

000
001
002
003
004
005
006
007
008
009
010
011054
055
056
057
058
059
060
061
062
063
064
065

Abstract

Attribute editing has become an important and emerging topic of computer vision. In this paper, we consider a task: given a reference garment image A and another image B with target attribute (collar/sleeve), generate photo-realistic image which combines the texture from reference A and the new attribute from reference B. There are two main difficulties in this task. First, it is hard to edit single attribute while keeping other highly entangled attributes unchanged. Second, many existing image generation methods fail to handle this task since there is no paired ground truth data. To overcome the those limitations, we propose a novel self-supervised model to synthesize garment images with disentangled attributes (e.g., collar and sleeves) without paired data. Our method consists of two training step: reconstruction learning and adversarial learning. Our model learns texture and location information through reconstruction learning. And the model capability is generalized to achieve single-attribute manipulation by adversarial learning. Current available datasets are mostly made of street images which include user's face and body parts but the gender and body shapes may entangle with the garment designs and textures. Therefore, we compose a new dataset, named GarmentSet, with annotation of landmarks of collar and sleeves on clean garment images. Thoughtful experiments on this dataset and real-world samples demonstrate that our method can synthesize significantly better results than the state-of-the-art methods in both quantitative and qualitative comparisons.

1. Introduction

Deep generative techniques [8, 15] have led to highly successful image/video generation, some focusing on style transfer [38], and others on synthesis of desired conditions [17, 26]. We propose a novel schema to disentangle attributes and synthesize high quality image with desired attribute while keeping other attribute unchanged. In this paper, we focus on the problem of fashion image attribute manipulation to demonstrate the capability of our method.



Figure 1. A graphical demonstration of our task and result. Given a reference fashion garment (on the left) and a desired design attribute (in the middle), we aim to generate a new fashion garment that seamlessly integrates the desired design attribute to the reference image. The first row shows collar-editing and the second row shows sleeve-editing. The results generated by our system are shown on the right.

By our method, users can switch a certain part of garments to the wanted designs (e.g., round collar to V-collar). The objective is to synthesize a photo-realistic new fashion image by combining different parts seamlessly together. Potential applications of such system range from item retrieval to professional fashion design assistant. In Fig. 1, we make a graphical demonstration of our task and result.

Current approaches useful for this task include conditional Generative Adversarial Networks (cGAN) [8, 17, 26], image-to-image translation [10], and CycleGAN [38]. While these models have been proved effective in generating photo-realistic images, image syntheses using such models usually involve highly entangled attributes/objects and may fail in editing target attribute/object separately [20, 31, 38]. Furthermore, the large variance in the garment textures causes additional problems. It is impossible to build a dataset that is large enough to approximate the distribution of all garment texture and design combinations, which serve as paired learning examples to train models such as [26, 10]. Novel learning paradigm is expected to overcome this difficulty.

To solve these challenges, we propose a novel self-supervised image generative model that can make disentangled attribute manipulations. The model is trained to handle the image translating tasks with no training paired images required. This new model, named TailorNet, which

108 exploits the latent space expression of input images. The
109 encoder-decoder architecture has been adapted to different
110 image synthesis tasks [26, 9, 16, 10]. We make improve-
111 ments to solve the texture-design entanglement problems in
112 our task. To isolate the structures from textures, we use the
113 image edge maps as inputs. By using edge maps instead
114 of RGB color images, we achieve good results with only
115 a small amount of data. We note that comparing to meth-
116 ods that explore the latent codes in GANs, such considera-
117 tion is explicit and data-parsimonious and achieves fashion
118 attribute-editing that is robust to various geometric trans-
119 formations between the reference and design attribute im-
120 ages. The model capacity is further generalized in a GAN
121 framework to achieve single-attribute manipulation using
122 random fashion inputs. Besides the reconstruction learning
123 and the adversarial learning, a attribute-aware discriminator
124 is weaved into the model to guide the image editing. This
125 attribute-aware discriminator helps in making high-quality
126 single attribute editing and guides a better self-supervised
127 learning process.

128 Currently available fashion image datasets (e.g., Deep-
129 Fashion [14]) are mostly made of street photos with com-
130 plex backgrounds or with user’s body parts presented. That
131 extra visual information may hinder the performance of
132 an image synthesis model. Thus, to simplify the training
133 data and screen out noisy backgrounds, we introduce a new
134 dataset, GarmentSet. The new dataset contains fashion im-
135 ages with no human user presented, and most images have
136 single color backgrounds.

137 Contributions made in this paper can be summarized as:

- 139 • We propose a new task in deep learning fashion studies.
140 Instead of generating image with text guidance,
141 virtual try on, or texture transferring, our task is to
142 make new fashion designs with disentangled user-
143 appointed attributes.
- 145 • We exploit a novel training schema consists of recon-
146 structive learning and adversarial leaning. The self-
147 supervised reconstructive learning guides the network
148 to learn shape, location information from edge map
149 and texture information from RGB image. The un-
150 paired adversarial learning gives the network gener-
151 alizability to synthesize image with new disentangled
152 attributes. Besides the reconstructive learning and
153 the adversarial learning, we propose a novel attribute-
154 aware discriminator, which helps high quality attribute
155 editing by isolating the design structures from the gar-
156 ment textures.
- 158 • A new dataset, GarmentSet, is introduced to serve our
159 attribute editing task. Unlike existing fashion datasets
160 in which most images illustrate the user’s face or body
161 parts with complex backgrounds, GarmentSet filters

162 out most redundant information and directly serves the
163 fashion design purpose.

164 The rest of this paper is organized as follows: A brief re-
165 view of related work is presented in Sec. 2. The details of
166 the proposed method are described in Sec. 3. We introduce
167 our new dataset in Sec. 4. Experimental details and results
168 are presented in Sec. 5. In Sec. 5.6, we further conduct ab-
169 lation studies to explicitly investigate the performances of
170 our model. And finally, Sec. 6 concludes the paper with
171 discussions of limitations.

2. Related Work

172 **Generative Adversarial Network (GAN)** [8] is one of the
173 most popular deep generative models and has shown im-
174 pressive results in image synthesis studies, like image edit-
175 ing [24, 30] and fine-grained objects generating [26, 13].
176 Training GAN based on conditions incorporates further
177 information to guide the image generating process. Re-
178 searchers utilize different conditions to generate images
179 with desired properties. Existing works have explored var-
180 ious conditions, from category labels [23], audio [6, 5],
181 text [26], skeleton[21, 11, 37, 25] to attributes [28]. There
182 are a few studies that investigate the task of image trans-
183 lations using cGAN [17, 26, 4]. In the context of fashion-
184 related applications, researchers apply cGAN in automated
185 garment textures filling [33], texture transferring [12] and
186 virtual try-on [39, 32] by replacing dress on a person with
187 a new one. A more related work is sequential attention
188 GAN proposed by Cheng et al. [7]. Their model uses text as
189 the guidance and continuously changes the fashion designs
190 based on user’s requests, but the attribute changes are highly
191 entangled. In contrast to this work, we propose a new train-
192 ing algorithm that combining self-supervised reconstruction
193 learning with adversarial learning to make disentangled at-
194 tribute manipulations with user-appointed images.

195 **Self-supervised generation** is recently introduced as a
196 novel and effective way to train generative models with-
197 out paired training data. Unpaired image-to-image trans-
198 lation framework such as CycleGAN [38] removes pixel-level
199 supervision. In CycleGAN, the unpaired image to image
200 translation is achieved by enforcing a bi-directional trans-
201 lation between two domains with an adversarial penalty on
202 the translated image in the target domain. The CycleGAN
203 variants [36, 18] are moving towards the direction of unsu-
204 pervised learning approaches. However, CycleGAN-family
205 models also create unexpected or even unwanted results
206 which we will show in our experiments. One reason for
207 such a phenomenon is due to the lack of straightforward
208 knowledge of the target translation domain in the circularity
209 training process. Inherent attributes of the source samples
210 may be changed in a translation process. To avoid such un-
211 wanted changes, we keep an image reconstruction penalty

216 in our image editing task.
 217

218 Attracted by the huge profit potentials in the fashion
 219 industries, numbers of researches using deep learning meth-
 220 ods have been conducted on fashion analysis and fashion
 221 image synthesis. Most existing researches focus on fash-
 222 ion trend prediction [2], clothing recognition with land-
 223 marks [14], clothing matching with fashion items in street
 224 photos [19] and fashion recommendation system [3, 34].
 225 Different from those research lines, we focus on fashion im-
 226 age synthesis task with single attribute manipulations.
 227

228 3. Methodology

229 This section presents implementation details of our
 230 method. There are two crucial steps in the model training:
 231 (1) self-supervised reconstruction learning and (2) general-
 232 ized attribute manipulations using adversarial learning. The
 233 model learns how to fill the correct texture and to locate
 234 the fashion pattern at the correct position. The second step
 235 helps the model generating high quality images with desired
 236 attributes. Although we use collar translating example in
 237 Fig. 2, we emphasize here that our model can be applied to
 238 attributes other than collar parts. We will show the sleeve
 239 editing results in the later section.
 240

241 3.1. TailorNet: Learning to Manipulate Designs

242 **Self-Supervised Learning Step.** The motivation of
 243 formulating a self-supervised model is the fact that it is almost
 244 impossible to collect paired training images for a fully sup-
 245 pervised model. Using the collar editing task as an exam-
 246 ple, for each image in a fully supervised training process,
 247 one needs to collect paired images for each collar type. In
 248 these paired images, only the collar parts are different while
 249 the other attributes, like body decorations, clothing textures,
 250 etc., must stay unchanged and match with other paired im-
 251 ages. Such data is usually unavailable. Also, the dataset
 252 size will increase exponentially when multiple attribute an-
 253 notations are needed for each image.
 254

255 Based on motivations discussed above, we employ an
 256 encoder-decoder structure for the self-supervised recon-
 257 struction training step. Given a masked garment image I^M
 258 (mask out the collar part) and edge map E^O , our recon-
 259 struction step reconstructs the original garment image (collar
 260 region). From daily experiences, certain fashion designs may
 261 highly entangle with textures or colors. For instances, light
 262 pink color is rarely used on men garment, leather is usually
 263 used to make jackets, etc. In our task, we want our model
 264 focusing only on the design structure editing rather than
 265 the clothing texture translating, which is inherited from the
 266 reference garment image. Specifically, the self-supervised
 267 learning step is defined as:

$$268 \hat{I}^R = SAM(\Psi(\Phi_{img}(I^M) \oplus \Phi_{edge}(E^O)), I^M), \quad (1)$$

270 where Φ_{img} and Φ_{edge} are image encoder and edge encoder,
 271 respectively. Φ_{img} and Φ_{edge} consist of several 2D con-
 272 volution layers and residual blocks. Ψ is the image decoder,
 273 which consists of several 2D transpose-convolution layers.
 274 \oplus is channel-wise concatenation. After encoding, we feed
 275 the concatenated latent vector to Ψ to output attention mask
 276 m and new pixel C . The SAM block (see Eq.3) outputs the
 277 reconstructed image \hat{I}^R based on m , C and I^M . Basically,
 278 this learning step learns how to reconstruct I^O according to
 279 the texture feature from I^M and shape feature from E^O .
 280

281 Different from other methods [6, 5], we only use percep-
 282 tual loss, which is computed by taking the L1 distance of
 283 the VGG features (first 16 layers) for this step, which yields
 284 better results (see Fig. 10). Specifically, the loss function
 285 for this training step is defined as:

$$286 \mathcal{L}_R = \mathcal{L}_{VGG}(I^O, \hat{I}^R) \\ 287 = \mathbb{E}_{\hat{I}^R} [\|\phi(I^O) - \phi(\hat{I}^R)\|_1^1], \quad (2)$$

288 where ϕ is a feature extractor pretrained from image clas-
 289 sification [29]. From the reconstruction training, the gener-
 290 ator learns to fill the garment texture and allocates the de-
 291 sired part at the correct position to output the original un-
 292 masked image I^O . In our empirical study, we observe that
 293 this learning step is critical to the full model since it learns
 294 how to synthesize collar part by reconstruction. In this step,
 295 the model learns how to do texture matching and do pattern
 296 locating. During the training process, we apply random ro-
 297 tations, translations and scale shift to the input edge maps.
 298 Thus, the model can be trained to handle potential geometric
 299 transformations and can allocate the desired fashion pattern
 300 at the correct position.
 301

302 **Self-Attention Mask Operation (SAM).** During the re-
 303 construction step, the model should learn to change only the
 304 target part and keep the rest parts of an image untouched.
 305 Thus, we introduce a self-attention mask to the generator.
 306 The self-attention mechanism can guide the model to focus
 307 on the target region. This subsection reveals how the self-
 308 attention mask helps in making high quality results.
 309

310 After the edge map and the masked image are encoded,
 311 the latent space vectors are concatenated and then are de-
 312 coded by the decoder network. The decoder produces two
 313 outputs: single channel self-attention mask m and new
 314 pixel C . The final output of the generator combines the
 315 masked color image with the input cropped image I^M . The
 316 combination step follows the equation:

$$317 \hat{I}^R_{i,j} = m_{i,j} \times C_{i,j} + (1 - m_{i,j}) \times I^M_{i,j}, \quad (3)$$

318 where $m_{i,j}$, $C_{i,j}$ and $\hat{I}^R_{i,j}$ are the pixel at i^{th} row and j^{th}
 319 column in the self-attention mask, the new pixel and the
 320 final output image. The self-attention mask layer and the
 321 color layer share the bottom transpose convolutional blocks
 322

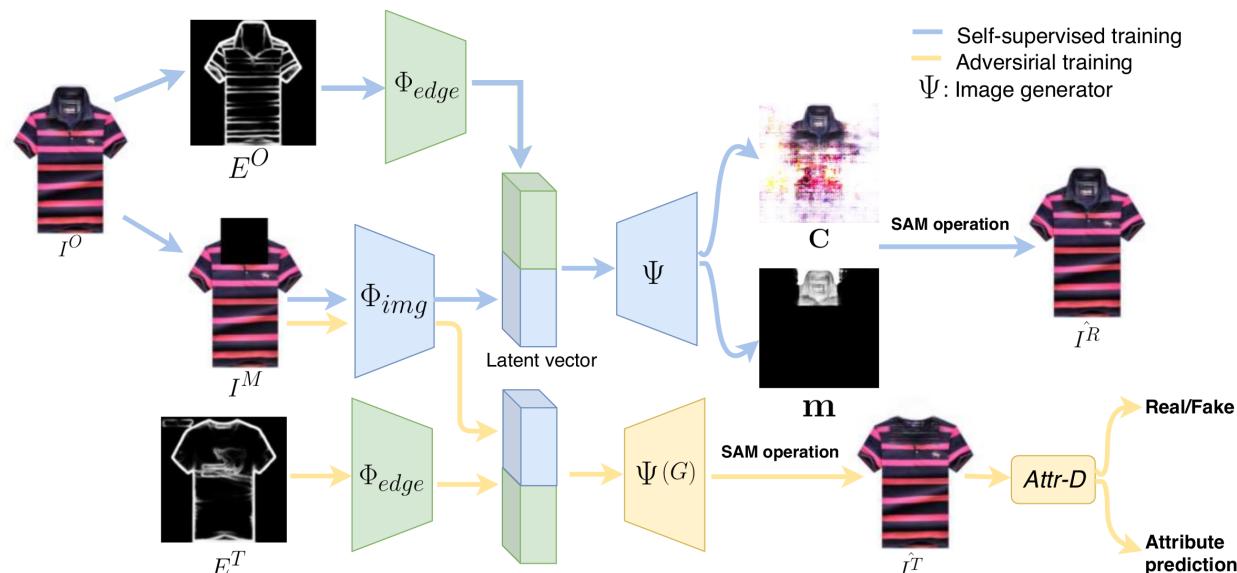


Figure 2. Network architecture of TailorNet. In the upper part (step 1: self-supervised reconstruction), we extract the edge feature (location and shape feature) by Φ_{edge} and extract image feature (texture feature) by Φ_{img} . Then we merge the two vector in latent space and pass through image generator Ψ to output mask and attention. In the lower part (step2: adversarial training), we extract the shape feature of target edge map E^T using Φ_{edge} and extract texture feature using Φ_{img} . The attribute discriminator (Attr-D) will output real/fake score and the attribute class of the fake image \hat{I}^T yield by Ψ . For better understanding the network details, we will release the code once the paper is accepted.

in the decoder. The output of the last transpose convolutional layer is fed into two activation layers: a Sigmoid layer with single-channel output (self-attention mask) and a hyperbolic tangent layer with three-channels output (the new pixel). The self-attention mask guides the network focusing on the attribute related region while training the network in a fully self-supervised manner.

3.2. Generalized Single Attribute Manipulations

We introduced the self-supervised reconstruction learning step in Sec. 3.1, which can reconstruct original image. However, our task is synthesizing new images by manipulating the attributes. The model trained with reconstruction step can not yield good results since it is not generalizable to synthesize new image with other new attribute types (e.g., new collar type and new sleeve type). Meanwhile the synthesized image with reconstruction step is blurry, which makes the results unrealistic. To tackle those problems, we have another adversarial learning step, which consists of the encoder-decoder network(see Sec. 3.1) and a novel attribute-aware discriminator.

In order to enforce the model to output image with correct attributes, we propose a discriminator with two different regression scores: a binary (real/fake) label and an attribute prediction vector. The attribute prediction vector is further optimized by cross entropy loss, which is defined as:

$$\ell_{attr}(I, \vec{V}^{T/O}) = -\vec{V}^{T/O} * \log(f(I)) - (1 - \vec{V}^{T/O}) * \log(1 - f(I)) , \quad (4)$$

where the vector $\vec{V}^{T/O}$ is the class vector of target attribute and f is a attribute classifier which outputs the class label vector of image I . During adversarial training, when we forward I^O and \vec{V}^O to discriminator, the parameters of discriminator will be updated to learn how to classify the attribute; when we forward \hat{I}^T and \vec{V}^T to discriminator, we will not update the parameters in discriminator, but update the parameters in generator according to the output attribute label. Thus, the full GAN loss can be formulated as:

$$\begin{aligned} \mathcal{L}_{GAN} = & \mathbb{E}_{\{E^T, I^M\}} [(1 - D(G(E^T, I^M))^2] + \\ & \frac{1}{2} \mathbb{E}_{\{E^T, I^M\}} [D(G(E^T, I^M))^2] + \\ & \frac{1}{2} \mathbb{E}_{\{I^O\}} [(D(I^O) - 1)^2] + \\ & \lambda_1 * \mathbb{E}_{\{E^T, I^M\}} [\ell_{attr}(G(E^T, I^M), \vec{V}^T)] + \\ & \lambda_1 * \mathbb{E}_{\{I^O\}} [\ell_{attr}(I^O, \vec{V}^O)] , \end{aligned} \quad (5)$$

where the λ_1 is a hyper-parameter to balance the loss terms. We set $\lambda_1 = 0.1$ in all experiments. Besides the GAN loss, we also apply perceptual loss in this training step. Thus, the loss function for this adversarial learning step is defined as:

$$\mathcal{L}_{adv} = \mathcal{L}_{GAN} + \lambda_2 * \mathcal{L}_{VGG} , \quad (6)$$

where λ_2 is a hyper-parameter. We set $\lambda_2 = 0.5$ to balance the loss terms. In our emppirical study, we observe that the



Figure 3. Samples of each collar type and sleeve type from GarmentSet dataset. The pie charts demonstrate the collar type and sleeve type distribution. There are in total 12 collar types and 2 sleeve types. That is the data distribution used in this paper. We will release the dataset once the paper is accepted.

Algorithm 1 Training steps

```

Require:  $\alpha$  is the learning rate.  $B$  is the batch size.
Require:  $\theta_G$  generator parameter.  $\theta_D$  is the discriminator parameter.
for number of iterations do
    Sample  $\{I_i^O, E_i^T, E_i^O, I_i^M, \vec{V}_i^T\}_{i=1}^B$ 
    Updating the generator G in reconstruction step:
     $\{I_i^R\}_{i=1}^B \leftarrow G_\theta(\{E_i^O\}_{i=1}^B, \{I_i^M\}_{i=1}^B)$ 
     $\theta_G \leftarrow Adam\{\nabla_{\theta_G}(\frac{1}{B} \sum_{i=1}^B (\mathcal{L}_R(I_i^O, I_i^R), \alpha))\}$ 
    Updating the discriminator D in adversarial step:
     $I_i^T \leftarrow G_\theta(E_i^T, I_i^M)$ 
     $\theta_D \leftarrow Adam\{\nabla_{\theta_D}(\frac{1}{B} \sum_{i=1}^B (\mathcal{L}_{GAN}(I_i^O, I_i^T, \vec{V}_i^T), \alpha))\}$ 
    Updating the generator G in adversarial step:
     $I_i^T \leftarrow G_\theta(E_i^T, I_i^M)$ 
     $\theta_G \leftarrow Adam\{\nabla_{\theta_G}(\frac{1}{B} \sum_{i=1}^B (\mathcal{L}_{adv}(I_i^O, I_i^T, \vec{V}_i^T), \alpha))\}$ 

```

model is sensitive to λ_2 and we need to choose different λ_2 if we use different VGG layer feature to compute perceptual loss.

3.3. Training Algorithm

Combining the self-supervised reconstruction learning step (Sec. 3.1) with the adversarial learning step (Sec. 3.2), our full model is trained in an alternative fashion. Specifically, we formulate the training algorithm in Algorithm 3.3. When training the generator G in reconstruction step with discriminator D fixed, the generator G will receive masked image I^M and an original attribute type E_o as input and it outputs the reconstructed image \hat{I}^R . We try to minimize the reconstruction loss \mathcal{L}^R to enforce the network to learn to generate correct texture and put it to geometry location. When training the discriminator D in adversarial step, the generator G will receive masked image I^M and a new attribute type E^T as input and it outputs the edited image \hat{I}^T . We try to minimize the GAN loss (\mathcal{L}_{GAN}) since there is no paired ground truth in this step. This step will enforce the

network to learn to synthesize images by manipulating the target attribute type E^T . Then we will update parameter in generator G in adversarial step by minimize the adversarial loss (\mathcal{L}_{adv}). By optimizing the loss in such a iterative manner, the TailorNet is able to learn realistic texture and geometry information to yield high-quality images with new attribute type.

4. GarmentSet dataset

This part serves as a brief introduction to GarmentSet dataset. Currently available datasets like DeepFashion [14] and FashionGen [27] are mostly made of images including user's face or body parts and street photos with noisy backgrounds. The redundant information raise unwanted hardness in the training process. To filter out such redundancy information in the images, we build our own dataset which will be published later.

In the dataset, we have 9636 images with collar part annotations and 8616 images with shoulder and sleeve annotations. Both classification types and landmarks are recorded. Although, in our studies, the landmark locations are only used in the image pre-processing steps, they still can be useful in the future researches like fashion item retrieve, clothing recognition, etc. The dataset keeps growing, and more attribute annotations will be added.

In Fig. 3, we present sample pictures for each collar type, each sleeve type and the overall data distributions in the dataset. Round collar, V-collar and lapel images together contribute over eighty percents of the total collar-annotation images. The sleeve-dataset only contains two types: short and long sleeves. Although not used in the training, the dataset also contains attribute landmark locations including collars, shoulders and sleeve ends. We keep collecting more data images and adding new attribute annotations. This distribution may change in the published version.

432 486
433 487
434 488
435 489
436 490
437 491
438 492
439 493
440 494
441 495
442 496
443 497
444 498
445 499
446 500
447 501
448 502
449 503
450 504
451 505
452 506
453 507
454 508
455 509
456 510
457 511
458 512
459 513
460 514
461 515
462 516
463 517
464 518
465 519
466 520
467 521
468 522
469 523
470 524
471 525
472 526
473 527
474 528
475 529
476 530
477 531
478 532
479 533
480 534
481 535
482 536
483 537
484 538
485 539

540
541
542
543
544
545

Type	Type 1 \Rightarrow Type 2			Type 2 \Rightarrow Type 1			Type 1 \Rightarrow Type 6			Type 6 \Rightarrow Type 1			Type 2 \Rightarrow Type 6			Type 6 \Rightarrow Type 2		
	C.E.	SSIM	PSNR															
CycleGAN	12.48	0.77	18.72	1.74	0.64	13.97	5.21	0.74	23.40	2.97	0.78	18.77	6.02	0.89	18.89	10.02	0.77	17.63
Pix2pix	20.01	0.87	22.75	12.73	0.89	21.42	21.62	0.88	17.63	15.79	0.89	21.88	15.12	0.77	23.15	19.67	0.88	23.04
Ours	11.04	0.93	25.52	1.20	0.91	24.08	8.34	0.92	25.68	2.44	0.93	24.06	6.57	0.92	24.44	6.89	0.92	23.88

594
595
596
597
598
599
600
601

Table 1. Measurements for all three models based on target translating types on testing set. In the table C. E. column is the average classification error scores for each paired type translation. The bold numbers in each column are the best scores.



Figure 4. Testing results of collar editing. The images are clustered based on the target collar types. The generated results of simple collar patterns (round and V-collar) are in general better than complicated ones (lapel). Our model is a single model and CycleGAN models are separated models trained with specific transformation.

	C.E.	SSIM	PSNR
Our model	6.78	0.9386	24.51
pix2pix	16.12	0.8923	22.83

Table 2. Comparing TailorGAN to pix2pix trained on all collar type inputs. The numbers are tested on testing set.

602
603
604

translations, we drop the cycleGAN model in this random translating comparison test. In the testing results, our models out-performances pix2pix model in all collar-type translating tasks. To compare with [38], we have trained three different CycleGAN models: collar type 1 (round collar) \Leftrightarrow type 2 (V-collar), collar type 1 (round collar) \Leftrightarrow type 6 (lapel) and collar type 2 \Leftrightarrow type 6. But our model is trained with all different types together.

Qualitative results. In Fig. 4 , we present testing results of three models. As one can see in the sample images, CycleGAN does not preserve garment textures. The trained pix2pix model performs badly in all examples. At the collar part, the pix2pix outputs only show color bulks with no structural patterns.

Quantitative results. For the quantitative comparisons, in measuring model performances, we use three metrics: classification cross entropy errors (C.E.), structure similarity index (SSIM) and peak signal to noise ratio (PSNR). The classification error is measured with a classifier pre-trained on GarmentSet dataset. We use a classification error since there is no paired testing images for the edited results. The classification error can measure the distance from the target collar designs. The SSIM and the PSNR scores are derived from the differences between the original image and the edited image. From the numerical results presented in Table.1, our model outperforms both CycleGAN and pix2pix in making high quality images with higher classification accuracy. We attribute this to the poor texture preserving ability of the CycleGAN/pix2pix model.

5.3. Synthesizing Unseen Collar Type

To test our model's capability of processing collar types that are missing in the dataset, we take one collar type out in the training stage and test the model's performance on this unseen collar type. In this test, we take collar type 1, 2 and 6 out. We also test the taken-one-out model with a fully trained model which meets all collar types in its' training process. The classification errors for each pair of models are calculated for each taken out collar type on testing set. The qualitative result is shown in Fig.5

605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

5. Experiments

5.1. Data Pre-processing

In this paper, we randomly sample 80% data for training and 20% for testing. Meanwhile, in Sec. 5.3, we keep one collar type out in the training set to demonstrate the robustness of our model. In the data pre-processing step, we generate mask-out images I^M and edge maps of data images. We use the left/right collar landmark locations to make a bounding box around the collar region. The bounding box size changes over landmark locations. A pretrained holistically edge detection (HED) model [35] takes charge of making edge maps. The HED model is trained on the BSDS dataset [22].

We notice that the image quality of the edited results depends on the input edge map resolution. Since the HED model used is pretrained on a general image dataset, it struggles in catching structural details and may also generate unwanted noisy maps. The HED model produces low-resolution edge maps for a target image with complicated detailed structures.

5.2. Comparative Studies

For the comparative studies, we compare our model with CycleGAN [38] and pix2pix [10]. In Table 2, our model is compared with pix2pix using random collar types. Since cycleGAN model can not handle two specific collar type



Figure 5. The testing results of unseen target collar type. The left side indicate which type we are generating (remove this type in training set). The 1th, 4th columns are the original reference images. The 2th, 5th columns are images with target collar types. The 3th, 6th columns are generated images with target collar types with texture/style of reference image.

C. E.	type 1 out	type 2 out	type 6 out
full model	2.05	7.51	9.27
one-out model	2.65	9.21	10.34

Table 3. Comparing one-out models to a fully trained model on testing set.

Based on the qualitative analysis and quantitative comparisons (see Tab.3), our model shows strong generalizing ability in synthesizing unseen collar types.

5.4. Sleeve Generation

In the previous discussions, we applied our model in collar part editing. GarmentSet dataset also contains sleeve landmarks and types information. Thus, we test our model's capability in editing sleeves and present testing results in Fig. 6. We used the same training scheme. Instead of collars, the sleeve parts are masked out. Due to simpler edge structures and better resolutions in the edge maps, the edited sleeve images have better image qualities and are close to the real images. We will discuss more details of sleeve editing results in the user evaluation section.

5.5. User Evaluations and Item Retrieves

To evaluate the performance in a human perceptive level, we conduct thoughtful user studies in this section. Human subjects evaluation (see Fig.7) is conducted to investigate the image quality and the attribute (collar) similarity of our generated results compared with [38, 10]. Here, we present the average scores for each model based on twenty users'



Figure 6. Testing results of editing sleeves using TailorGAN. The 1th, 4th columns are the original reference images. The 2th, 5th columns are images with target sleeve types. The 3th, 6th columns are generated images with target sleeve types with texture/style of reference image.

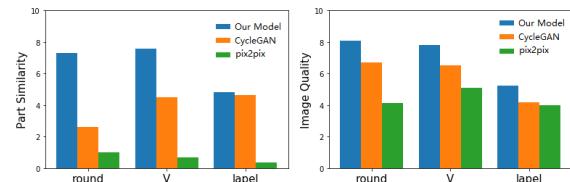
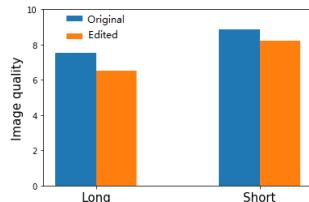


Figure 7. Average user evaluation scores based on image quality and target part similarity. The part similarity is based on users' scores on the edited part structure similarity between the edited image and the target image.

evaluations. The maximum score is ten. As showed in Fig. 7, our model receives best scores in both image quality and similarity evaluations.

We also collected users' feedback to the sleeve changing results and the feedback show that users can not distinguish real/fake between our generated images and real images. The sleeve-tests only evaluate the image quality. In each test case, there are two pictures: (1) the original picture and (2) the edited picture. Users decide scores ranging from zero to ten to both pictures based on the image quality. Translating from short to long sleeves is in general a harder task due to auto-texture filling and background changing. Users may find that it is harder to distinguish the real images from the edited ones for short sleeve garments. Those observations are reflected in the evaluation scores.

To prove that TailorGAN can be useful in image item retrieves, we upload the edited images to a searching-based

756
757
758
759
760
761
762763 Figure 8. User feedback to sleeve editing results based on image
764 quality.765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781782 Figure 9. Top 5 matching items from [1]. The first column (green
783 box) are generated images and the rest of the column (yellow) are
784 retrieved images based on the generated images from [1].785 website[1] and show sample search results in Fig.9.
786

5.6. Ablation Studies

790
791
792
793
794
795
796
797
798
799
800

In this section, we want to clarify our choices of using edge map inputs and picking VGG perceptual loss. As we argued, to disentangle the texture and the design structure is crucial in making high quality fashion attribute editing. Also, using VGG perceptual loss makes sharper results. A popular loss function for image synthesis is the L1 error between a synthesized and target image. Our experiments show that using L1 loss produces blurry images. For a qualitative analysis, We plot testing results of RGB inputs and results of using L1 loss and compare those changes with our full model.

801
802
803
804
805
806
807
808
809

Fig.10 shows edited image results of our current model versus L1 loss and RGB input results. L1 loss model doesn't prioritize high frequency details and tends to average the pixel values in the editing region. On the other hand, VGG layers captures various features from edge, color features in the starting layers to the texture and common image structures in the high level layers. On the other hand, using RGB images as inputs, the results show unexpected effects. The RGB-input results include extra structures that

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863Figure 10. Results of using L1 loss ($w/o \mathcal{L}_{VGG}$) and using RGB pictures ($w/o E^T$) as inputs instead of edge maps. We also present results of Pix2Pix in the last column.

	C.E.	SSIM	PSNR
Full model	6.78	0.9386	24.51
w/o perceptual loss	7.62	0.9121	22.23
w/o edge input	7.21	0.9082	23.58

836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

Table 4. Measurements of the final model versus two variants on

testing set. w/o perceptual loss represents we use pixel wise L1

loss rather than perceptual loss. w/o edge input indicates that we

use RGB image as input rather than grey-scale edge map.

do not belong to neither target images nor original images. Our hypothesis is that those unexpected structures are from the texture-design entanglement. We measure the classification error, SSIM and PSNR for three variants. Since the major part of the image is left untouched in the result, the leading score may not be impressive in numbers. But, through the qualitative analysis, we can confirm that our current model can generate high quality images with mode details preserved.

6. Conclusion

In this paper, we introduce a novel task to the deep learning fashion field. We investigate the problem of doing image editing to a fashion item with user defined attribute. Such methods can be useful in real world applications. We propose a novel training schema that can do single attribute manipulations to an arbitrary fashion image. To serve a better model training, we collect and build our own dataset. Our model out-performances the baseline models and successfully generates photo realistic images with desired attribute.

864

References

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

- [1] Taobao. www.TaoBao.com. Accessed: 2019-07-30.
- [2] Z. Al-Halah, R. Stiefelhagen, and K. Grauman. Fashion forward: Forecasting visual style in fashion. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 388–397, 2017.
- [3] H. Chen, A. C. Gallagher, and B. Girod. Describing clothing by semantic attributes. In *Computer Vision - ECCV 2012 - 12th European Conference on Computer Vision, Florence, Italy, October 7-13, 2012, Proceedings, Part III*, pages 609–623, 2012.
- [4] L. Chen, Z. Li, R. K. Maddox, Z. Duan, and C. Xu. Lip movements generation at a glance. In *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part VII*, pages 538–553, 2018.
- [5] L. Chen, R. K. Maddox, Z. Duan, and C. Xu. Hierarchical cross-modal talking face generation with dynamic pixel-wise loss. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7832–7841, 2019.
- [6] L. Chen, S. Srivastava, Z. Duan, and C. Xu. Deep cross-modal audio-visual generation. In *Proceedings of the on Thematic Workshops of ACM Multimedia 2017, Mountain View, CA, USA, October 23 - 27, 2017*, pages 349–357, 2017.
- [7] Y. Cheng, Z. Gan, Y. Li, J. Liu, and J. Gao. Sequential attention gan for interactive image editing via dialogue. *arXiv preprint arXiv:1812.08352*, 2018.
- [8] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. C. Courville, and Y. 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.
- [9] X. Huang, Y. Li, O. Poursaeed, J. Hopcroft, and S. Belongie. Stacked generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5077–5086, 2017.
- [10] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017.
- [11] N. Jetchev and U. Bergmann. The conditional analogy gan: Swapping fashion articles on people images. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2287–2292, 2017.
- [12] S. Jiang and Y. fu. Fashion style generator. *IJCAI*, 2017.
- [13] N. Kato, H. Osone, K. Oomori, C. W. Ooi, and Y. Ochiai. Gans-based clothes design: Pattern maker is all you need to design clothing. In *Proceedings of the 10th Augmented Human International Conference 2019*, page 21. ACM, 2019.
- [14] M. H. Kiapour, X. Han, S. Lazebnik, A. C. Berg, and T. L. Berg. Where to buy it: Matching street clothing photos in online shops. In *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015*, pages 3343–3351, 2015.
- [15] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

- [16] A. B. L. Larsen, S. K. Sønderby, H. Larochelle, and O. Winther. Autoencoding beyond pixels using a learned similarity metric. In *Proceedings of the 33nd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016*, pages 1558–1566, 2016.
- [17] C. Lassner, G. Pons-Moll, and P. V. Gehler. A generative model of people in clothing. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 853–862, 2017.
- [18] M. Li, H. Huang, L. Ma, W. Liu, T. Zhang, and Y. Jiang. Unsupervised image-to-image translation with stacked cycle-consistent adversarial networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 184–199, 2018.
- [19] S. Liu, Z. Song, G. Liu, C. Xu, H. Lu, and S. Yan. Street-to-shop: Cross-scenario clothing retrieval via parts alignment and auxiliary set. In *2012 IEEE Conference on Computer Vision and Pattern Recognition, Providence, RI, USA, June 16-21, 2012*, pages 3330–3337, 2012.
- [20] Y. Lu, Y. Tai, and C. Tang. Attribute-guided face generation using conditional cyclegan. In *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part XII*, pages 293–308, 2018.
- [21] L. Ma, X. Jia, Q. Sun, B. Schiele, T. Tuytelaars, and L. Van Gool. Pose guided person image generation. In *Advances in Neural Information Processing Systems*, pages 406–416, 2017.
- [22] D. Martin, C. Fowlkes, D. Tal, and J. Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *Proc. 8th Int'l Conf. Computer Vision*, volume 2, pages 416–423, July 2001.
- [23] A. Nguyen, A. Dosovitskiy, J. Yosinski, T. Brox, and J. Clune. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. In *Advances in Neural Information Processing Systems*, pages 3387–3395, 2016.
- [24] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- [25] A. Raj, P. Sangkloy, H. Chang, J. Lu, D. Ceylan, and J. Hays. Swapnet: Garment transfer in single view images. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 666–682, 2018.
- [26] S. E. Reed, Z. Akata, X. Yan, L. Logeswaran, B. Schiele, and H. Lee. Generative adversarial text to image synthesis. In *Proceedings of the 33nd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016*, pages 1060–1069, 2016.
- [27] N. Rostamzadeh, S. Hosseini, T. Boquet, W. Stokowiec, Y. Zhang, C. Jauvin, and C. Pal. Fashion-gen: The generative fashion dataset and challenge. *arXiv preprint arXiv:1806.08317*, 2018.
- [28] W. Shen and R. Liu. Learning residual images for face attribute manipulation. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 1225–1233, 2017.

- 972 [29] K. Simonyan and A. Zisserman. Very deep convolutional
973 networks for large-scale image recognition. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. 1026
974
975
976
- 977 [30] X. Wang and A. Gupta. Generative image modeling using
978 style and structure adversarial networks. In *European Conference on Computer Vision*, pages 318–335. Springer, 2016. 1027
979
980 [31] Z. Wang, X. Tang, W. Luo, and S. Gao. Face aging
981 with identity-preserved conditional generative adversarial
982 networks. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 7939–7947, 2018. 1028
983
984
985 [32] R. Y. X. Han, Z. Wu and L. S. Davis. Viton: an image based
986 virtual try-on network. *CVPR*, 2018. 1029
987
988 [33] W. Xian, P. Sangkloy, V. Agrawal, A. Raj, J. Lu, C. Fang,
989 F. Yu, and J. Hays. Texturegan: Controlling deep image synthesis
990 with texture patches. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 8456–8465, 2018. 1030
991
992 [34] T. Xiao, T. Xia, Y. Yang, C. Huang, and X. Wang. Learning
993 from massive noisy labeled data for image classification. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, pages 1031
994 2691–2699, 2015. 1032
995
996 [35] S. Xie and Z. Tu. Holistically-nested edge detection. In *Proceedings of the IEEE international conference on computer vision*, pages 1395–1403, 2015. 1033
997
998
999
- 1000 [36] Y. Yuan, S. Liu, J. Zhang, Y. Zhang, C. Dong, and L. Lin.
1001 Unsupervised image super-resolution using cycle-in-cycle
1002 generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 701–710, 2018. 1034
1003
1004 [37] B. Zhao, X. Wu, Z.-Q. Cheng, H. Liu, Z. Jie, and J. Feng.
1005 Multi-view image generation from a single-view. In *2018 ACM Multimedia Conference on Multimedia Conference*,
1006 pages 383–391. ACM, 2018. 1035
1007
1008 [38] J. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image
1009 translation using cycle-consistent adversarial networks. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 1036
1010 2242–2251, 2017. 1037
1011
1012 [39] S. Zhu, S. Fidler, R. Urtasun, D. Lin, and C. C. Loy. Be your
1013 own prada: Fashion synthesis with structural coherence. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 1038
1014 1689–1697, 2017. 1039
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025