://mobilefirst.me. This enabled multiple images to be downloaded at once, saving immense amounts of time.

Various keywords and phrases were used an in attempt to download as many images as possible for each emotion. This both meant using synonyms for each of the emotions, such as: enraged, annoyed, frustrated, furious, heated, or irritated when looking for images representing the 'Angry' class. But also using phrases like "Angry man", "Angry boy", "Angry woman", "Angry girl", and "Angry family". This meant that more images could be found and that the images collected were quite varied in terms of the ages and number of people present.

Some emotions were much easier to find images for than others, resulting in quite an in-balance for the final dataset. Compared to other emotion datasets, however, this is the norm since some emotions are simply more represented then others. The easiest emotion to find images for was Happy, and so that is the largest of the categories in the dataset with a total of $559$ images. Then Neutral, Sad, and Angry were all fairly easy to find images for and have a similar amount of representation in the dataset. The remaining emotions (Disgust, Fear, and Surprise), however, were very difficult to find good images for. Few of the websites available for creative commons images have particularly good search engines, and the use of keywords and short phrases had no significant impact when it came to these emotions – resulting in them all being quite under-represented in the final database.

The total number of images in this newly created database is $1,463$. As this is clearly not enough to use as a training set, it seemed a better idea to combine these images with another emotion dataset in order to have a more reasonable sized training set. The dataset chosen to combine this newly created dataset with was FER2013 [15], containing a total of $35,887$ images. Despite the overall size of the dataset, it only contains $547$ images for the Disgust class, showing how common the under-representation of that more unusual emotion is.

The images in the FER2013 dataset have the same 7 emotion categories as the dataset created throughout this project, and the images are grey-scale crops of faces with dimension $48 \times 48$. In order to combine the two datasets, some pre-processing was required (Section 7.2).

## 7.2. Pre-Processing

The new emotion dataset contained images with various resolutions and number of people per picture, so in order to combine this with the FER2013 dataset, the faces needed to be cropped from each image. This was done by using a Tensorflow [1] implementation of multitask cascaded convolutional networks (MTCNN) for face detection [3] and alignment [71]. This code was slightly modified to crop each image around the bounding box detected for each face, and then run on every single image in the collected emotion dataset.

Next, each of the cropped faces were manually inspected to remove any anomalies such as non-faces or re-categorise faces from the background of an image with a different emotion than that of the class it was originally placed in. Due to the vast amount of cropped images, this was a tedious process and prone to human error.
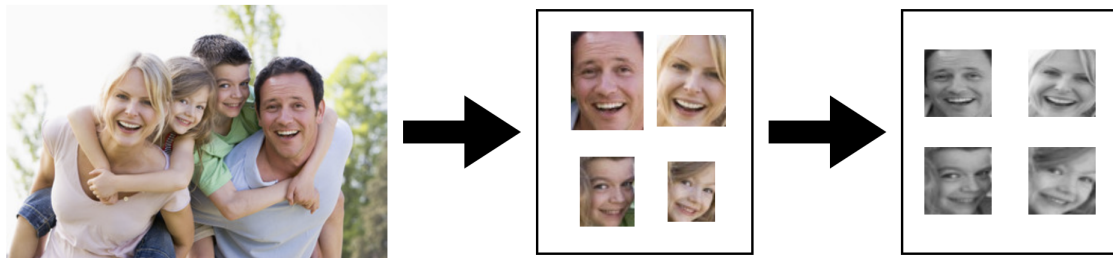
**Figure 7.1:** The steps taken to pre-process images from the newly created emotion dataset. First the MTCNN face detector was used to detect all the faces in an image, and crop each one out. Next, simple OpenCV functions `resize()` and `cvtColor()` were used to resize the cropped images into $48 \times 48$ and make them grayscale.

Finally, the OpenCV [5] functions `resize()` and `cvtColor()` were applied to each cropped image to ensure they were $48 \times 48$ in size and grayscale. Figure 7.1 outlines the brief steps of this pre-processing pipeline. The final images were then combined with the FER2013 images and used in the implementation of the StarGAN, discussed in Chapter 8.

# Chapter 8

# Implementation

For the implementation part of this project, the StarGAN [11] was chosen as it has proven to be very effective with image-to-image emotion generation, and because the authors' PyTorch [53] code base is surprisingly well commented and easy to understand compared to other official GAN implementations. As mentioned in Chapter 7, the emotion dataset compiled throughout this project was not large enough to be used as a training set on its own, so instead it was pre-processed and combined with the FER2013 [15] dataset. The total size of the combined training data was $30,100$ images, split into folders for each of the 7 basic emotions.

The original StarGAN implementation was trained on images of size $256 \times 256$, so in order for the code to run for the significantly smaller $48 \times 48$ images a few modifications were made at the command line. The default for the number of convolutional filters in the first layer of the discriminator and generator was $64$, so for this implementation that was decreased to $32$ to avoid a PyTorch division by zero error. The default number of residual blocks in the generator and strided convolutional layers in the discriminator was set to $6$, however this still caused a PyTorch division by zero error so for this implementation the number was varied between $3$ and $5$.

The Imperial Department of Computing GPU cluster (`firecrest.doc.ic.ac.uk`) was used to run the jobs for training which significantly decreased training time. On average each experiment took around 2 hours to train, partly due to the small image size, and partly thanks to the GPU cluster.

## 8.1. Analysing Results

The combined dataset was trained for the default $200,000$ iterations[1] for each of the models in this experiment. It was trained on models with 3, 4, and 5 layers[2], and a few of the resulting sample images generated can be seen in Figure 8.1. Due to the bad results obtained from the 3 layer model it was run for 200,000 iterations longer to see if extending the training time would help the model converge. All this

---

[1] An iteration for this implementation is the total number of iterations for training the discriminator, with 5 discriminator updates for every 1 generator update.

[2] A layer here refers to the number of residual blocks in the generator and the number of strided convolutional layers in the discriminator.

accomplished, however, was creating a clearer face outline in the generated image, but it still did not manage to generate photo-realistic results.

Looking at the generated sample images it becomes evident that the greater the number of layers of the model the more photo-realistic the generated images are. However, although the images generated by the 5 layer model are the most realistic, very few actually resemble the various emotions they are aiming to portray, and instead all 7 images look almost identical. The third row of the 5 layer model comes the closest to matching the emotions with a clear change to the mouth for happy and eyes/mouth for surprised. Overall though, none of the attempted models for this manage to generate images that are both photo-realistic and reflect each of the 7 emotions.

Although the images are not particularly impressive, they do show that the StarGAN has indeed learned things about the 7 basic emotions and certain characteristics they have in common. Even for the quite terrifying and non-realistic faces generated by the 3 layer model there is a clear o-shaped mouth added to all the surprised images indicating that the GAN has learned that the most important feature indicating surprise is that the mouth makes an o-shape.

A reason why the GAN seems to struggle with generating images could partly be due to the 'in the wild' nature of them. When the GAN with such limited training data is trying to learn important features indicating each emotion and figure out which part of the face is the mouth/nose/eyes simultaneously, its performance suffers. The fact that many images show a hand over the mouth, or a face turned sideways makes it nearly impossible for the GAN to even locate where the facial features [37, 40] are for it to modify, resulting in the non-facelike images generated by the 3 layer model.

Another reason why this GAN struggles with training is most likely due to the image size. When the images are smaller, specific details become harder to notice in an image as a single pixel represents more than it would if the same image was larger. This means that many of the fine details regarding how a mouth looks at the corners or how eyes are placed and shaped on a persons face become a lot harder for the GAN to learn. This makes it almost impossible for the GAN to successfully generate images that look human-like as it has not been able to properly learn what that looks like.

An obvious obstacle is also the fact that these models are trained on a combination of two datasets, one much smaller than the other and possibly acting more as noise [55] than as a contribution to the larger one (FER2013). Due to this the better models (4 and 5 layers) were also trained on the FER2013 dataset alone (see Figure 8.2).

As the sample images show, when the models were trained on solely FER2013 they performed much better. The generated images are still far from perfect, but looking at the first row of generated images for the 5 layer model shows 7 photo-realistic images with decent resemblance to their conditioned emotions. The input image on

the first row (of the 5 layer model) definitely makes the generation task easier by not having the face turned sideways or having anything else to distract in the image (unlike the image in row 2 where the face is rotated and a hand covers part of it, and as a result none of the emotions are evident). This confirms the suspicion that combining the two datasets made the new dataset act more like noise than a helpful addition to the training data, hindering the model from learning crucial features of each emotion.

Another clear observation is that increasing the number of layers improves the overall performance. Due to the small image size the layers cannot be increased much further without reducing the filter size in the first layer of the discriminator and generator. If future work were to be conducted on this project it would be interesting to see the effects of decreasing the filter size and increasing the number of layers, and whether this would continue the trend of improving the overall performance. It would also be interesting to increase the overall image size to $96 \times 96$ or $128 \times 128$ as this would most likely help improve performance as well.
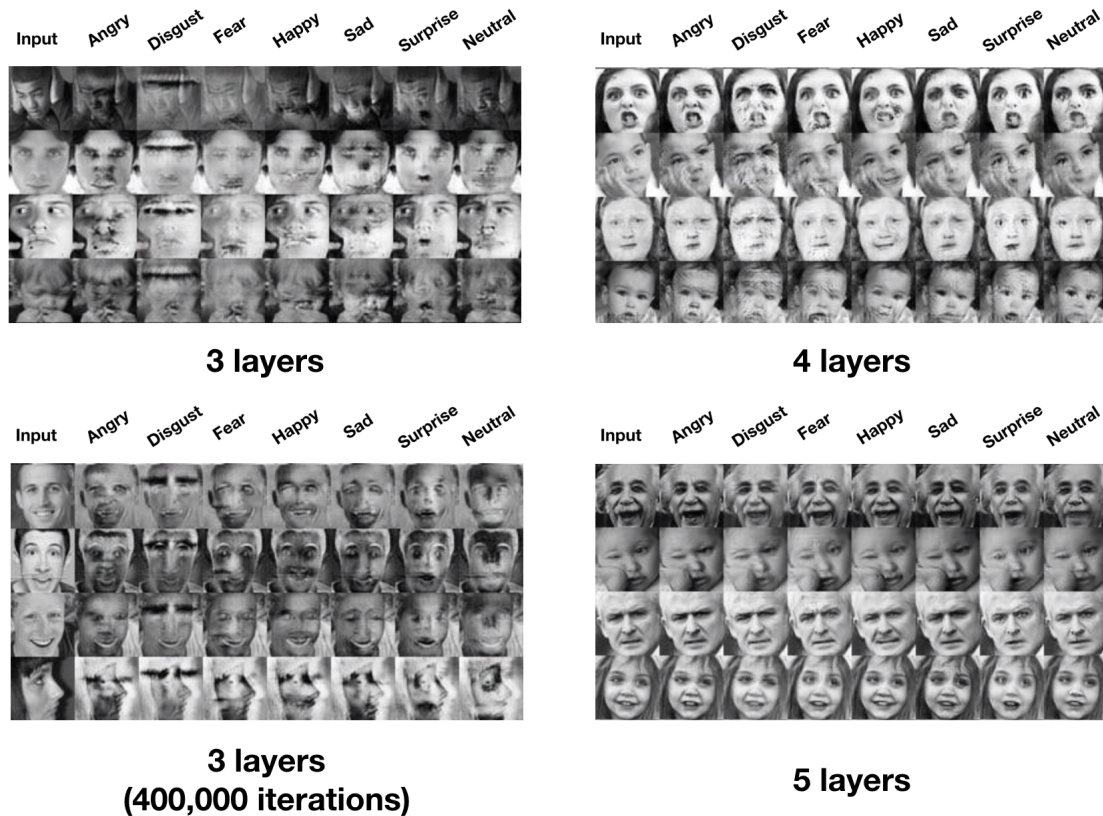
**Figure 8.1:** A few of the sample images generated after training the StarGAN with various numbers of layers on the FER2013 dataset combined with the emotion dataset created for this project. The number of layers under each group of images refers to the number of residual blocks in the generator and strided convolutional layers in the discriminator (these are set to the same). Each of these sample images were collected after 200,000 iterations, other than the bottom left one where the images were collected after 400,000 iterations.
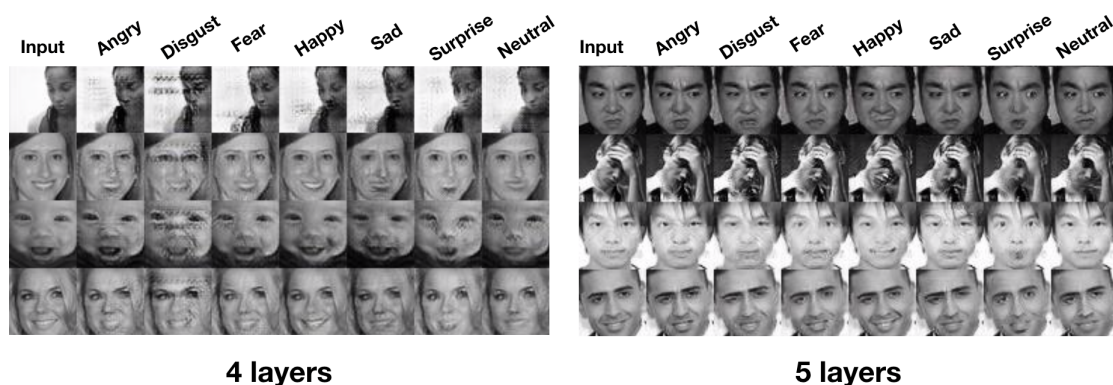


**Figure 8.2:** A few of the sample images generated after training the StarGAN with various numbers of layers on the FER2013 dataset. The number of layers under each group of images refers to the number of residual blocks in the generator and strided convolutional layers in the discriminator (these are set to the same). Each of these sample images were collected after 200,000 iterations.

# Chapter 9

# Conclusion

The main goal of this project was to undertake a survey of GAN research to date and determine possible areas for future research. This was discussed thoroughly in Chapter 6, with the main suggestions for future work being:

- **Evaluating GANs** - there is still no one recommended method for evaluating the overall performance of GANs, making it increasingly difficult to compare the performance between various architectures.

- **Understanding GANs** - how and why GANs work and what they actually learn remains an open question. Although some questions were answered by Bau et al. [4], many still remain.

- **Techniques for training GANs** - although some progress has been made with regards to training GANs in a more stable manner (for example by introducing the Wasserstein distance as a loss function [2]), stability in training remains a large problem for GANs. This was amplified by recent proposals such as the BigGAN [6] where scaling up introduced numerous training instability issues.

These areas for future work suggest that the better we can understand how GANs work, and comparatively measure them against each other, the better future GAN models will perform.

The second part of this project entailed compiling a dataset of emotion images 'in the wild' of all the 7 basic emotions. This was done successfully with a resulting database of $1,463$ images. These images were combined with the FER2013 images and used to train a StarGAN to generate emotions on images of faces. As the results from this implementation (Chapter 8) showed this was not particularly successful. Although the StarGAN seemed to learn certain features about the emotions, such as an o-shaped mouth being common for surprise, it was often unable to generate photo-realistic results of the various emotions.

This was partly to do with the emotion dataset created during this project acting more as noise than a helpful addition to the FER2013 dataset, and partly due to the images being too small. The size of the images meant that each pixel contained more details than it would in a larger image, making it much more difficult for the GAN to

distinguish smaller details in the image and thereby recreate them in a realistic way. When training the StarGAN on only the FER2013 dataset this helped solve the issue of the new dataset acting as noise, but as the image size remained the same, it was still not able to learn the finer details of the images.

# Bibliography

[1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org. 23

[2] Martín Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein GAN. *CoRR*, abs/1701.07875, 2017. 9, 29

[3] Yannis Avrithis, Nicolas Tsapatsoulis, and Stefanos Kollias. Broadcast news parsing using visual cues: A robust face detection approach. In *2000 IEEE International Conference on Multimedia and Expo. ICME2000. Proceedings. Latest Advances in the Fast Changing World of Multimedia (Cat. No. 00TH8532)*, volume 3, pages 1469–1472. IEEE, 2000. 23

[4] David Bau, Jun-Yan Zhu, Hendrik Strobelt, Bolei Zhou, Joshua B. Tenenbaum, William T. Freeman, and Antonio Torralba. GAN dissection: Visualizing and understanding generative adversarial networks. *CoRR*, abs/1811.10597, 2018. 20, 29

[5] G. Bradski. The OpenCV Library. *Dr. Dobb's Journal of Software Tools*, 2000. 24

[6] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. *CoRR*, abs/1809.11096, 2018. 1, 5, 11, 13, 14, 15, 17, 19, 20, 29

[7] Andrew Brock, Theodore Lim, James M. Ritchie, and Nick Weston. Neural photo editing with introspective adversarial networks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, 2017. 13

[8] G Caridakis, A Raouzaiou, K Karpouzis, and S Kollias. Synthesizing gesture expressivity based on real sequences. In *Workshop on Multimodal Corpora. From Multimodal Behaviour Theories to Usable Models. 5th International Conference on Language Resources and Evaluation (LREC2006)*, pages 19–23, 2006. 16

[9] Xi Chen, Yan Duan, Rein Houthooft, John Schulman, Ilya Sutskever, and Pieter Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. *CoRR*, abs/1606.03657, 2016. 3, 4, 5, 12

[10] Jianpeng Cheng, Li Dong, and Mirella Lapata. Long short-term memory-networks for machine reading. *CoRR*, abs/1601.06733, 2016. 7

[11] Yunjey Choi, Min-Je Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. *CoRR*, abs/1711.09020, 2017. i, 2, 4, 5, 6, 7, 25

[12] Antonia Creswell, Tom White, Vincent Dumoulin, Kai Arulkumaran, Biswa Sengupta, and Anil A. Bharath. Generative adversarial networks: An overview. *IEEE Signal Process. Mag.*, 35(1):53–65, 2018. 1, 3, 4, 8, 9, 12, 16, 17, 20

[13] Anastasios D Doulamis, Yannis S Avrithis, Nikolaos D Doulamis, and Stefanos D Kollias. Interactive content-based retrieval in video databases using fuzzy classification and relevance feedback. In *Proceedings IEEE International Conference on Multimedia Computing and Systems*, volume 2, pages 954–958. IEEE, 1999. 16

[14] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. http://www.deeplearningbook.org. 8

[15] Ian J. Goodfellow, Dumitru Erhan, Pierre Luc Carrier, Aaron C. Courville, Mehdi Mirza, Benjamin Hamner, William Cukierski, Yichuan Tang, David Thaler, Dong-Hyun Lee, Yingbo Zhou, Chetan Ramaiah, Fangxiang Feng, Ruifan Li, Xiaojie Wang, Dimitris Athanasakis, John Shawe-Taylor, Maxim Milakov, John Park, Radu Tudor Ionescu, Marius Popescu, Cristian Grozea, James Bergstra, Jingjing Xie, Lukasz Romaszko, Bing Xu, Chuang Zhang, and Yoshua Bengio. Challenges in representation learning: A report on three machine learning contests. In *Neural Information Processing - 20th International Conference, ICONIP 2013, Daegu, Korea, November 3-7, 2013. Proceedings, Part III*, pages 117–124, 2013. i, 2, 23, 25

[16] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pages 2672–2680, 2014. i, 1, 8, 9, 17, 20

[17] Ian J. Goodfellow, David Warde-Farley, Mehdi Mirza, Aaron C. Courville, and Yoshua Bengio. Maxout networks. In *Proceedings of the 30th International Conference on Machine Learning, ICML 2013, Atlanta, GA, USA, 16-21 June 2013*, pages 1319–1327, 2013. 8

[18] Georgios Goudelis, Konstantinos Karpouzis, and Stefanos Kollias. Exploring trace transform for robust human action recognition. *Pattern Recognition*, 46(12):3238–3248, 2013. 16

[19] Ishaan Gulrajani, Faruk Ahmed, Martín Arjovsky, Vincent Dumoulin, and Aaron C. Courville. Improved training of wasserstein gans. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, pages 5769–5779, 2017. 9

[20] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997. 12

[21] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *CoRR*, abs/1502.03167, 2015. 4

[22] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A. Efros. Image-to-image translation with conditional adversarial networks. *CoRR*, abs/1611.07004, 2016. 4, 5, 9, 12, 17

[23] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *CoRR*, abs/1710.10196, 2017. 5, 10

[24] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. *CoRR*, abs/1812.04948, 2018. 5, 13, 14

[25] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. 8

[26] Ilianna Kollia, Nikolaos Simou, Andreas Stafylopatis, and Stefanos Kollias. Semantic image analysis using a symbolic neural architecture. *Image Analysis & Stereology*, 29(3):159–172, 2010. 12

[27] Ilianna Kollia, Nikolaos Simou, Giorgos Stamou, and Andreas Stafylopatis. Interweaving knowledge representation and adaptive neural networks. In *Workshop on Inductive Reasoning and Machine Learning on the Semantic Web*, 2009. 4

[28] Dimitrios Kollias, Shiyang Cheng, Maja Pantic, and Stefanos Zafeiriou. Photorealistic facial synthesis in the dimensional affect space. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 0–0, 2018. 16

[29] Dimitrios Kollias, Shiyang Cheng, Evangelos Ververas, Irene Kotsia, and Stefanos Zafeiriou. Generating faces for affect analysis. *arXiv preprint arXiv:1811.05027*, 2018. 16

[30] Dimitris Kollias, George Marandianos, Amaryllis Raouzaiou, and Andreas-Georgios Stafylopatis. Interweaving deep learning and semantic techniques for emotion analysis in human-machine interaction. In *2015 10th International Workshop on Semantic and Social Media Adaptation and Personalization (SMAP)*, pages 1–6. IEEE, 2015. 12

[31] Dimitrios Kollias, Mihalis A Nicolaou, Irene Kotsia, Guoying Zhao, and Stefanos Zafeiriou. Recognition of affect in the wild using deep neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 26–33, 2017. 22

[32] Dimitrios Kollias, Athanasios Tagaris, and Andreas Stafylopatis. On line emotion detection using retrainable deep neural networks. In *2016 IEEE Symposium Series on Computational Intelligence (SSCI)*, pages 1–8. IEEE, 2016. 12

[33] Dimitrios Kollias, Athanasios Tagaris, Andreas Stafylopatis, Stefanos Kollias, and Georgios Tagaris. Deep neural architectures for prediction in healthcare. *Complex & Intelligent Systems*, 4(2):119–131, 2018. 16

[34] Dimitrios Kollias, Panagiotis Tzirakis, Mihalis A Nicolaou, Athanasios Papaioannou, Guoying Zhao, Björn Schuller, Irene Kotsia, and Stefanos Zafeiriou. Deep affect prediction in-the-wild: Aff-wild database and challenge, deep architectures, and beyond. *International Journal of Computer Vision*, 127(6-7):907–929, 2019. 22

[35] Dimitrios Kollias, Miao Yu, Athanasios Tagaris, Georgios Leontidis, Andreas Stafylopatis, and Stefanos Kollias. Adaptation and contextualization of deep neural network models. In *2017 IEEE Symposium Series on Computational Intelligence (SSCI)*, pages 1–8. IEEE, 2017. 1

[36] Dimitrios Kollias and Stefanos Zafeiriou. Aff-wild2: Extending the aff-wild database for affect recognition. *arXiv preprint arXiv:1811.07770*, 2018. 22

[37] Dimitrios Kollias and Stefanos Zafeiriou. A multi-component cnn-rnn approach for dimensional emotion recognition in-the-wild. *arXiv preprint arXiv:1805.01452*, 2018. 26

[38] Dimitrios Kollias and Stefanos Zafeiriou. A multi-task learning & generation framework: Valence-arousal, action units & primary expressions. *arXiv preprint arXiv:1811.07771*, 2018. 22

[39] Dimitrios Kollias and Stefanos Zafeiriou. Training deep neural networks with different datasets in-the-wild: The emotion recognition paradigm. In *2018 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2018. 12

[40] Dimitrios Kollias and Stefanos Zafeiriou. Exploiting multi-cnn features in cnn-rnn based dimensional emotion recognition on the omg in-the-wild dataset. *arXiv preprint arXiv:1910.01417*, 2019. 26

[41] Dimitrios Kollias and Stefanos Zafeiriou. Expression, affect, action unit recognition: Aff-wild2, multi-task learning and arcface. *arXiv preprint arXiv:1910.04855*, 2019. 22

[42] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates, Inc., 2012. 4

[43] Chuan Li and Michael Wand. Precomputed real-time texture synthesis with markovian generative adversarial networks. *CoRR*, abs/1604.04382, 2016. 5

[44] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2015. 3

[45] Andrew L Maas, Awni Y Hannun, and Andrew Y Ng. Rectifier nonlinearities improve neural network acoustic models. In *Proc. ICML*, volume 30, 2013. 8

[46] Elman Mansimov, Emilio Parisotto, Jimmy Lei Ba, and Ruslan Salakhutdinov. Generating images from captions with attention. *arXiv preprint arXiv:1511.02793*, 2015. 12, 16

[47] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. In *1st International Conference on Learning Representations, ICLR 2013, Scottsdale, Arizona, USA, May 2-4, 2013, Workshop Track Proceedings*, 2013. 12

[48] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *CoRR*, abs/1411.1784, 2014. 3, 4, 16

[49] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. *CoRR*, abs/1802.05957, 2018. 10

[50] Takeru Miyato and Masanori Koyama. cgans with projection discriminator. *CoRR*, abs/1802.05637, 2018. 3

[51] Michael A. Nielsen. *Neural Networks and Deep Learning*. Determination Press, 2015. 8

[52] Ankur P. Parikh, Oscar Täckström, Dipanjan Das, and Jakob Uszkoreit. A decomposable attention model for natural language inference. *CoRR*, abs/1606.01933, 2016. 7

[53] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. In *NIPS-W*, 2017. 25

[54] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *CoRR*, abs/1511.06434, 2015. 4, 8, 9, 12, 16

[55] Konstantinos A Raftopoulos, Stefanos D Kollias, Dionysios D Sourlas, and Marin Ferecatu. On the beneficial effect of noise in vertex localization. *International Journal of Computer Vision*, 126(1):111–139, 2018. 26

[56] Scott E. Reed, Zeynep Akata, Santosh Mohan, Samuel Tenka, Bernt Schiele, and Honglak Lee. Learning what and where to draw. *CoRR*, abs/1610.02454, 2016. 4, 5, 17

[57] Scott E. Reed, Zeynep Akata, Xinchen Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. Generative adversarial text to image synthesis. *CoRR*, abs/1605.05396, 2016. 4, 16

[58] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. *CoRR*, abs/1505.04597, 2015. 4

[59] Tim Salimans, Ian J. Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. *CoRR*, abs/1606.03498, 2016. 3, 8, 9, 10

[60] Nikos Simou, Th Athanasiadis, Giorgos Stoilos, and Stefanos Kollias. Image indexing and retrieval using expressive fuzzy description logics. *Signal, Image and Video Processing*, 2(4):321–335, 2008. 16

[61] Nikolaos Simou and Stefanos Kollias. Fire: A fuzzy reasoning engine for impecise knowledge. In *K-Space PhD Students Workshop, Berlin, Germany*, volume 14. Citeseer, 2007. 20

[62] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin A. Riedmiller. Striving for simplicity: The all convolutional net. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Workshop Track Proceedings*, 2015. 4

[63] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15:1929–1958, 2014. 8

[64] Athanasios Tagaris, Dimitrios Kollias, and Andreas Stafylopatis. Assessment of parkinsons disease based on deep neural networks. In *International Conference on Engineering Applications of Neural Networks*, pages 391–403. Springer, 2017. 16

[65] Athanasios Tagaris, Dimitrios Kollias, Andreas Stafylopatis, Georgios Tagaris, and Stefanos Kollias. Machine learning for neurodegenerative disorder diagnosissurvey of practices and launch of benchmark dataset. *International Journal on Artificial Intelligence Tools*, 27(03):1850011, 2018. 16

[66] T. Tieleman and G. E. Hinton. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural Networks for Machine Learning*, 4(2), 2012. 9

[67] Jiajun Wu, Chengkai Zhang, Tianfan Xue, Bill Freeman, and Josh Tenenbaum. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 82–90, 2016. 4, 12, 13

[68] Stefanos Zafeiriou, Dimitrios Kollias, Mihalis A Nicolaou, Athanasios Papaioannou, Guoying Zhao, and Irene Kotsia. Aff-wild: Valence and arousal'in-the-wild'challenge. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 34–41, 2017. 22

[69] Han Zhang, Ian J. Goodfellow, Dimitris N. Metaxas, and Augustus Odena. Self-attention generative adversarial networks. *CoRR*, abs/1805.08318, 2018. 1, 5, 7, 10

[70] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaolei Huang, Xiaogang Wang, and Dimitris N. Metaxas. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. *CoRR*, abs/1612.03242, 2016. 4, 5, 6, 13, 17

[71] Kaipeng Zhang, Zhanpeng Zhang, Zhifeng Li, and Yu Qiao. Joint face detection and alignment using multi-task cascaded convolutional networks. *CoRR*, abs/1604.02878, 2016. 23

[72] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. *CoRR*, abs/1703.10593, 2017. 5, 6, 17, 18, 19

[73] Jun-Yan Zhu, Richard Zhang, Deepak Pathak, Trevor Darrell, Alexei A. Efros, Oliver Wang, and Eli Shechtman. Toward multimodal image-to-image translation. *CoRR*, abs/1711.11586, 2017. 5, 6, 13, 14, 15