

Fashion Style Generation: Evolutionary Search with Gaussian Mixture Models in the Latent Space

Imke Grabe¹, Jichen Zhu¹, and Manex Aguirrezabal²

¹ IT University of Copenhagen, Copenhagen, Denmark
{imgr,jicz}@itu.dk

² University of Copenhagen, Copenhagen, Denmark
manex.aguirrezabal@hum.ku.dk

Abstract. This paper presents a novel approach for guiding a Generative Adversarial Network trained on the *FashionGen* dataset to generate designs corresponding to target fashion styles. Finding the latent vectors in the generator’s latent space that correspond to a style is approached as an evolutionary search problem. A Gaussian mixture model is applied to identify fashion styles based on the higher-layer representations of outfits in a clothing-specific attribute prediction model. Over generations, a genetic algorithm optimizes a population of designs to increase their probability of belonging to one of the Gaussian mixture components or styles. Showing that the developed system can generate images of maximum fitness visually resembling certain styles, our approach provides a promising direction to guide the search for style-coherent designs.

Keywords: Intelligent fashion · Generative adversarial networks · Genetic algorithm · Gaussian mixture model

1 Introduction

In many areas of music, arts, and design, *artificial intelligence* (AI) can procedurally generate new cultural artifacts [21,6,27,31,30]. With the recent developments of *Generative Adversarial Networks* (GANs), AI technology can become a powerful asset for human creators to design complex objects. For example, researchers have explored how to use GANs to design fashion artifacts in fashion design. The emerging area of *intelligent fashion* investigates the detection and recommendation of clothing items, the analysis of style trends, and the synthesis of clothing [4]. As part of fashion synthesis, GANs have been applied in the creation of new items [16], the simulation of try-on scenarios [28], or for personalized design [32]. Because fashion plays a fundamental part in human culture, it is essential to investigate how generative AI might contribute to its creation.

This paper focuses on the generation of fashion styles, defined as visual themes, such as the combination of clothing artifacts or attributes as part of an outfit (e.g., blue pants and a vertically striped shirt). By extracting the representation of images on higher-level layers of an attribute prediction model,

styles can be identified based on high-dimensional visual themes [22]. Existing work in intelligent fashion has analyzed styles [25,17,12], ultimately allowing for the prediction of trend behavior [1]. Research on fashion style generation with GANs focuses on controlling specific features [29,13], conditioning design with a text encoding [32], or transferring an exact outfit to other poses [32,28]. The control with regards to fashion styles, as defined above, has, however, not been considered in the generative process.

An open problem in generative design is how to guide the generation of GANs towards desirable outcomes. The GANs’ generator network learns to map random input variables to the complex output features resembling the training data during the training. The procedure creates an entangled latent space, making it impossible to inspect how changes in the latent code affect the semantic output features. This entanglement impedes the control of the generated designs towards a desired look. Differentiated control of the latent features of generative clothing models has been achieved by conditioning the generator with the encoding of a text description in the latent vector [32], or by disentangling color, texture, and shape inputs through separate losses in the loss function during training [29]. Approaches like *StyleGAN* [15] aim at controlling certain stylistic features in images, such as the transfer of a complete outfit [28]. However, guiding the generation of designs towards fashion styles consisting of broader visual themes remains an open research problem. Fashion describes the specific category of clothing driven by the developments of style trends. As fashion styles capture meaningful temporal and local developments with societies [19,1,22], responding to such themes matters for the generative process.

This paper presents a new method to guide GANs using a *Gaussian mixture model* (GMM). While GMMs have been used in the field of intelligent fashion to identify fashion styles [22,1], they have not been applied to support generative purposes. To better control the GANs’ entangled latent space, we utilize the GMM to find the latent vectors that correspond to *stylistic* designs in an evolutionary search problem. More specifically, our method combines (1) generative deep learning, (2) the analysis of fashion styles, and (3) the application of genetic optimization algorithms to search a design space.

The main contribution of this work is a new method to guide GANs using a GMM. This paper presents the method proposed in our prior work [10]. It allows GANs’ generative process to be guided based on higher-level themes, instead of separate attributes as in existing approaches [29,13]. We tested our method in the context of generating fashion styles. We found that while our proposed framework generally supports the generation of designs according to target styles, the GMM-based fitness measure is not always aligned with visual coherency to the styles. Some generated designs reveal that the fitness measure relies on a machine-specific understanding of style. Further investigation into the interplay between style model and the exploration of the latent space and the parameter setting of the genetic algorithm is required to align the results with a human understanding of style.

The remainder of the paper is organized as follows. After presenting related work, we introduce our dataset and proposed model, consisting of a GAN model, a style model, and the evolutionary search connecting the former. Next, we present our experimental results and conclude with discussions.

2 Related Work

This section presents related work in the three research areas relevant to this study, namely within (1) GANs, (2) analyzing fashion styles, and (3) applying evolutionary search to steer the generation process.

2.1 GANs

The introduction of GANs revolutionized the creation of computer-generated content [9]. GANs, consisting of two neural networks competing against each other as generator and a discriminator, learn to generate outputs by resembling a training distribution. For example, GANs can generate clothing artifacts when trained with a dataset of those. Notably, the training method of *Progressively growing GANs* (P-GANs) supports the generation of high-resolution images, as was demonstrated by Rostamzadeh et al. [23] with their introduction of a fashion dataset. In fashion generation, different objectives have been guiding the training of GANs, such as conditioning their output with the text description of desired looks [32], or color, texture, and shape [29].

Notice that for GANs, the term *style* is typically used to describe the manipulation of a design with regards to specified visual attributes. StyleGANs [15] provide an architecture that disentangles the latent space of the generator network, allowing for a targeted modification of high-level to low-level attributes corresponding to different resolutions in the network. Yildirim et al. [28] trained a StyleGAN to transfer a complete outfit to other models, as in a try-on scenario. This notion sets focus on certain visual attributes, similar to Jiang et al. [13], who apply GANs to transfer patterns to shirt designs.

2.2 Fashion styles

Studies have addressed the phenomenon of styles in fashion from different angles, varying from weak [18,8] to strong [25,17] style annotations, as well as their unsupervised discovery [12,1,22]. Drawing on the latter, this paper builds on research that approaches fashion styles as a “mode in the data capturing a distribution of attributes” [1, p.7]. Hence, they can be discovered unsupervised by clustering attribute predictions of images. Clustering models such as a Gaussian mixture model (GMM) have previously been applied to find recurring components in the attribute embedding of images [22,1]. After identifying styles with a GMM, projecting the embedding of any (generated) image onto the GMM can measure of how well it fits into the discovered styles. In that way, we can evaluate how an image resembles a particular style.

Instead of using the prediction scores of an attribute prediction model as the embedding, previous layers of the model also capture meaningful themes in images. Matzen et al. [22] introduce a method that finds styles as Gaussian mixture components in the projections of images onto the model’s penultimate feature space. The second-last layer captures learned features that are more specific than the output of the final fully-connected layer. It serves as a valuable embedding for recognizing high-level themes beyond attribute scores [20]. The resulting clusters can be understood as a “global visual vocabulary for fashion” [22, p.7]. Letting this *language of fashion* inform the generative process by GANs is the objective of our suggested framework. We adapt Matzen et al.’s [22] procedure to find styles in the dataset *FashionGen* to use them for guiding the generation towards a chosen style.

To sum up, as the subject of fashion analysis, *style* refers to the broader sense of visual themes in outfits. Unless otherwise specified, the rest of the paper will use the term *style* to refer to this definition. Going beyond the attributes corresponding to the different layers of the GAN, specific items, or textures, the concept has not been considered in generation yet. We suggest a system affording the generation of designs based on fashion styles discovered in an unsupervised manner by a GMM. Embedded into an evolutionary search, the GMM guides the generation of designs.

2.3 Evolutionary search of GANs’ latent space

Our project aims to generate designs of various styles with visually diverse attributes with GANs. Navigating the networks’ complex latent space with this objective requires changing minor and major visual features. Drawing on the smooth characteristic of GANs’ latent space, evolutionary search can optimize its latent variables [3]. More specifically, a genetic algorithm alters a population of latent vectors by recombining and mutating them [5]. Through selection based on a fitness objective, such as maximizing the probability of belonging to a target style, latent vectors improve over many generations.

Previous work has applied genetic algorithms to explore the latent space with different goals. While Roziere et al. [24] improve the quality of only one fashion image generated by a P-GAN, Fernandes et al. [7] evolve a set of latent vectors to increase the diversity among the corresponding set of generated images. Instead of a pre-defined fitness objective, interactive genetic algorithms use user evaluation as a measure of fitness. Designing smaller clothing items in this interactive manner has been proposed in *DeepIE* [2], further developed into *StyleIE* [26]. Where an interactive genetic algorithm uses a human-in-the-loop to assess fitness, we apply a GMM to measure the fitness of a generated image, following the goal of generating designs that belong to an automatically discovered target style. Our contribution is to create a fully automated method for style generation by adding a GMM to the evolutionary loop. Instead of a human who chooses images of desired properties, the selection is based on the images’ projection onto the clustering space.

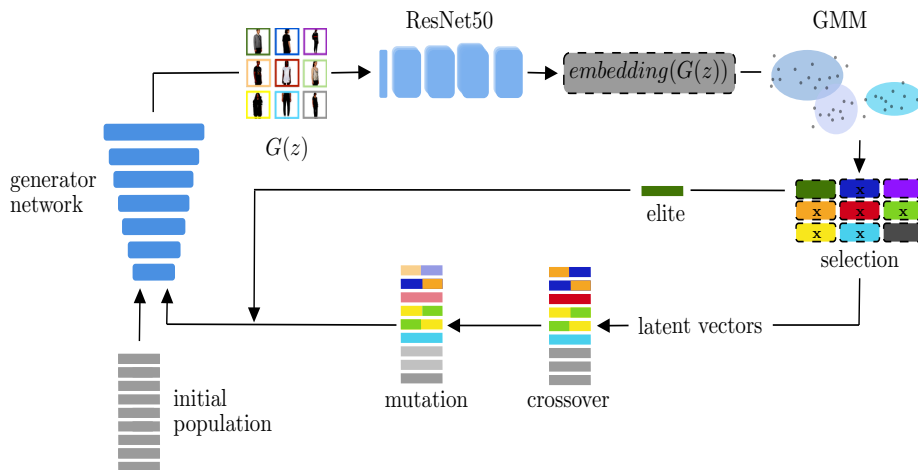


Fig. 1. Model of the genetic algorithm for searching the generative model’s latent space with the help of the clustering model. First, the trained generator network creates images based on randomly initialized latent vectors. Next, the generated images are represented as their embedding of the ResNet50 and then projected onto the clustering space of the GMM. The fittest ones are selected based on their posterior probability of belonging to the target cluster. Additionally, the elite, is preserved for the next generation. The others are recombined and new individuals are added before being mutated. The resulting latent vectors lay the basis for the next generation of images.

3 Dataset

FashionGen’s [23] clothing partition in a resolution of 256×256 pixels is used as a dataset for training the generative model and the style model. Consisting of over 200,000 images of outfits worn by a model, each outfit is represented in four different poses in the subset. The dataset mainly consists of images of whole-body outfits taken under consistent lighting conditions in front of white backgrounds, providing the ideal conditions for generation and clustering to focus purely on fashion attributes and identify styles across a combination of artifacts.

4 Model

Our proposed framework consists of three parts: A generative model, a style model, and the evolutionary search connecting the two. To make the (1) generative model output designs of a style that is discovered by the (2) style model, parts of the two models are combined in an (3) evolutionary search algorithm, as illustrated in Fig. 1. The details of the components are presented in this order.³

³ The code is available: <https://github.com/imkegrave/fashionstyle-generation-GMM>

Table 1. GAN training parameters. Note that the batch size was changed to 8 for the last training epoch.

Parameter	Setting
Optimizer	Adam
Activation function	leakyRELU
Learning rate	Equalized learning rate
Batch size	16 (8)
Loss function	WGAN-GP
Noise distribution	$\mathcal{N}(\mu = 0, \sigma = 1)$

4.1 Generative model

As the first part of the framework, a GAN is trained on the complete dataset. The P-GAN architecture and training procedure are utilized, which applies the *Wasserstein GAN* loss with gradient penalty (WGAN-GP) [14].⁴ Hence, the goal of the training procedure is to minimize the following function with respect to G , and maximize it with respect to D :

$$\min_G \max_D \mathbb{E}_{x \sim p_{data}} [D(x)] - \mathbb{E}_{z \sim p_z} [D(G(z))] + \lambda \mathbb{E}_{\hat{x} \sim p(\hat{x})} [(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2] \quad (1)$$

By competing against a discriminator network D , a generator network G learns how to map an input vector z of 512 variables,⁵ randomly sampled from a standard normal distribution $\mathcal{N}(\mu = 0, \sigma = 1)$, to output images resembling the real data distribution in the dataset. The model was trained on a *NVIDIA Tesla V100-SXM2-16GB* GPU for seven days to create images of resolution 256×256 . Table 1 provides an overview of the training parameters.

4.2 Style model

The goal of the second part of the framework is to discover styles in the dataset. Informed by the method in [22], the visual embedding learned by an attribute prediction model is leveraged to cluster outfits into styles.

We create image representations using a *ResNet50*.⁶ The network was pre-trained to map input images of *DeepFashion* to 1000 clothing-specific attributes instead of more general ones as in models pre-trained on ImageNet. The pre-training makes it a robust feature basis for analyzing datasets containing a wide range of clothing items.

⁴ We use the implementation made available by Facebook Research: https://github.com/facebookresearch/pytorch_GAN_zoo

⁵ The implementation expands the latent codes with 20 additional variables based on the outfits' item representation. As initial experimentation did not make their effect on the evolutionary search behaviour apparent, we treated them as the other latent variables. Their role should be further examined in future experiments.

⁶ The backbone of the attribute prediction model provided by the open-source toolbox for visual fashion analysis by the Multimedia Lab, Chinese University of Hong Kong is adapted: <https://github.com/open-mmlab/mmfashion>

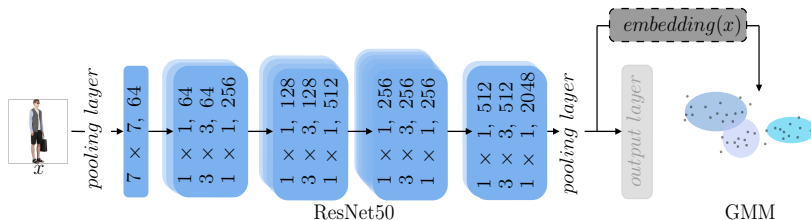


Fig. 2. Style model consisting of a ResNet50 as feature extraction model and a GMM as clustering model. The ResNet consists of several units containing 1-6 layer(s). A layer contains convolution(s) in the form of kernels, here notated as $r \times r, n$ with kernel size r and n channels. The output layer is discarded to extract the embedding for image x from the pooling layer. Based on the embedding for all outfits in the dataset, k Gaussian mixture components are determined.

As we are concerned with finding subtle styles, we approach them as a visual concept situated between pre-defined coarse categories and low-level single-feature attribution. To discover coherent style clusters representing these subtle visual themes, the penultimate layer of the ResNet50 serves as a meaningful embedding space. The last layer before the linear layer outputting the prediction scores for attributes captures high-level features of fashion images, though not directly class-specific to DeepFashion. In practice, the output of the last convolutional unit, more precisely after it has been converted to a 2048-dimensional vector by the pooling layer, is extracted as the *embedding(x)* of image x , as can be seen in Fig. 2.

Based on the embedding, clustering is employed to find recurring themes in the embedding space. To capture the distribution of all outfits in the style model while at the same time making the clustering computationally efficient, we chose only images in pose 4, which show most full-body images, chosen to build the style model. Hence, an embedding is retrieved for one image per outfit in the dataset. The embedding vectors are scaled to zero mean before principal component analysis (PCA) projects them onto the 135 principal components capturing 90% of the embedding’s variance.

A GMM is applied to find a mixture of Gaussian probability distributions that represent the data distribution. Through experimentation in line with Matzen et al. [22], we find that 150 mixture components seem to capture visually coherent fashion styles in the dataset. These style clusters range from 60 to 1100 ($\mu = 314.2$) images in size. Any image’s posterior probability p_t of belonging to a component t depicts how well it represents a style. The probabilistic model guarantees that increasing an image’s posterior probability of belonging to a cluster component reduces its probability of belonging to other clusters.

4.3 Evolutionary search

The trained generator of the GAN acts as a genotype-to-phenotype mapping, where the latent encoding, interpreted as the genotype, governs the appearance

of the output designs, the phenotype. Adjusting the genotypes should move the corresponding phenotypes projected onto the GMM closer to a target style. Over several generations, the proposed genetic algorithm alters a population of N_{pop} latent vectors initialized from the distribution as the latent variables. The algorithm selects latent vectors for the next generation by evaluating their probability of falling into the targeted style cluster. The goal is to arrive at a set of latent vectors representing designs of the desired style. While the following sections explain the details, Algorithm 1 shows the pseudo-code for the procedure.⁷

Representation and transformation A generated design is represented by its latent vector $z = \langle v_1, \dots, v_l \rangle$ with $v \in \mathbb{R}$ consisting of l latent variables.

While an initial latent vector z is sampled randomly from the underlying distribution (see section 4.1), the genetic algorithm aims to transform z towards z^* with $G(z^*)$ representing a design of the targeted style. This transformation is guided by the fitness objective defined below.

Fitness and selection To arrive at the desired goal, individuals are evaluated against a fitness measure. Recall from section 4.2, that an image’s posterior probability of belonging to a style component t is p_t . Following the goal of resembling a certain target style t , the fitness criterion \mathcal{F}_t is defined by an individual’s posterior probability of belonging to the respective target style cluster, referred to as f_t . Applying the measure brought forward by the GMM, $f_t = p_t$.

To assess the fitness of a latent vector, image $G(z)$ is generated. The generated image is then projected onto the embedding space to retrieve $embedding(G(z))$. Finally, the $embedding(G(z))$ is projected onto the GMM, where the probability of belonging to target cluster t is extracted as p_t , representing the fitness criterion f_t . That defines the fitness f of an individual, or latent vector z , of belonging to a target cluster t as

$$f_t(z) = p_t(embedding(G(z))). \quad (2)$$

For the generated images to fit into a style cluster, the fitness \mathcal{F}_t needs to be maximized to obtain a latent vector z^* of target style t :

$$z^* = \arg \max_z \mathcal{F}_t(z) \quad (3)$$

By maximizing the fitness function, the population of latent vectors should improve towards the defined requirement through selection and variation.

At the beginning of each generation, we preserve a copy of the best N_{elite} individuals to save them from alteration. Inheriting them to the next round guarantees that we do not destroy the fittest individuals during a generation. N_{pop} individuals are selected by conducting N_{pop} tournaments, where N_{ts} randomly chosen individuals compete against each other based on their fitness. Two

⁷ The genetic algorithm was implemented using the evolutionary computation framework DEAP: <https://deap.readthedocs.io/en/master/index.html>

Algorithm 1: Evolutionary Search of GANs’ latent space using GMM.

Result: $G(z^*)$ closest to style cluster centroids

```

1 Function fitness( $z, t$ ):
2   | return  $p_i(\text{embedding}(G(z)))$ 
3 Function selection( $population, N$ ):
4   | for  $i \leftarrow 1$  to  $N$  do
5     |  $best_i \leftarrow$  winner out of  $N_{ts}$  randomly chosen  $z$  with fitness( $z, t$ )
6   | end
7   | return  $best_1, \dots, best_N$ 
8 Function crossover( $a, b$ ):
9   | for  $i$  in  $a$  do
10    |  $a_i = \alpha a_i + (1 - \alpha)b_i$  for  $\alpha \sim \text{Bernoulli}(0.5)$ 
11    |  $b_i = \alpha b_i + (1 - \alpha)a_i$  for  $\alpha \sim \text{Bernoulli}(0.5)$ 
12  | end
13  | return  $a, b$ 
14 Function mutation( $a$ ):
15  |  $noise \leftarrow$  vector of length  $l$  where  $noise_i \sim \mathcal{N}(\mu = 0, \sigma = 1)$ 
16  | return  $a + noise$ 
17  $population \leftarrow N_{pop} \times z$ 
18 for  $g$  in  $N_{gen}$  do
19  |  $elite \leftarrow N_{elite} \times z$  with maximum fitness ( $z, t$ ) in  $population$ 
20  |  $population = \text{selection}(population, N_{pop})$ 
21  | for  $a, b$  in  $population$  do
22    | if  $random < p_{cx}$  then
23      | |  $a, b = \text{crossover}(a, b)$ 
24    | end
25  | end
26  |  $population = population + N_{new} \times z$ 
27  | for  $a$  in  $population$  do
28    | if  $random < p_{mut}$  then
29      | |  $a = \text{mutation}(a)$ 
30    | end
31  | end
32  |  $population = population + elite$ 
33 end

```

kinds of variation operations, recombination and mutation, are applied to the population resulting from the tournament selection.

Recombination *Uniform crossover* is applied to generate new offspring as in [2,7]. Two latent vectors a and b are recombined to produce two new individuals \hat{a} and \hat{b} , with their i th attribute randomly chosen from either a or b :

$$\hat{a}_i = \alpha a_i + (1 - \alpha)b_i \text{ and } \hat{b}_i = \alpha b_i + (1 - \alpha)a_i \text{ with } \alpha = \text{Bernoulli}(0.5) \quad (4)$$

Table 2. GA parameters under variation.

Parameter	Setting
Crossover rate p_{cx}	0.7, 0.9
Mutation rate p_{mut}	0.2, 0.5
Population size p_{pop}	100, 200
Tournament size p_{ts}	3, 6

Table 3. Constant GA parameters.

Parameter	Setting
Size of individual	512
N_{gen}	500
N_{elite}	1
N_{new}	10
Recombination operator	uniform crossover ($\alpha = 0.5$)
Mutation operator	nonuniform mutation ($\mu = 0, \sigma = 1$)

If recombined, an individual is replaced by its child. The crossover rate p_{cx} defines the chance for an individual of the population to participate in recombination. As commonly applied, we choose high rates (0.7 and 0.9) to ensure the continuing development of fit individuals [11]. To introduce new gene material, we add N_{new} new random vectors to the population in every generation in addition to the crossover.

Mutation Following the objective of moving closer to a target cluster, we apply *nonuniform mutation* [2]. With a probability of 0.5, a variable of a latent vector z is mutated by adding some *noise* drawn from the original distribution:

$$z_i = z_i + noise \sim \mathcal{N}(\mu = 0, \sigma = 1) \quad (5)$$

The mutation rate p_{mut} defines the chance for an individual of the population to be mutated. While the mutation rate of a genetic algorithm is usually set to a few percent [11], we consider both low (0.2) and high (0.5) mutation rates, as initial runs showed low diversity among the population.

In accordance with typical tuning methods of evolutionary algorithms, we consider different tournament sizes ($N_{ts} = \{3, 6\}$) and population sizes ($N_{pop} = \{100, 200\}$) adhering to the commonly used parameter choices [5,11]. The parameters under variation are summarized in Table 2. The number of generations is set to $N_{gen} = 500$ as a compromise of running time and complexity of the problem. Table 3 displays the constant parameters. Due to the stochasticity underlying evolutionary systems, some runs naturally never achieve any fitness due to an ‘unlucky’ initialization of the population. Therefore, we test the model for different styles per parameter combination. From a random selection of styles presented to the experimenter, they chose five distinct ones to ensure visual diversity as a basis for the experiments.

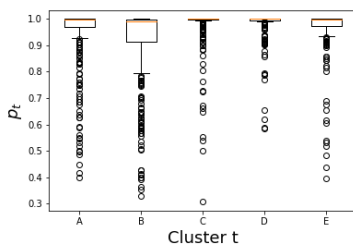


Fig. 3. Boxplot of the original images’ posterior probability for the five target styles chosen for the experiment. The five clusters range from 228 to 367 in size. The images contained by them have a mean p_t between 0.92 and 0.98.

Table 4. Average maximum fitness across all five runs per parameter combination.

		$p_{mut} = 0.2$		$p_{mut} = 0.5$	
		$N_{ts} = 3$	$N_{ts} = 6$	$N_{ts} = 3$	$N_{ts} = 6$
$p_{cx} = 0.7$	$N_{pop} = 100$	0.5746	0.2	0.6021	0.2119
	$N_{pop} = 200$	0.4031	0.5491	0.6341	0.4028
$p_{cx} = 0.9$	$N_{pop} = 100$	0.4538	0.3855	0.735	0.2982
	$N_{pop} = 200$	0.8073	0.2	0.7043	0.5248

5 Results

We tested our model in each parameter combination to generate designs for five different styles. Fig. 3 displays the distribution of the images’ posterior probability per target style cluster. A comparison of the mean maximum fitness reached across all five runs is presented in Table 4. As we ran the whole system for five different times, or style clusters, the results that we include are averages over those five runs. Recall that the fitness to be maximized is defined as an image’s posterior probability of belonging to a GMM component. Hence, it can vary between a minimum of 0 and a maximum of 1.

As the comparison of the results reveals, the algorithm finds individuals of highest fitness for a crossover rate of $p_{cx} = 0.9$ and a tournament size of $N_{ts} = 3$. In particular, the highest fitness is reached in combination with a population size of $N_{pop} = 200$ and a low mutation rate of $p_{mut} = 0.2$, followed by the second highest fitness with $N_{pop} = 100$ and $p_{mut} = 0.5$.

To better understand the fitness measure in relation to the produced designs, an exemplary case from the experiments for style cluster A, B, C, and D each is presented in Fig. 4. A plot of average and maximum fitness of each exemplary run is shown in Fig. 5 for inspection of the search behavior. For a generated design for the style displayed in Fig. 4A, the similarity to the target cluster is visible, as they share the features of an open denim jacket with bleached platts and darker pants. This is aligned with the high fitness of the generation. Interestingly, the random sampling achieved no fitness despite the similarity of the displayed pants. In comparison, some generated designs of lower fitness, such as the example in

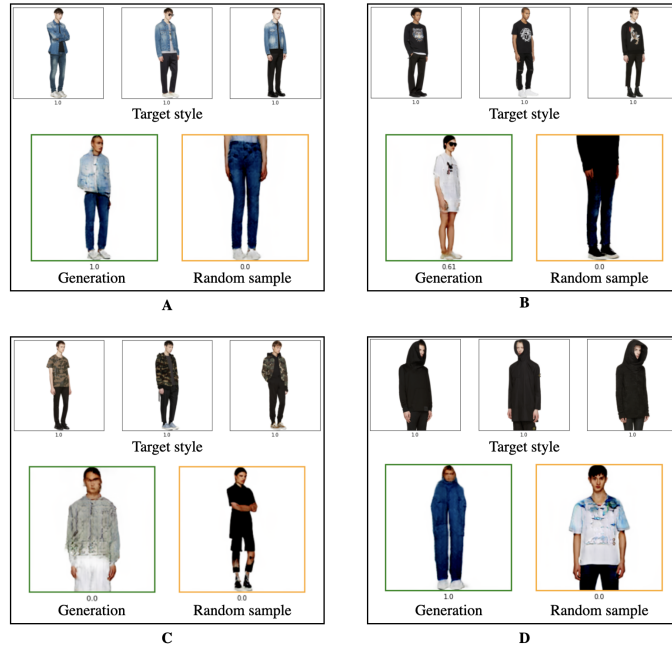


Fig. 4. Four exemplary designs for different clusters. See Fig. 5 for the parameter settings underlying the runs. The figure displays the three images with the highest posterior probability in the style cluster, the fittest generated image, and the fittest out of $N_{pop} \times N_{gen}$ randomly sampled images for comparison. The titles of the images display their fitness. Zoom in for detail.

Fig. 4B, show less visual similarity. Specifically, only one feature, namely the pattern on the chest, seems to resemble the target style characterized by a black outfit with a white pattern on the shirt.

The same concern arises for a generated design for the target style characterized by a camouflaged patterned upper part and black pants, shown in Fig. 4C. While the image shows an upper body with a pattern similar to the target style, it achieved no fitness, as the corresponding plot in Fig. 5C shows. While this could be caused by an unlucky initialization, it might also point to flaws in the fitness measure. That the assigned fitness does not seem to be in alignment with visual coherency is also the case for the example given in Fig. 4D. Here, high fitness is not aligned with a legit outfit design.

6 Discussion

For some of the analyzed designs, the fitness measure aligns with the visual coherency to the target style. In Fig. 4A, high fitness represents a design of high visual similarity. In contrast, lower fitness describes a design of little visual

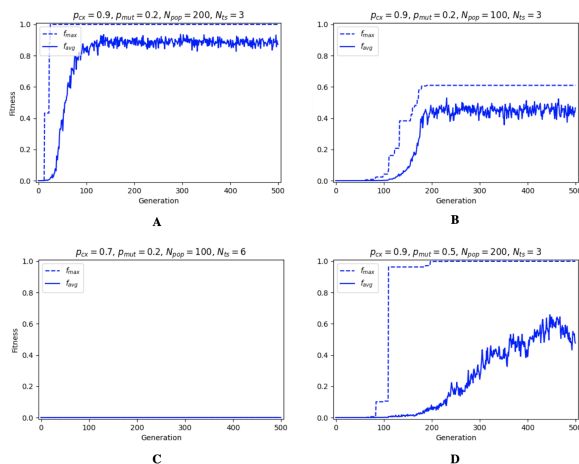


Fig. 5. Maximum and average fitness over generations of the exemplary runs in Fig. 4.

coherency in Fig. 4B. Taking a look at the search behavior behind the latter, displayed in Fig. 5B, an explanation could be the following. The population converged to a local maximum after around half of the generations, only allowing for optimizing a minor feature in the following generations. Though convergence to a local maximum is a desirable end stadium for genetic algorithms, it poses a problem in this scenario. Due to weak style matches dominating the population, the subsequent evolution underlies the prevailing look. A solution to introduce more diversity to the population could be to increase mutation or the number of new individuals when fitness stagnates.

For other examples, the fitness measure does not represent the visual coherency. Fig. 4C displays a design with a pattern similar to the target style, but achieved no fitness. This example raises the question of whether other factors, such as the models' posing, play a role in defining a style. At the same time, the example shown in Fig. 4D achieved the maximum fitness, for a design that a human judge might not assign validity. As our style model is based on the representation of an attribute model trained to recognize clothing features, it might lack the sensitivity of body shapes. However, coming with a background in clothing and physiology guides, such notions guide our perception of the designs. Hence, the suggested system might find some completely different attributes behind what we identify as decisive properties for a style. We hypothesize that the designs found by the evolutionary search reveal a machine-specific understanding of fashion style. However, to further investigate that claim, we need to find out whether our model is sufficient to capture styles. As part of our future work, further insight into the implementation and its parameters would show if they optimally support this goal.

The difference in perception could also be because computational networks are not able to capture subtle differences in styles (yet) [25]. Detecting style-

coherent similarity of clothing images remains a challenge, in which also the attribute models used to retrieve the visual representation, including the training data, play an essential role [12]. Even though the embedding of images was used instead of the ResNet’s final layer’s attribute prediction scores, the categories underlying its training influence what the model *sees* on every layer. To better understand the style-defining features, a qualitative analysis of the ResNet’s performance, e.g. with the help of a *Class Activation Map* (CAM), could be performed. Such an analysis could shed insight on the critical features for a particular fashion style [25]. With that knowledge, the latent space of the GAN might be controllable in a more targeted way.

To quantify the relation between the clusters found by the evolutionary search and human judgement, a human could be added to the loop. To guide the system into the direction of human-compatible styles, the approach could eventually be combined with interactive evolution like in DeepIE or StyleIE, in order to integrate how humans see style. Additionally, we could consider if exploring the latent space can even tell us more about what the model interprets behind styles, providing insight into clustering components. Following that path contributes to understanding how the proposed system sees visual themes emerge, eventually leading to a better understanding of its functioning.

As a final remark, it is essential to reflect on the datasets used for training generative models. As Takagi et al. [25] point out, fashion photographs do not represent what people actually wear, hence might not give an actual representation of current styles. Most datasets used for training generative fashion design models consist of catalog or social media images, highly biased to the presented populations. Therefore, a crucial task is to discuss how data affects the system’s stability with regard to design diversity, such as achieving desirable silhouettes while still considering diverse body forms.

7 Conclusion and Future Work

This research aimed to extend the classical generative deep learning approach to facilitate the generation of designs that respond to fashion styles. We investigated the application of a genetic algorithm and a GMM to guide the generation of images with regard to previously identified style clusters. Our suggested framework facilitates the search for images responding to certain style clusters. The experimental results indicate that the proposed procedure provides a promising direction to guide the search for style-coherent designs. Further research is required to establish a robust and reliable exploration process.

While the system can generate images of maximum fitness, the designs do not necessarily correspond to the target styles in the way one would expect. We hypothesize that the generations reveal a machine-specific understanding of fashion style. While some of the generated images exhibit similarity to the target cluster visible to the human eye, other generated outputs raise the question of how the algorithm might understand styles differently, outlining the need to improve the fitness measure.

Integrating fashion style analysis and fashion generation opens up new design possibilities, such as extending trends forecasting to generating trending designs. Future work could investigate how the ability to react to stylistic developments through unsupervised learning expands the capabilities of generative models as creative design tools.

References

1. Al-Halah, Z., Grauman, K.: Modeling Fashion Influence from Photos. *IEEE Transactions on Multimedia* (2020). <https://doi.org/10.1109/TMM.2020.3037459>
2. Bontrager, P., Lin, W., Togelius, J., Risi, S.: Deep interactive evolution. In: *International Conference on Computational Intelligence in Music, Sound, Art and Design*. pp. 267–282. Springer (2018)
3. Bontrager, P., Roy, A., Togelius, J., Memon, N., Ross, A.: Deepmasterprints: Generating masterprints for dictionary attacks via latent variable evolution. In: *2018 IEEE 9th International Conference on Biometrics Theory, Applications and Systems (BTAS)*. pp. 1–9. IEEE (2018)
4. Cheng, W.H., Song, S., Chen, C.Y., Hidayati, S.C., Liu, J.: *Fashion Meets Computer Vision: A Survey* (2021)
5. Eiben, A.E., Smith, J.E., et al.: *Introduction to evolutionary computing*, vol. 53. Springer (2003)
6. Elgammal, A., Liu, B., Elhoseiny, M., Mazzone, M.: CAN: Creative Adversarial Networks, Generating "Art" by Learning About Styles and Deviating from Style Norms. *arXiv preprint arXiv:1706.07068* (2017)
7. Fernandes, P., Correia, J., Machado, P.: Evolutionary latent space exploration of generative adversarial networks. In: *International Conference on the Applications of Evolutionary Computation (Part of EvoStar)*. pp. 595–609. Springer (2020)
8. Ge, Y., Zhang, R., Wang, X., Tang, X., Luo, P.: Deepfashion2: A versatile benchmark for detection, pose estimation, segmentation and re-identification of clothing images. In: *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. pp. 5332–5340 (2019). <https://doi.org/10.1109/CVPR.2019.00548>
9. Goodfellow, I.J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: *Generative Adversarial Networks*. *arXiv preprint arXiv:1406.2661* (2014)
10. Grabe, I.: *Evolutionary Search for Fashion Styles in the Latent Space of Generative Adversarial Networks*. Master's thesis, University of Copenhagen (2021)
11. Hassanat, A., Almohammadi, K., Alkafaween, E., Abunawas, E., Hammouri, A., Prasath, V.: Choosing mutation and crossover ratios for genetic algorithms—a review with a new dynamic approach. *Information* **10**(12), 390 (2019)
12. Hsiao, W.L., Grauman, K.: Learning the Latent "Look": Unsupervised Discovery of a Style-Coherent Embedding from Fashion Images (2017)
13. Jiang, S., Li, J., Fu, Y.: Deep learning for fashion style generation. *IEEE Transactions on Neural Networks and Learning Systems* (2021)
14. Karras, T., Aila, T., Laine, S., Lehtinen, J.: Progressive Growing of GANs for Improved Quality, Stability, and Variation. *arXiv preprint arXiv:1710.10196* (2017)
15. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 4401–4410 (2019)

16. Kato, N., Osone, H., Oomori, K., Ooi, C.W., Ochiai, Y.: Gans-based clothes design: Pattern maker is all you need to design clothing. In: Proceedings of the 10th Augmented Human International Conference 2019. pp. 1–7 (2019)
17. Kiapour, M.H., Yamaguchi, K., Berg, A.C., Berg, T.L.: Hipster wars: Discovering elements of fashion styles. In: European conference on computer vision. pp. 472–488. Springer (2014)
18. Liu, Z., Luo, P., Qiu, S., Wang, X., Tang, X.: Deepfashion: Powering robust clothes recognition and retrieval with rich annotations. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 1096–1104 (2016). <https://doi.org/10.1109/CVPR.2016.124>
19. Mackinney-Valentin, M.: On the Nature of Trends: A Study of Trend Mechanisms in Contemporary Fashion. Ph.D. thesis, Royal Danish Academy (Jun 2010)
20. Mahmood, A., Ospina, A.G., Bennamoun, M., An, S., Sohel, F., Boussaid, F., Hovey, R., Fisher, R.B., Kendrick, G.A.: Automatic hierarchical classification of kelps using deep residual features. *Sensors* **20**(2), 447 (2020)
21. Marchetti, F., Wilson, C., Powell, C., Minisci, E., Riccardi, A.: Convolutional generative adversarial network, via transfer learning, for traditional scottish music generation. In: International Conference on Computational Intelligence in Music, Sound, Art and Design (Part of EvoStar). pp. 187–202. Springer (2021)
22. Matzen, K., Bala, K., Snavely, N.: StreetStyle: Exploring world-wide clothing styles from millions of photos (2017)
23. Rostamzadeh, N., Hosseini, S., Boquet, T., Stokowiec, W., Zhang, Y., Jauvin, C., Pal, C.: Fashion-Gen: The Generative Fashion Dataset and Challenge. arXiv preprint arXiv:1806.08317 (2018)
24. Roziere, B., Teytaud, F., Hosu, V., Lin, H., Rapin, J., Zameshina, M., Teytaud, O.: EvolGAN: Evolutionary Generative Adversarial Networks. In: Proceedings of the Asian Conference on Computer Vision (2020)
25. Takagi, M., Simo-Serra, E., Iizuka, S., Ishikawa, H.: What makes a Style: Experimental Analysis of Fashion Prediction. In: Proceedings of the IEEE International Conference on Computer Vision Workshops. pp. 2247–2253 (2017)
26. Tejada-Ocampo, C., López-Cuevas, A., Terashima-Marin, H.: Improving deep interactive evolution with a style-based generator for artistic expression and creative exploration. *Entropy* **23**(1) (2021). <https://doi.org/10.3390/e23010011>, <https://www.mdpi.com/1099-4300/23/1/11>
27. Xin, C., Arakawa, K.: Object design system by interactive evolutionary computation using gan with contour images. In: International Conference on Human-Centered Intelligent Systems. pp. 66–75. Springer (2021)
28. Yildirim, G., Jetchev, N., Vollgraf, R., Bergmann, U.: Generating high-resolution fashion model images wearing custom outfits. In: Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops. pp. 0–0 (2019)
29. Yildirim, G., Seward, C., Bergmann, U.: Disentangling Multiple Conditional Inputs in GANs. arXiv preprint arXiv:1806.07819 (2018)
30. Zhu, J., Liapis, A., Risi, S., Bidarra, R., Youngblood, G.M.: Explainable ai for designers: A human-centered perspective on mixed-initiative co-creation. In: 2018 IEEE Conference on Computational Intelligence and Games. pp. 1–8. IEEE (2018)
31. Zhu, J., Ontañón, S.: Shall i compare thee to another story?—an empirical study of analogy-based story generation. *IEEE Transactions on Computational Intelligence and AI in Games* **6**(2), 216–227 (2013)
32. Zhu, S., Urtasun, R., Fidler, S., Lin, D., Change Loy, C.: Be your own Prada: Fashion Synthesis with Structural Coherence. In: Proceedings of the IEEE international conference on computer vision. pp. 1680–1688 (2017)