

---

# GANsformer Reproducibility Challenge

---

**Giorgia Adorni**  
giorgia.adorni@usi.ch

**Felix Boelter**  
felix.boelter@usi.ch

**Stefano Carlo Lambertenghi**  
stefano.carlo.lambertenghi@usi.ch

## Abstract

The field of image generation through generative modelling is abundantly discussed nowadays. It can be used for a highly varied range of applications such as up-scaling already existing images, creating non-existing objects, such as interior design scenes, products or even human faces, and achieving transfer-learning processes. In this context, Generative Adversarial Networks (GANs) are a class of widely studied machine learning frameworks first appearing in the paper “*Generative adversarial nets*” by Goodfellow et al. [1] that achieve the aforementioned goal. In our work, we reproduce and evaluate a novel variation of the original GAN network, the GANformer, proposed in “*Generative adversarial Transformers*” by Hudson and Zitnick [2]. The goal of this project was to recreate the methods presented by this paper to reproduce the original results and comment on the authors’ claims. Due to resources and time limitations, we had to constraint the networks training times and dataset types and sizes. Our research successfully recreated both variations of the proposed GANformer model and found extreme differences between the authors’ and our results. Moreover, discrepancies between the publication methodology and the one implemented, made available in the code, allowed us to study two undisclosed variations of the presented procedures.

## 1 Introduction

This project aims at investigating the reliability and reproducibility of a paper accepted for publication in a top machine learning conference. The models have been implemented using code and information provided by the authors.

With this work, we are going to verify the empirical results and claims of the paper “*Generative adversarial Transformers*” by Hudson and Zitnick [2], by reproducing three of the computational experiments performed by the authors: (1) the *StyleGAN2* by Karras et al. [3, 4], a GAN network that uses one global latent style vector to modulate the features of each layer, hence controlling the style of all image features globally. (2) the *GANformer* with *Simplex Attention* by Hudson and Zitnick [2], which generalises the StyleGAN design with  $k$  latent vectors that cooperate through attention. Thus, allowing for a spatially finer control over the generation process since multiple style vectors impact different regions in the image concurrently, particularly permitting communication in one direction, in the generative context – from the latent vectors to the image features. (3) the *GANformer* with *Duplex Attention* by Hudson and Zitnick [2], which is based on the same principles as the previous but propagating information from global latent vectors to local image features, enabling both top-down and bottom-up reasoning to coincide.

The first model is used as a baseline, while the remaining are the architectures introduced by the authors. They consider the GANformer as “a novel and efficient type of transformer” which demonstrates its strength and robustness over a range of tasks of visual generative modelling — simulated multi-object environments (real-world indoor and outdoor scenes) — achieving state-of-the-art results in terms of both image quality and diversity while benefiting from fast learning and better data-efficiency.

## 2 Methodology: Generative Adversarial Transformers

The Generative Adversarial Transformers (GANsformer), introduced by Hudson and Zitnick [2], are models which combine GANs and the transformers to generate better and more realistic examples.

GANs, and in particular the StyleGAN2 model [3], presented in Section ??, is used as a starting point for the GANformer design for the properties it owns as CNN: they are powerful generators for the overall style of the image, since by nature they merge the local information of the pixels together with the general information regarding the image. However, they are less powerful with respect to small details of localised regions within the generated image itself, since they miss out on the long range interaction of the faraway pixel.

Accordingly, GANsformers take advantage of the transformers attention mechanism to make the StyleGAN2 architecture even more powerful: the integration of attention in the architecture allows the network to draw global dependencies between input and output, and understand the context of the image thanks to the transformer's strength for long-range interactions. Thus, rather than focusing on using global information and controlling all features globally, the transformer uses attention to propagate information from the local pixels to the global high-level representation and vice versa.

The *bipartite transformer* structure computes *soft attention*, iteratively aggregating and disseminating information between the generated image features and a compact set of *latent variables* enabling bidirectional interaction between these dual representations. This architecture offers a solution to the StyleGAN limitation in its ability to perform spatial decomposition which leads to the impossibility of controlling the style of a localised region within the generated image.

The *transformer network* corresponds to the *multi-layer bidirectional transformer encoder* (BERT), introduced by Devlin et al. [6], which interleaves *multi-head self-attention* and *feed-forward layers*.

The discriminator model performs multiple layers of convolution down-sampling on the image, reducing the representation's resolution gradually until making final prediction. Optionally, attention can be incorporated into the discriminator as well where it has multiple  $k$  aggregator variables, that use attention to adaptively collect information from the image while being processed.

The generator likewise is composed of two parts, a mapping network and a synthesis network. The mapping network of a GANformer is the same as that of the StyleGAN2. In the synthesis network, while in the StyleGAN2, a single global  $w$  vector controls all the features equally, the GANformer uses attention so that the  $k$  latent components specialise to control different regions in the image to create it cooperatively, and therefore perform better especially in generating images depicting multi-object scenes, also allowing for a flexible and dynamic style modulation at the region level.

Hudson and Zitnick [2] have applied some adaptations to the structure of the GANformer as presented here, in order to foster an interesting communication flow. The details are provided later in Section 3.2. Rather than densely modelling interactions among all the pairs of pixels in the images, instead it supports *adaptive long-range interaction* between far away pixels in a moderated manner, passing through a compact and global latent bottleneck that selectively gathers information from the entire input and distributes it back to the relevant regions.

There are two attention operations that could be computed over the bipartite graph, depending on the direction in which information propagates, (1) *simplex attention* permits communication either in one way only, in the generative context, from the latent vectors to the image features, and (2) *duplex attention*, which enables it both top-down and bottom-up.

## 2.1 Simplex attention

The simplex and duplex attention layers and formulas are provided in the original article

I'm not sure about keeping all the details regarding simplex and duplex attention since we can find these stuff exactly the same in the original paper.

As already mentioned, simplex attention distributes information in a single direction over the bipartite transformer graph.

Formally, let  $X^{n \times d}$  denote an input set of  $n$  vectors of dimension  $d$  — where, for the image case,  $n = W \times H$  — and  $Y^{m \times d}$  denote a set of  $m$  aggregator variables — the latent variables, in the generative case. Specifically, the attention is computed over the derived bipartite graph between these two groups of elements, as in Equation (??), moreover:

$$a(X, Y) = \text{Attention}(q(X), k(Y), v(Y)), \quad (1)$$

where  $q(\cdot)$ ,  $k(\cdot)$ ,  $v(\cdot)$  are functions that respectively map elements into queries, keys, and values, all maintaining dimensionality  $d$ . The mappings are provided with positional encodings to reflect the distinct position of each element (e.g. in the image). This bipartite attention is a generalisation of self-attention, where  $Y = X$ .

Standard transformers implement an additive update rule of the form:

$$u^a(X, Y) = \text{LayerNorm}(X + a(X, Y)), \quad (2)$$

however, [2] used the retrieved information to control both the scale as well as the bias of the elements in  $X$ , in line with the practice promoted by the StyleGAN model [4]:

$$u^s(X, Y) = \gamma(a(X, Y)) \odot \omega(X) + \beta(a(X, Y)), \quad (3)$$

where  $\gamma(\cdot)$ ,  $\beta(\cdot)$  are mappings that compute multiplicative and additive styles (gain and bias), maintaining dimensionality  $d$ , and  $\omega(X) = X - \mu(X)$  normalises each element with  $\sigma(X)$  respect to the other features. By normalising  $X$  (image

features), and then letting  $Y$  (latent vectors) control the statistical tendencies of  $X$ , the information propagation from  $Y$  to  $X$  is enabled, allowing the latent vectors to control the visual generation of spatial attended regions within the image, so as to guide the synthesis of objects or entities. The multiplicative integration permits significant gains in the model performance.

## 2.2 Duplex attention

Duplex attention can be explained by taking into account the variables  $Y$  to set their own key-value structure:  $Y = (K^{n \times d}, V^{n \times d})$ , where the values store the content of the  $Y$  variables, as before (e.g. the randomly sampled latent vectors in the case of GANs) while the keys track the centroids  $K$  of the attention-based assignments between  $Y$  and  $X$ , which can be computed as  $K = a(Y, X)$  — namely, the weighted averages of the  $X$  elements using the bipartite attention distribution derived through comparing it to  $Y$ . Consequently, the new update rule is defined as follows:

$$u^d(X, Y) = \gamma(A(X, K, V)) \odot \omega(X) + \beta(A(X, K, V)), \quad (4)$$

where, two attention operations are compound on top of each other: first compute the *soft attention* assignments between  $X$  and  $Y$ , by  $K = a(Y, X)$ , and then refine the assignments by considering their centroids, by  $A(X, K, V)$ . This is analogous to the *k-means algorithm* and works more effectively than the simpler update  $u^a$  defined above in Equation (3).

Finally, to support bidirectional interaction between  $X$  and  $Y$  (the image and the latent vectors), two reciprocal simplex attentions are chained from  $X$  to  $Y$  and from  $Y$  to  $X$ , obtaining the duplex attention, which alternates computing  $Y := u^a(Y, X)$  and  $X := u^d(X, Y)$ , such that each representation is refined in light of its interaction with the other, integrating together bottom-up and top-down interactions.

## 3 Implementation

All this section should be revised since we are using also another dataset

The code from the authors' has been merged with the code provided by StyleGAN2 to obtain a hybrid version of the StyleGAN2 and the GANformer. In addition, we created a simplified version of the code which removed unnecessary operations in the creation of the network that were used for other model architectures. Furthermore, we implemented a Google Colab Pro version of the authors' code, as we had access to a more powerful GPU.

### 3.1 Datasets

The original paper [2] explored the GANformer model on four datasets for images and scenes: CLEVR [7], LSUN-Bedrooms [8], Cityscapes [9] and FFHQ [4].

Initially, we used the Cityscapes dataset since it is smaller among the four: it contains 25k images with 256x256 resolution. However, the memory required to complete the training was too high on this dataset (more than 25Gb). Even if we had more memory available, Colab Pro's limitation of 24 hours sessions would have interrupted our experiments prematurely.

For this reason, we switch to another dataset, the Google Cartoon Set [10]<sup>1</sup>, containing 10k 2D cartoon avatar images with 64x64 resolution, composed of 16 components that vary in 10 artwork attributes, 4 colour attributes, and 4 proportion attributes (see Table 6 in Appendix A).

After an initial examination of the result obtained with this dataset, we decided to proceed further using a more challenging dataset, the FFHQ, exploited in fact from the authors of the reproduced paper. This dataset, presented by Karras et al. [4], is a collection 70k high-quality images of human faces at a  $1024 \times 1024$  resolution, meaning that it offers a much higher quality and a vastly more variation in terms of age, ethnicity, image background and coverage of accessories such as eyeglasses, sunglasses, hats, etc., than existing high-resolution datasets.

### 3.2 Hyper-parameters and design choices

In this section we present the relevant hyper-parameters used in our experimentation, both for training and also in terms of layer sizes and technical choices.

<sup>1</sup><https://google.github.io/cartoonset>

Table 1 contains a comparison between StyleGAN2 (the baseline) and the novel networks proposed in the original paper.

Note that in the code provided by the author [2], the hyper-parameters are not the same as mentioned in the article.

**Table 1: Comparison of the hyper-parameters given in the code with those mentioned in the paper statements.** GANformer<sub>s</sub> refers to the GANformer with Simplex attention, while GANformer<sub>d</sub> refers to the GANformer with duplex attention.

|                | StyleGAN2 | GANformer <sub>s</sub><br>(code) | GANformer <sub>d</sub><br>(code) | GANformer <sub>s</sub><br>(article) | GANformer <sub>d</sub><br>(article) |
|----------------|-----------|----------------------------------|----------------------------------|-------------------------------------|-------------------------------------|
| latent_size    | –         | 32                               | 32                               | 32                                  | 32                                  |
| dlatent_size   | –         | 32                               | 32                               | 32                                  | 32                                  |
| components_num | –         | 16                               | 16                               | 16                                  | 16                                  |
| beta1          | 0.0       | 0.0                              | 0.0                              | 0.9                                 | 0.9                                 |
| beta2          | 0.99      | 0.99                             | 0.99                             | 0.999                               | 0.999                               |
| epsilon        | 1e-8      | 1e-8                             | 1e-8                             | 1e-3                                | 1e-3                                |

A *kernel size* of  $k = 3$  is used after each application of the attention, together with a *Leaky ReLU non-linearity* after each convolution and then up-sample or down-sample the features  $X$ , as part of the generator or discriminator respectively, as in e.g. StyleGAN2 [3]. To account for the features location within the image, we use a sinusoidal positional encoding along the horizontal and vertical dimensions for the visual features  $X$ , and trained positional embeddings for the set of latent variables  $Y$ . Overall, the bipartite transformer is thus composed of a stack that alternates attention (simplex or duplex), convolution, and up-sampling layers, starting from a  $4 \times 4$  grid up to the desirable resolution.

Both the simplex and the duplex attention operations enjoy a bi-linear efficiency of  $\mathcal{O}(mn)$  thanks to the network’s bipartite structure that considers all pairs of corresponding elements from  $X$  and  $Y$ . Since, as we see below, we maintain  $Y$  to be of a fairly small size, choosing  $m$  in the range of 8–32, this compares favourably to the prohibitive  $\mathcal{O}(n^2)$  complexity of self-attention, which impedes its applicability to high-resolution images.

As to the loss function, optimisation and training configurations, we adopt the settings and techniques used in StyleGAN2 [3], including in particular style mixing, stochastic variation, exponential moving average for weights, and a non-saturating logistic loss with a lazy R1 regularisation.

### 3.3 Experimental setup

The source code of our work is available at the following GitHub repository: <https://github.com/GiorgiaAuroraAdorni/gansformer-reproducibility-challenge>.

The approaches proposed in both the original paper codebase by Karras et al. [3] and by Hudson and Zitnick [2] have been implemented in Python using TensorFlow [11], so, according to that, we used the same setup. We created a Jupyter Notebook which runs all the experiments in Google Colaboratory, which allows us to write and execute Python in the browser.

All the models have been trained on a Tesla P100-PCIE-16GB (GPU) provided by Google Colab Pro.

### 3.4 Computational requirements

In the original paper [2], they evaluate all models under comparable conditions of training scheme, model size, and optimisation details, implementing all the models within the codebase introduced by the StyleGAN authors [3]. All models have been trained with images of  $256 \times 256$  resolution and for the same number of training steps, roughly spanning a week on 2 NVIDIA V100 GPUs per model (or equivalently 3–4 days using 4 GPUs).

Considered that we had available just one GPU and not enough time to reproduce this settings, we decided to resize the images from  $256 \times 256$  to  $64 \times 64$  resolution for the Google Carton Set and to  $128 \times 128$  for FFHQ dataset.

For the GANformer we select  $k = 32$  number of latent variables.

All models have been trained for the same number of steps, 300 000 image training samples (300 kimg) while the paper present results after training 100, 200, 500, 1000, 2000, 5000 and 10000 kimg samples

For the StyleGAN2 model we present results after training 300 kimg, obtaining good results. Note that the original StyleGAN2 model has been trained by its authors [3] for up to 70000 kimg samples, which is expected to take over 90 GPU-days for a single model.

For the GANformer, the authors [3] show impressive results, especially when using duplex attention: the model manages to learn a lot faster than competing approaches, generating astonishing images early in the training. This model is expected to take 4 GPU-days.

However, we are not able to replicate this achievements, first because this model learns significantly slower than the StyleGAN2, which is able to train approximately 1.3 times faster than the GANformer in terms of time per kimg. (in the paper they reach better results with the GANformer with 3-times less training steps than the StyleGAN2, but they don't specify the time required for a step) Secondly, the GANformer with simplex attention seems to be as slow if not slower to achieve qualitative results in terms of training steps when compared to StyleGAN2 and if Duplex attention is selected, qualitative results are never obtained.

As previously mentioned we trained using Colab Pro which enabled us to access a Tesla P100 GPU by Nvidia with 16Gb of vram and 25Gb of RAM.

For StyleGAN2, with the given resources, for 300 kimg, training took around 8h while for all variations of the GANformer, training took around 10 h.

## 4 Results

The all section should be revised since we are using also another dataset

This section shows and comments on our results while also comparing them to the original GANformer paper. We start by evaluating the two presented variations of the GANformer and StyleGAN2 over a set of four metrics: Frechet Inception Distance (FID), Inception Score (IS), Precision and Recall. The FID is one of the most popular metrics for evaluating GANs, providing stable and reliable image fidelity and diversity indications. It is a measure of similarity between curves that considers the location and ordering of the points along the curves. FID is used to measure the feature distance between the real and the generated images for this specific application. However, it can also be used to measure the distance between two distributions. For this reason, we have decided to use it as a reference metric for all the following analyses.

## 5 Google Cartoon Set results

Starting from the Google Cartoon Set, in Table 2, we compared the GANformer (Simplex and Duplex) with the competing StyleGAN2 model. To have a fair comparison, the three image synthesis methods are run for the same amount of iterations and the same dataset. To express the improvement over StyleGAN2, a difference factor of FID is positioned alongside the scores.

**Table 2: Comparison between the GANformer (Simplex and Duplex) and competing StyleGAN2.** In the last column, is reported the percentage of improvement all the models, in terms of FID score, with respect to the baseline StyleGAN2 architecture.

| Model                        | FID ↓        | IS ↑        | Precision↑    | Recall↑       | FID Improvement (%) |
|------------------------------|--------------|-------------|---------------|---------------|---------------------|
| StyleGAN2                    | <b>24.77</b> | 2.50        | <b>0.0018</b> | <b>0.0211</b> | 0 %                 |
| GANformer, Simplex attention | 28.11        | <b>2.58</b> | 0.0015        | 0.0076        | -13.48 %            |
| GANformer, Duplex attention  | 27.08        | 2.47        | <b>0.0018</b> | 0.0090        | -9.33 %             |

Unexpectedly, the novel GANformer with Duplex attention, is worse than the baseline on all aspects, except the precision metric, with a staggering -9.33 % deterioration in FID score. To investigate this further, a similar representation is recreated in Table 3 using the original paper findings for the same metrics and with a mean of the scores spanning the four used datasets.

It is noteworthy to state that, in the code given by the authors, attention seems to be only optionally used in the discriminator and when analysing the pre-trained models provided, it is never used, prompting us to implement two

Table 3: **Original paper’s reported results (mean of results over the 4 datasets used by the authors)**. In the last column, is reported the percentage of improvement all the models, in terms of FID score, with respect to the baseline StyleGAN2 architecture.

| Model                        | FID ↓       | IS ↑        | Precision↑   | Recall ↑     | FID Improvement (%) |
|------------------------------|-------------|-------------|--------------|--------------|---------------------|
| StyleGAN2                    | 11.29       | 2.74        | 52.02        | 23.98        | 0 %                 |
| GANformer, Simplex attention | 10.29       | <b>2.82</b> | <b>56.76</b> | 18.21        | +8.86 %             |
| GANformer, Duplex attention  | <b>7.22</b> | 2.78        | 55.45        | <b>33.94</b> | <b>+36.11 %</b>     |

variations of the GANformer not openly discussed in [2]. We use the GANformer paradigm for the generator and a vanilla StyleGAN2 discriminator and obtain the results visible in Table 4.

Table 4: **Comparison between the GANformer (Simplex and Duplex) both with and without attention on the Discriminator and competing StyleGAN2**. In the last column, is reported the percentage of improvement all the models, in terms of FID score, with respect to the baseline StyleGAN2 architecture.

| Model  | FID ↓        | IS ↑        | Precision↑    | Recall ↑      | FID Improvement (%) |
|--|--------------|-------------|---------------|---------------|---------------------|
| StyleGAN2                                      | 24.77        | 2.50        | 0.0018        | 0.0211        | 0 %                 |
| GANformer, Simplex attention                   | 28.11        | 2.58        | 0.0015        | 0.0076        | -13.48 %            |
| GANformer, Duplex attention                    | 27.08        | 2.47        | 0.0018        | 0.0090        | -9.33 %             |
| GANformer, Simplex attention (StyleGAN2 disc.) | <b>19.09</b> | <b>2.62</b> | <b>0.0035</b> | <b>0.0476</b> | <b>+22.93 %</b>     |
| GANformer, Duplex attention (StyleGAN2 disc.)  | 24.81        | <b>2.62</b> | <b>0.0035</b> | 0.0211        | -0.16 %             |

Due to the relevance of the FID metric, we compare all the created models score over iterations number to both show quality of results over training steps in Figure 1.

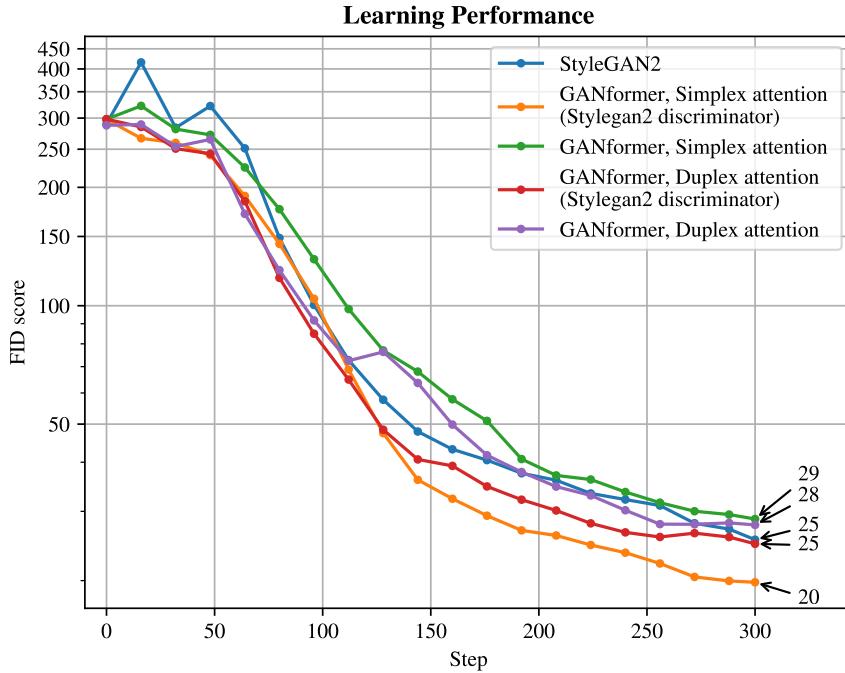


Figure 1: **Comparison between the StyleGAN2 and the GANformer models**. We evaluate the models according to FID score along 300k image samples. The score is computed every sixteen checkpoints.

As mentioned before, unexpectedly, the model that yields the better results in the original paper, is not the best in our experiments. Not only the FID score is worse than the baseline StyleGAN2 at the final training step, both for simplex

and duplex attention implementations, but also claims about efficiency cannot be verified. Models which include attention on the generator only, however, are faster in terms of steps to reach a qualitative results when compared to the baseline. We believe that this behaviour explains the choices made by the authors in their GitHub publication of the code. Moreover, a comment has to be made on the claim of efficiency: both with and without attention on the discriminator, a training step is considerably slower to be completed on the same resources when compared to StyleGAN2. While the latter is capable of having a training speed of 10.9 images generated a second, all flavours of the GANformer, in the best case, only yield 8.3 images generated a second.

Finally, to visualise our findings less empirically, we have used random seeds to create latent inputs and shown the resulting images generated by the baseline StyleGAN2 and all four presented variations of the GANformer in Figures 5 and 6 in Appendix A.

An image generated by each variation is visible in Figure 2.

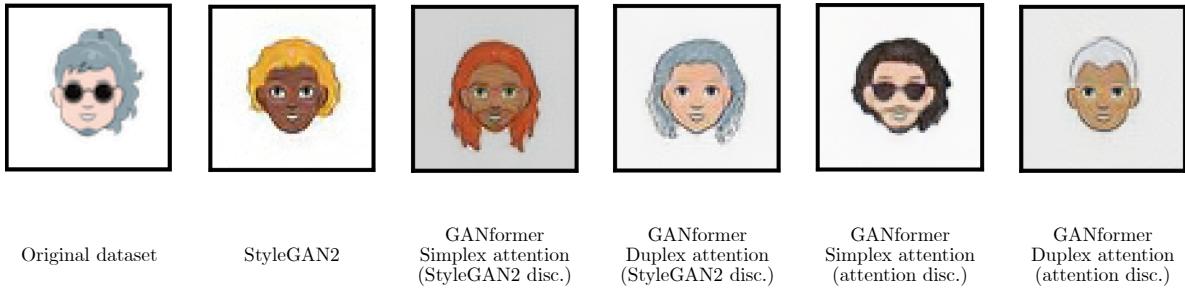


Figure 2: ... using of the various models.

It is blatantly obvious that the Duplex GANformer implementation is not capable of performing the task at hand, while all other models are at least able to produce visually compelling results. Comparing all outputs to the real images in the original dataset, all except Duplex GANformer tend to create similar results with only background colour differing for models with attention on the generator.

not true

## 6 FFHQ dataset results

For the FFHQ dataset, in Table 5 we compared the two GANformer models, both with Duplex attention on the generator, one with attention on the discriminator too and the other with a vanilla StyleGAN2 discriminator. To have a fair comparison, the three image synthesis methods are run for the same amount of iterations.

Table 5: Comparison between the GANformer (**Simplex** and **Duplex**) both with and without attention on the Discriminator and competing StyleGAN2. In the last column, is reported the percentage of improvement all the models, in terms of FID score, with respect to the baseline StyleGAN2 architecture.

| Model  | FID ↓        | IS ↑        | Precision↑  | Recall↑       | FID Improvement (%) |
|--|--------------|-------------|-------------|---------------|---------------------|
| StyleGAN2                                      | 40.22        | 3.21        |             |               | 0 %                 |
| GANformer, Simplex attention                   | 45.71        | 3.35        | 0.55        | 0.0055        | -13.65 %            |
| GANformer, Duplex attention                    | 53.48        | 3.35        | 0.48        | 0.0033        | -32.97 %            |
| GANformer, Simplex attention (StyleGAN2 disc.) | <b>39.73</b> | 3.49        |             |               | <b>+1.22 %</b>      |
| GANformer, Duplex attention (StyleGAN2 disc.)  | 43.66        | <b>3.61</b> | <b>0.55</b> | <b>0.0078</b> | -8.55 %             |

Once again, we can see that the model with multiple attention is considerably worse than the model with attention only in the generator.

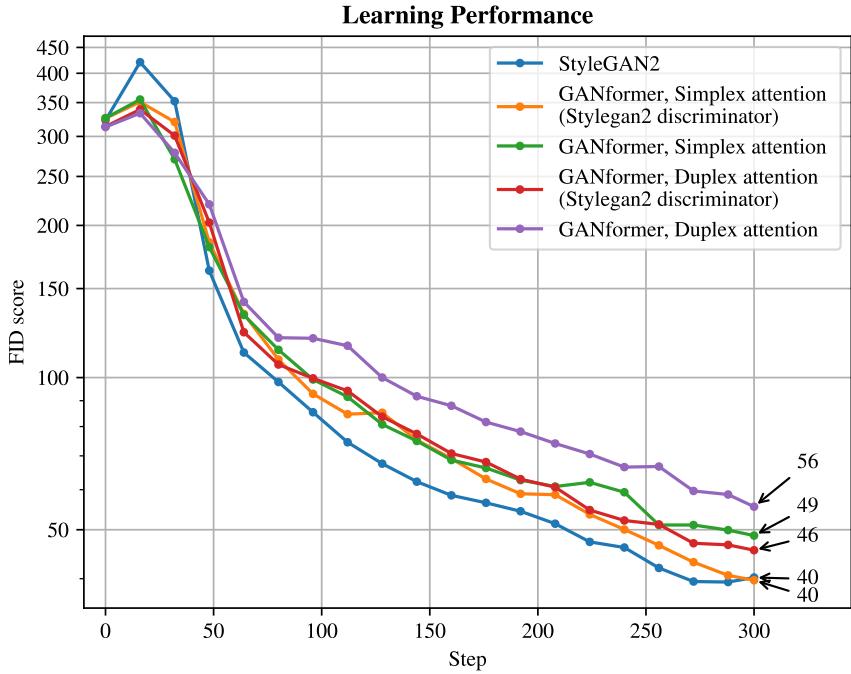


Figure 3: **Comparison between the StyleGAN2 and the GANformer models.** We evaluate the models according to FID score along 300k image samples. The score is computed every sixteen checkpoints.

Finally, to visualise our findings less empirically, as for the previous dataset, we have used random seeds to create latent inputs and shown the resulting images generated by the baseline StyleGAN2 and all four presented variations of the GANformer in Figures 7 and 8 in Appendix B.

An image generated by each variation is visible in Figure 2.



Figure 4: ... using of the various models.

## 7 Discussion and conclusion

[...]

All this section should be revised since we are using also another dataset

In this project we have discussed the replication of “*Generative adversarial Transformers*” by Hudson and Zitnick [2]. Contrary to our initial expectations, in our opinion, some of the claims made by the authors were misleading or unproven in a scientific manner. Once we made this realisation, our methodology shifted from a pure reproduction

of the results to a search of undisclosed variations of the novel GANformer system in the published material trying to obtain results that were comparable to the given ones. As mentioned in Section 3.4, our limited resources forced us to reduce the depth of our experiments compared to the authors'. We had to shift from four datasets to a single, size-reduced one: Forrester et al. [10], significantly decreasing the time of execution and the complexity of our research. For the same reasons, we have discarded four out of the five baselines: GAN, k-GAN, SAGAN and VQGAN in favour of just StyleGAN2. The choice of StyleGAN2 as the baseline was prompted by a similarity of its execution to the novel GANformer, where parts of the author's code are a carbon copy of the StyleGAN2 implementation first proposed in Karras et al. [4].

In this study, we have firstly implemented a Google Colab Pro compatible version of the authors code, enabling us to train and test StyleGAN2 and the two flavours of GANformer (Duplex and Simplex attention) in less than 54 hours in total. As presented in Section 4, we were not able to achieve the improvements declared by the authors but instead found a decrease both in time performance and quality. Believing it was an error arisen by our adaptation we have analysed the pre-trained models provided and noticed an alarming discrepancy between the code and the methodology discussed in the paper. In "*Generative adversarial Transformers*" by Hudson and Zitnick [2], the attention component is said to be placed both on the generator and the discriminator sub-networks, while in all the pre-trained models provided, it seems to never be used on the discriminator.

Believing that all of the authors experiments were performed on a different network than the one presented in the publication, we have tried to recreate the structure that could have yielded the claimed qualitative results. To perform such a task we have decomposed the GANformer network in two sections by keeping the presented Generator but substituting the hypothetically valid discriminator network with a vanilla StyleGAN discriminator.

To our surprise, the StyleGAN/GANformer hybrid performed significantly better than the baseline. In the case of duplex attention on the discriminator, we believe that introducing k-means will find the centroid/mean of multiple regions of the cartoon face and merge them into one representative region of the cartoon face, which will make the discriminator unable to accurately tell the difference between a fake image and a real image. Furthermore, the generator tries to maximise the discriminator's output for the generated fake image [12, 13], which will create shapes that are close to what the discriminator sees from the real images confusing the discriminator from differentiating fake and real images.

In conclusion, we have successfully reproduced the "*Generative Adversarial Transformers*" by Hudson and Zitnick [2] and found unexpected results. We have then modified the proposed methodologies to obtain an image generation network that takes inspiration from the novel GANformer and adapts it to produce images that score significantly higher than the baseline over four quality metrics.

## 8 Authors contribution

review

use the standard CRediT authorship contribution statement

All the authors analysed the original article and the code provided by the creators. They took together all the important decisions, either regarding the organisation of the work and the problems that have arisen.

Stefano dealt with the dataset preparation and preprocessing phase. He also adapted the StyleGAN2 baseline for Colab and performed the training loop.

Felix dealt with the given code and take charge of merging the StyleGAN2 and GANformer models, adapting the code so that it reproduced carefully what was written in the paper.

Stefano and Giorgia collaborated on the adaptation of the style transformer for StyleGAN2 and later on the adaptation of the image visualisation for the GANformer, and in general with the generation of images and videos.

Stefano and Felix worked together on the style mixing adaptation for the GANformer, and since only they had available the GPU resources from Google Colab, they also handled with the models training.

Stefano also concentrated on the code refactoring phase, mainly cleaning the code related to the discriminator GANformer.

Giorgia focused on the main conceptual ideas behind the original work, and performed a theoretical research on the topic and related works. Later she drafted this manuscript, with inputs from the other authors which provided critical feedback on it.

All the authors discussed and interpret the results and finally contribute to the final version of this article.

## References

- [1] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2, NIPS'14*, page 2672–2680. MIT Press, 2014. URL <http://arxiv.org/abs/1406.2661>.
- [2] Drew A. Hudson and C. Lawrence Zitnick. Generative Adversarial Transformers, 2021. URL <http://arxiv.org/abs/2103.01209>.
- [3] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8110–8119, 2020. URL <http://arxiv.org/abs/1912.04958>.
- [4] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4401–4410, 2019. URL <http://arxiv.org/abs/1812.04948>.
- [5] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. URL <http://arxiv.org/abs/1706.03762>.
- [6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics, 2019. URL <http://arxiv.org/abs/1810.04805>.
- [7] Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2901–2910, 2017. URL <http://arxiv.org/abs/1612.06890>.
- [8] Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop, 2015. URL <http://arxiv.org/abs/1506.03365>.
- [9] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016. URL <http://arxiv.org/abs/1604.01685v2>.
- [10] Cole Forrester, Mosseni Inbar, Krishnan Dilip, Sarna Aaron, Maschinot Aaron, Freeman Bill, and Fuman Shiraz. Cartoon set. <https://google.github.io/cartoonset/>, 2018.
- [11] 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. URL <https://www.tensorflow.org/>. Software available from tensorflow.org.
- [12] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron Courville. Improved Training of Wasserstein GANs. *arXiv:1704.00028 [cs, stat]*, December 2017. URL <http://arxiv.org/abs/1704.00028>. arXiv: 1704.00028 version: 3.
- [13] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein GAN. *arXiv:1701.07875 [cs, stat]*, 2017. URL <http://arxiv.org/abs/1701.07875>.

## A Google Cartoon Set results

The dataset used in this work is the Google Cartoon Set [10] introduced in Section 3.1, containing 10k 2D cartoon avatar. These images are composed of 16 components that vary in 10 artwork attributes, 4 colour attributes, and 4 proportion attributes (see Table 6).

Table 6: Attributes of the Cartoon Set.

|                    | # Variants           | Description   |
|--------------------|----------------------|---|
| <b>Artwork</b>     | chin_length          | 3 Length of chin (below mouth region)                                   |
|                    | eye_angle            | 3 Tilt of the eye inwards or outwards                                   |
|                    | eye_lashes           | 2 Whether or not eyelashes are visible                                  |
|                    | eye_lid              | 2 Appearance of the eyelids   |
|                    | eyebrow_shape        | 14 Shape of eyebrows  |
|                    | eyebrow_weight       | 2 Line weight of eyebrows   |
|                    | face_shape           | 7 Overall shape of the face   |
|                    | facial_hair          | 15 Type of facial hair (type 14 is no hair)                             |
|                    | glasses              | 12 Type of glasses (type 11 is no glasses)                              |
| <b>Colors</b>      | hair                 | 111 Type of head hair   |
|                    | eye_color            | 5 Color of the eye irises   |
|                    | face_color           | 11 Color of the face skin   |
|                    | glasses_color        | 7 Color of the glasses, if present                                      |
| <b>Proportions</b> | hair_color           | 10 Color of the hair, facial hair, and eyebrows                         |
|                    | eye_eyebrow_distance | 3 Distance between the eye and eyebrows                                 |
|                    | eye_slant            | 3 Similar to eye_angle, but rotates the eye and does not change artwork |
|                    | eyebrow_thickness    | 4 Vertical scaling of the eyebrows                                      |
|                    | eyebrow_width        | 3 Horizontal scaling of the eyebrows                                    |

## B FFHQ dataset results



(a) StyleGAN2.



(b) GANformer with Simplex attention and vanilla StyleGAN2 discriminator.



(c) GANformer with Simplex attention.



(d) GANformer with Duplex attention and vanilla StyleGAN2 discriminator.



(e) GANformer with Duplex attention.

Figure 5: Visualisation of 9 images generated with the various models.



(a) StyleGAN2.



(b) GANformer with Simplex attention and vanilla StyleGAN2 discriminator.



(c) GANformer with Simplex attention.



(d) GANformer with Duplex attention and vanilla StyleGAN2 discriminator.



(e) GANformer with Duplex attention.

Figure 6: **Simple z interpolation using of the various models.**



(a) StyleGAN2.



(b) GANformer with Simplex attention and vanilla StyleGAN2 discriminator.



(c) GANformer with Simplex attention.



(d) GANformer with Duplex attention and vanilla StyleGAN2 discriminator.



(e) GANformer with Duplex attention.

Figure 7: Visualisation of 9 images generated with the various models.



(a) StyleGAN2.



(b) GANformer with Simplex attention and vanilla StyleGAN2 discriminator.



(c) GANformer with Simplex attention.



(d) GANformer with Duplex attention and vanilla StyleGAN2 discriminator.



(e) GANformer with Duplex attention.

Figure 8: **Simple z interpolation using of the various models.**