

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

# CT-Net: Complementary Transferring Network for Garment Transfer with Arbitrary Geometric Changes

Anonymous CVPR 2021 submission

Paper ID 4041

## Abstract

Most existing virtual try-on methods attempt to transfer a stand-alone piece of clothing to an image of a person. While garment transfer shows greater potential in realistic applications with the goal of transferring outfits directly across different people images. However, transferring between images with heavy misalignments or severe occlusions still remains as a challenge. In this work, we propose Complementary Transferring Net (CT-Net) to adaptively model different levels of geometric changes and transfer outfits between different people. In specific, CT-Net consists of three modules: i) a complementary warping module first estimates two complementary warpings to transfer the desired clothes and segmentations in different granularities. ii) The coarse transferred segmentations is further refined leveraging a mask refine module, which is then used to guide the generation. iii) A fuse module adaptively integrates all information to render the garment transfer results with well-preserved characteristics of clothes and identities of human. Extensive experiments conducted on DeepFashion dataset demonstrate that our model synthesizes high quality virtual try-on images and significantly outperforms all state-of-art methods both qualitatively and quantitatively.

## 1. Introduction

Most existing virtual try-on methods are based on simplifying assumptions: (i) Pure clothing images or 3D information are available as inputs; (ii) Pose changes are simple without heavy misalignments or severe occlusions of the clothes. We argue that these simplifying assumptions greatly limit the application scope of these methods in the realistic virtual try-on scenarios. To address this issue, we propose Complementary Transferring Network (CT-Net), a novel image-based garment transfer network that does not rely on pure clothing images or 3D information while capable to adaptively deal with different levels of geometric changes. As shown in Figure 1, given a target person

image  $I^T$  and a model image  $I^M$ , without any restriction to the poses or shapes of  $I^T$  and  $I^M$ , our CT-Net synthesizes photo-realistic try-on results, in which the person in  $I^T$  wearing the clothes depicted in  $I^M$  with well-preserved details.

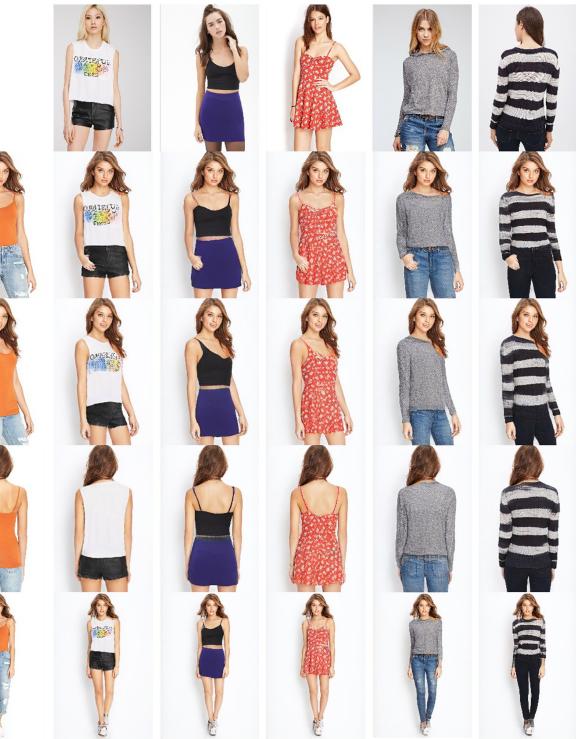


Figure 1. Virtual try-on results generated by CT-Net. First row: model images. First column: target person images. As shown above, CT-Net naturally transfers clothes across different people images with arbitrary poses or shapes and synthesizes photo-realistic try-on results with well-preserved characteristics of the desired clothes and identities of humans.

Despite various methods have been proposed to realize virtual try-on in different settings [13, 36, 41, 40, 29, 6, 11, 39], there is still a gap between these methods and the unlimited realistic scenarios. Some methods [10, 4, 28] in-

108 involve 3D information to deal with misalignments or occlusions, but they are greatly limited by expensive devices and  
 109 high computational costs. [13, 36, 41, 40] rely on stand-  
 110 alone clothing images, which are not easy to get timely on-  
 111 line. Moreover, most of them attempt to model the geo-  
 112 metric changes of the clothes utilizing a Thin Plate Spline (TPS)  
 113 warping. However, limited by a small number of parame-  
 114 ters, TPS warping is only capable to shape simple deforma-  
 115 tions, which makes their methods fail to deal with complex  
 116 cases when heavy misalignments or severe occlusions oc-  
 117 cur. Garment transfer methods aim to transfer outfits across  
 118 different people images. Although prior arts [29, 11, 39]  
 119 has achieved considerable progress, none of them address  
 120 the issue of heavy misalignments and severe occlusions.  
 121

We aim to fulfill this gap by proposing a novel garment transfer network, Complementary Transferring Network (CT-Net), which precisely transfers outfits across different people while tolerating heavy misalignments and severe occlusions. As shown in Figure 2, CT-Net has three sub-modules. First, a ComplementaryWarp Module(CWM) is introduced to warp the desired clothes and segmentation into target region. **Specially**, we estimate two complementary warpings with different levels of freedom: (i) Dense correspondence warping. (ii) Thin Plate Spline (TPS) warping. Dense correspondence warping has a high degree of freedom and is utilized to map pixels to be well-aligned with the target pose to guide the generation; **While** limited by a small number of parameters, TPS warping roughly transfers the desired clothes into target region with well-preserved details, which is utilized to further refine the synthesized result. Second, a MaskRefine Module(MRM) is introduced to predict transferred segmentation, in which the target person wearing the desired clothes. **Opposed** to prior works, which may suffer from the large misalignments between inputs [29, 11, 40], we adopt a transfer-refine strategy by first warping the original segmentation to be well-aligned with the target pose and then refine it to predict precise transferred segmentation. Our MaskRefine module not only works well in complex cases with heavy misalignments or severe occlusions, but also adds structural constraints to the learning of dense correspondence warping, encouraging the warping results to be more coherent with the source image. Third, a Fuse Module(FM) adaptively integrates all the information provided by previous modules to render the garment transfer results.

Extensive experiments conducted on DeepFashion dataset demonstrate the superiority of our method compared to the state-of-art methods. To summarize, our main contributions can be summarized as follows:

- We propose a novel image-based garment transfer network, which adaptively models different levels of geometric changes by estimating two complementary warpings and synthesizes photo-realistic garment transfer re-

- sults with well-preserved details of the clothes and human parts. 162  
 • We introduce a novel ComplementaryWarp module to es- 163  
 timate **two precise complementary warpings simultaneously**. 164  
 • A novel MaskRefine Module is proposed to predict pre- 165  
 cise **transferred segmentation**, which also adds structural 166  
 constraints to the training of warping module, forcing the 167  
 warping results to be more coherent with target pose. 168  
 • Evaluated on DeepFashion [23] dataset, CW-Net outper- 169  
 forms all the state-of-art methods by a large margin and 170  
 synthesizes photo-realistic images with well-preserved 171  
 clothing and body details. 172  
 173  
 174  
 175  
 176  
 177  
 178  
 179  
 180  
 181  
 182  
 183  
 184  
 185  
 186  
 187  
 188  
 189  
 190  
 191  
 192  
 193  
 194

## 2. Related Work

**Generative Adversarial Networks.** Generative Adversarial Networks (GANs) [9] have been demonstrated very effective in generating fake images, which are indistinguishable from the real ones in the original dataset. To further control the generation results, conditional GAN (cGAN) [26] inputs extra information to guide the generation, which promotes the development of many applications, such as image editing [5, 43, 19, 15] and image synthesis [16, 44, 37, 27, 1, 12]. Specially, Isola *et al.* [16] proposed an image-to-image translation network to transfer images from one domain to another (*e.g.* segmentation to photos, edges to photos), which shows convincing ability to handle cross-domain relationships. Similarly, we also employ a cGAN to synthesize photo-realistic virtual try-on results conditioned on desired clothing and target pose representations.

**Pose-guided Human Image Generation.** Ma *et al.* [24] made an early attempt to generate human images conditioned on pose utilizing a two-stage network. They first generate a coarse image according to the target pose and then refine it with adversarial loss. Esser *et al.* [7] proposed a conditional U-Net to disentangle the pose and appearance. However, their model still suffers from misalignments between source and target pose, result in unsatisfying generation when large spatial transformation occurs. More recent methods propose to solve this problem utilizing warp-based methods. Siarohin *et al.* [33] introduces deformable skip connections to spatially transform the features. Zhu *et al.* [45] employs a sequence of pose-attentional transfer blocks to progressively deal with large pose discrepancies. [11, 30, 21] estimate appearance flow to transfer the pixels or features to be aligned with the target view to facilitate the generation. Zhang *et al.* [42] for the first time introduces cross-domain semantic matching, which precisely learns a dense mapping between cross-domain inputs. Inspired by [42], we also estimate a dense semantic warping. However, we focus on the exact problem of garment transfer and learn

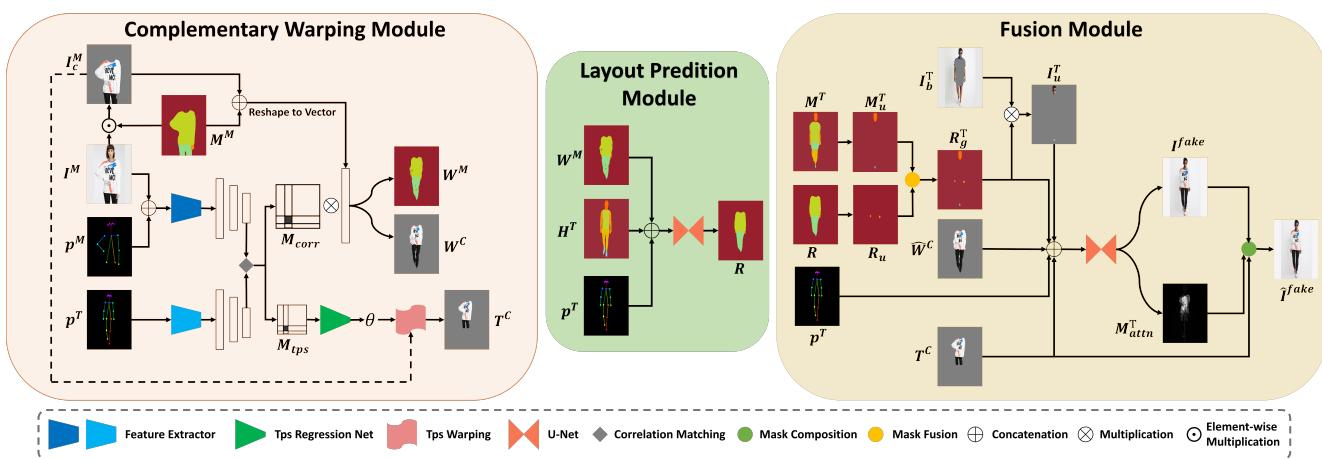
216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230

Figure 2. The overall architecture of CT-Net, which consists of three modules. (i) Complementary Warp Module estimates two complementary warpings to transfer the desired clothes in different granularities. (ii) MaskRefine Module predicts precise transferred segmentations based on coarse transferred result  $W^M$ . (iii) Fuse Module adaptively integrates all information to render photo-realistic garment transfer results.

a semantic mapping between pose representations. Benefited from the joint training of all the modules, we achieve more precise warping results on the clothing items.

**Virtual Try-on.** Virtual Try-on has been paid more and more attention as the fast development of e-commerce. Many conventional virtual try-on methods rely on 3D information, such as pre-defined 3D clothing template [10], 3D body model [4] or depth information captured by 3D sensors [28]. Along with the advances of deep neural networks, more recent works attempt to synthesize try-on results based on 2D images. Han *et al.* [13] proposed VITON to transfer clothes in a stand-alone pure clothing image to the target person. Variations [36, 41, 40] are proposed to improve VITON with clearer clothing details and well-preserved human identities. However, all of them rely on simple assumptions that the clean clothing images are available and geometric changes are simple enough to be modeled by a thin plate spline warping (TPS) with a small number of parameters (*e.g.* 6 for affine and  $2 \times 5 \times 5$  for TPS as in [36]). These assumptions greatly hinder the application of these methods in the realistic virtual try-on scenarios. Wu *et al.* [39] proposed to use densepose [2] descriptor to warp the desired clothes onto the target person. But mapping estimated by densepose descriptor can be very sparse when there are large occlusions, leading to unconvincing synthesized results. Methods mentioned above only focus on the transfer of upper clothes. SwapNet [29] employs a two-stage network to transfer the entire outfit of the inputting images. To deal with the misalignments of features, they adopt ROI pooling and encode each clothing regions into high-dimensional features before feeding them into the generator. However, the encoded features are inadequate to preserve the details of the local textures, which leads

to blurry synthesized results. Different from these methods, we adaptively combine two complementary warpings to deal with different levels of deformations and synthesize photo-realistic images with well-preserved characteristics of clothes and identities of humans.

### 3. Complementary Transferring Network

Given the image  $I^M$  depicting model wearing desired clothes, the image  $I^P$  depicting target person, assuming clothes, poses and shapes of  $I^M$  and  $I^P$  can be arbitrary, our goal is to synthesize high quality try-on results with well-preserved textures. To achieve our goal, we adopt a warp-refine strategy, presenting Complementary Transferring Network (CT-Net). As shown in Figure 2, CT-Net consists of three modules. First, we introduce a Complementary Warp Module (CWM), in which a dense correspondence warping and a Thin Plate Spline warping are estimated. Dense correspondence warping is utilized to warp the clothes and segmentations of  $I^M$  to be aligned with target poses to guide the generation, and the synthesized result will be further refined leveraging the complementary warping result from TPS (Section 3.1). Second, we introduce a MaskRefine Module (MRM), in which we refine the transferred mask with extra clothing-agnostic pose representations (Section 3.2). The third module is a Fuse Module (FM), which adaptively integrates all the information to render photo-realistic garment transfer results (Section 3.3).

#### 3.1. Complementary Warp Module

To synthesize garment transfer results, one of the main challenges is to combine the clothes of the model with the target person's pose. Most of the prior works feed the target

270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323

324 pose and model clothes into conditional generative adversarial netwrok (cGAN). However, in cases with heavy misalignments or severe occlusions, they can only retain the global style and fail to preserve local details. A good practice to overcome this issue is to utilize warping methods to align the clothes with the target pose first before feeding them into generation. Similarly, we propose ComplementaryWarp Module, which simultaneously estimates two complementary warpings, to transfer the information in two different granularities to facilitate the generation.

325 As shown in Figure 2, we first employ two seperate feature extractors to extract high-level features. Then we match the features to calculate correspondence matrixs,  $\mathcal{M}_{corr}$  and  $\mathcal{M}_{tps}$ , which are used to estimate dense correspondence warping and TPS warping. Given a model image  $I^M$ , the corresponding clothes  $I_c^M$  and the segmentation  $M^M$ , dense correspondence warping is utilized to transfer  $I_c^M$  and  $M^M$  in finer granularity to get  $W^C$  and  $W^M$ , which are well-aligned to target pose. TPS warping transfers  $I_c^M$  in a coarse granularity to get  $T^C$ , which is roughly aligned with target region.

326 **Correspondence Matrix.** We adopt the state-of-art pose estimation method [3] to estimate pose segmentations of the target person and the model image. And we further convert the sparse pose segmentations into dense by replacing each zero pixel with its distance to the joints, denoted as  $p^T, p^M$ . 327 Different from [42], which aims to learn a correspondence matrix from multi-modal inputs, we simplify the task by directly inputing pose representations in both sides. To be specific, let  $\mathcal{F}_A, \mathcal{F}_B$  denote the seperate feature extractors, we first extract high-level features  $m_f \in \mathbb{R}^{H \times W \times C}$  and  $t_f \in \mathbb{R}^{H \times W \times C}$  as follows:

$$m_f = \mathcal{F}_A(I^M, P^M) \quad (1)$$

$$t_f = \mathcal{F}_B(P^T) \quad (2)$$

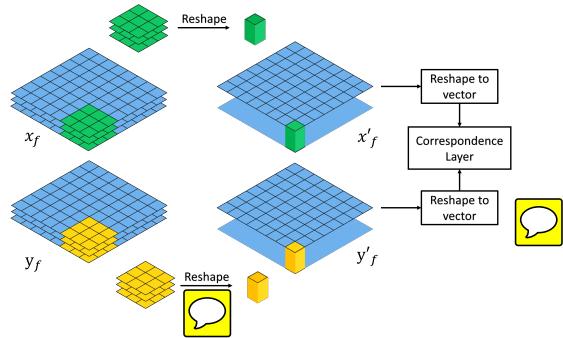
328 To estimate the correspondence matrixs, we first aggregate the features into different scales with different kernels, which is illustrated in Figure 3. Specifically, we use kernel of size 3, wiht stride 1 and padding size 1 to estimate the correspondence matrix  $\mathcal{M}_{corr} \in \mathbb{R}^{HW \times HW}$  for 329 dense correspondence warping and we select kernel of size 4, with stride 4 and padding size 0 to estimate  $\mathcal{M}_{tps} \in \mathbb{R}^{HW/16 \times HW/16}$  for TPS warping.

330 We employ the same correspondence layer as [42] to 331 match the aggregated features  $m'_f$  and  $t'_f$ , which can be 332 formulated as:

$$\mathcal{M}(i, j) = \frac{(m'_f(i)^T - u_m)(t'_f(j) - u_t)}{\|m'_f(i) - u_m\| \|t'_f(j) - u_t\|} \quad (3)$$

333 where  $u_m$  and  $u_t$  represent mean vectors.

334 **Dense Correspondence Warping.** We calculate the weighted average to estimate the dense correspondence



335 Figure 3. Illustration of the correspondence matching process with 336 different kernels. This figure shows an example of kernel with size 3. 337

338 warping [42]:

$$\mathcal{W}^X(u) = \sum_v softmax(\alpha \mathcal{M}_{corr}(u, v)) \cdot Y^X(v) \quad (4)$$

339 where  $\alpha$  is a hyper-parameter controling the sharpness of 340 the softmax. We set it as 100 here. Dense correspondence 341 warping learns a dense mapping between two images with 342 a high degree of freedom, which is capable to handle large 343 geometric changes. We utilize it to transfer the clothes and 344 segmentions of the model image to be well-aligned with the 345 target person's pose and shape, which provides important 346 guides for the generator to synthesize high-quality try-on 347 results. However, learning a perfect dense semantic matching 348 is hard, some important textures may get lost in the warping. 349 To deal with this issue, we involve TPS warping, which 350 has very limited parameters but is capable to transfer with 351 little information loss, to refine the final result.

352 **TPS Warping.** We estimate the TPS warping from 353  $\mathcal{M}_{tps}$ . As shown in Figure 2, we first employ a regression 354 net [6, 40] to predict the corresponding control points and 355 then calculate the parameters  $\theta$ .

356 Given the model clothes  $I_c^M$ , we transfer it with deformation 357 shaped by TPS to be roughly aligned with the target 358 person  $I^T$ . For training, we adpot the second-order 359 constraint from Yang *et al.* [40] to restrict the TPS warping 360 from generating unnatural deformations or mess textures, 361 which can be formulated as:

$$\begin{aligned} \mathcal{L}_2 = & \sum_{p \in P} \lambda_r (\|pp_0\|_2 + \|pp_1\|_2) + (\|pp_2\|_2 - \|pp_3\|_2) \\ & + \lambda_s (|S(p, p_0) - S(p, p_1)| + |S(p, p_2) - S(p, p_3)|) \end{aligned} \quad (5)$$

362 Where  $\lambda_r$  and  $\lambda_s$  are hyper-parameters;  $S(p, p_i) = \frac{y_i - y}{x_i - x}$  ( $i = 0, 1, 2, 3$ ) is the slope between two points. The 363 total loss can be represented as  $\mathcal{L}_{tps}$ :

$$\mathcal{L}_{tps} = \lambda_1 \mathcal{L}_1 + \lambda_2 \mathcal{L}_2 \quad (6)$$

432 where  $\mathcal{L}_1 = \|I_c^T - T^C\|_1$ ;  $\lambda_1$  and  $\lambda_2$  are the weights for  
 433 two loss terms. Both of them are set to 10, respectively, in  
 434 our experiments.  
 435

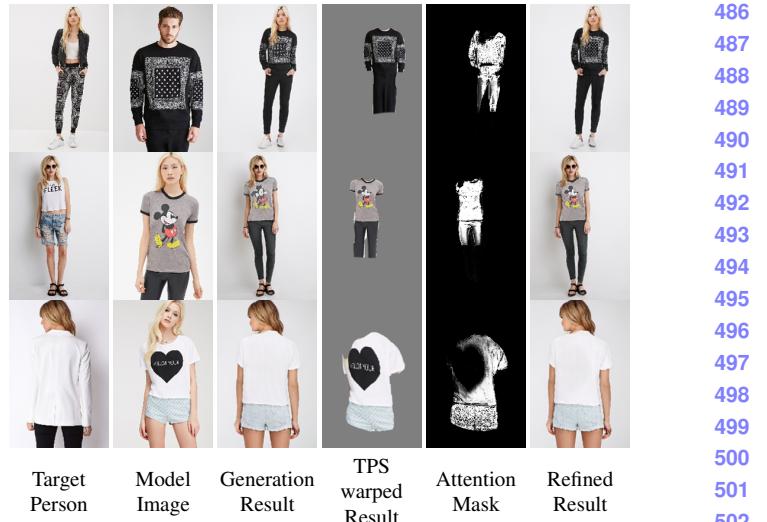
### 3.2. MaskRefine Module

436 We propose MaskRefine Module to generate a precise  
 437 transferred segmentation, in which the target person wearing  
 438 the desired clothes depicted in the model image. Prior  
 439 works mostly feed the pose representations of the target person  
 440 and the model clothes into generative adversarial networks  
 441 (GANs) directly to generate the transferred segmentations.  
 442 However, these methods greatly suffer from the misalignments  
 443 between the target pose and desired clothes thus fail to predict  
 444 convincing results while facing with large geometric changes.  
 445 Opposed to these methods, we explore a new warp-refine strategy.  
 446 As shown in Figure 2, we first warp the original segmentation in CWM  
 447 to make it well-aligned with the target person’s pose, then we utilize a U-net  
 448 structure to refine it with extra clothing-agnostic representations.  
 449 In this way, we eliminate the misalignments between two inputs and make it much easier for our network to  
 450 predict precise transferred layouts.  
 451

452 We get the clothing segmentations utilizing the state-of-  
 453 art human parsing network [22] and apply the segmentation  
 454 of densepose descriptor [2] as our extra clothing-agnostic  
 455 representations  $H^T$ . For training this module, we adopt the  
 456 pixel-level cross-entropy loss, denoted as  $\mathcal{L}_{mask}$ . Since the  
 457 head and shoes are not in the transfer list of our task, we  
 458 remove these areas in the segmentation before feeding it  
 459 into our network. Through joint training, MaskRefine module  
 460 also adds structural constraints to the training of CWM,  
 461 forcing warping results to be more coherent with target pose  
 462 and have clearer boundaries.  
 463

### 3.3. Fuse Module

464 Fuse Module is proposed to adaptively integrate all the  
 465 information provided by previous stages to render garment  
 466 transfer results. As shown in Figure 2, we first fuse the  
 467 head and shoes masks in the original segmentation with the  
 468 refined arm and leg masks to form the guide segmentation  
 469  $R_g^T$ . Leveraging  $R_g^T$ , we extract the non-target body parts  
 470 from the masked body region  $I_u^P$ . We then adopt conditional  
 471 generative network (cGAN) to integrate the untarget  
 472 body parts  $I_u^P$ , the guided mask  $R_g^T$ , warping results  $\widehat{W}^C$   
 473 and  $T^C$  and the pose segmentations of the target person  $p^T$   
 474 to render the try-on result  $I^{fake}$ . Note  $\widehat{W}^C = W^C \odot R_c$ ,  
 475 where  $R_c$  denotes the refined clothing mask. By feeding  
 476 the non-target body parts, the generator learns to preserve  
 477 all the details in the non-target region and synthesize  
 478 results with distinct human identities [40]. To further refine  
 479 the generation result, we simultaneously predict a compo-  
 480 sition mask  $M_{attn}^T$  utilizing the generator. The generated  
 481 result  $I^{fake}$  is composited with the warping result of TPS  
 482



503 Figure 4. Examples of our attention mechanism. From left to right:  
 504 target Person  $I^T$ , Model Image  $I^M$ , Generation Result  $I^{fake}$ , TPS  
 505 warped Result  $T^C$ , Attention Mask  $M_{attn}$ , Refined Result  $\hat{I}^{fake}$ .  
 506

507 to synthesize the final garment transfer result  $\hat{I}^{fake}$ .  
 508

$$\hat{I}^{fake} = T^C \odot M_{attn}^T + I^{fake} \odot (1 - M_{attn}^T) \quad (7)$$

509 Leveraging the attention mask, we adaptively deal with  
 510 different cases by combining the advantages of the two  
 511 warping methods. For example, shown by the third row  
 512 in Figure 4, when there are heavy misalignments or severe  
 513 occlusions, the final generation result may rely more on the  
 514 warping result from dense correspondence warping, ignor-  
 515 ing the large logos or unreasonable textures; When the geo-  
 516 metric change is simple and can be shaped by the TPS war-  
 517 ping, as the first two rows in Figure 4, the attention mask may  
 518 select more regions on the warping result from TPS to re-  
 519 fine the generation results. In this way, we adaptively deal  
 520 with different levels of geometric changes by combing the  
 521 complementary warping results and expand the application  
 522 scope of our model to wilder scenarios.  
 523

### 3.4. Loss Functions

524 To encourage the training of different modules benefit  
 525 each other, we train our model in a joint style. We com-  
 526 bine several different losses to produce high-quality try-on  
 527 results, which will be introduced in details in the following  
 528 sections.  
 529

530 **Perceptual Loss.** Based on differences between high-  
 531 level features, perceptual loss has been proved efficiently in  
 532 the image generation tasks [18]. To pose perceptual con-  
 533 straints on the synthesized results, we adopt a pre-trained  
 534 VGG network [34] to extract multi-level features  $\phi_j$  to com-

540 put the perceptual loss as:

$$\mathcal{L}_{perceptual}(G) = \sum_{j=1}^N \lambda_j \|\phi_j(\hat{x}_t) - \phi_j(x)\|_2 \quad (8)$$

541 **Style Loss.** We further adopt the style loss [8] to penalize  
 542 the statistic error between high-level features, which can be  
 543 formulated as:

$$\mathcal{L}_{style} = \sum_{j=1}^N \|G_j^\phi(\hat{x}_t) - G_j^\phi(x_t)\|_2 \quad (9)$$

544 where  $G_j^\phi$  denotes the Gram matrix estimated from  $\phi_j$ .

545 **Contextual Loss.** To encourage our network to preserve  
 546 more details from the model clothes  $I_c^M$ , we employ the  
 547 contextual loss proposed in [25], which can be formulated  
 548 as:

$$\begin{aligned} \mathcal{L}_{contextual} &= \\ &\sum_{l=1} \lambda_l \left[ -\log \left( \frac{1}{n_l} \sum_i \max_j A^l(\phi_i^l(\hat{x}_B), \phi_j^l(y_B)) \right) \right], \end{aligned} \quad (10)$$

549 where  $A^l$  denotes the pairwise affinities between features.

550 **Adversarial Loss.** To force the generator to learn the  
 551 real distributions of the dataset and generate realistic human  
 552 images, we deploy a discriminator given in [37] to discrim-  
 553 inate the generated fake images from the real samples in the  
 554 dataset. The loss can be formulated as:

$$\begin{aligned} \mathcal{L}_{adv}^D &= \mathbb{E}_{x,y}[\log(\mathcal{D}(x,y))] + \mathbb{E}_x[\log(1 - \mathcal{D}(y, \mathcal{G}(x)))] \\ \mathcal{L}_{adv}^G &= -\mathbb{E}[\mathcal{D}(\mathcal{G}(x,y))] \end{aligned} \quad (11)$$

555 where  $x$  represents the input and  $y$  is the ground-truth.

556 **Objective Function.** Besides losses above, we apply a L1  
 557 regularization  $\mathcal{L}_{reg} = \|1 - M\|_1$  on  $M_{attn}$  to encourage the  
 558 generator to learn more from the TPS warped clothes. We  
 559 also take L1 loss to stabilize our training process, which can  
 560 be defined as  $\mathcal{L}_{L1}(G) = \|\hat{x} - x\|_1$ . Our objective function  
 561 is a weighted sum of above terms:

$$\begin{aligned} \mathcal{L}(G) &= \alpha_1 \mathcal{L}_1 + \alpha_2 \mathcal{L}_{tps} + \alpha_3 \mathcal{L}_{mask} + \alpha_4 \mathcal{L}_{perceptual} + \\ &\alpha_5 \mathcal{L}_{style} + \alpha_6 \mathcal{L}_{contextual} + \alpha_7 \mathcal{L}_{adv} + \alpha_8 \mathcal{L}_{reg} \end{aligned} \quad (12)$$

574 where  $\alpha_i, (i = 1, \dots, 8)$  are hyper-parameters controlling  
 575 the weights of each loss.

## 4. Experiments

### 4.1. Dataset

586 We evaluate our model on the In-shop Clothes Retrieval  
 587 Benchmark of DeepFashion dataset [23], which contains

594 52,712 fashion images of resolution  $256 \times 256$ . For train-  
 595 ing, we select 37,836 pairs of images depicting the same  
 596 person wearing the same outfit with different poses. At test  
 597 stage, we select 4,932 pairs of images which are not over-  
 598 laped with the training set. As the realistic virtual try-on  
 599 scenarios, each testing pair contains two different people  
 600 with different clothes and poses.

### 4.2. Implementation Details

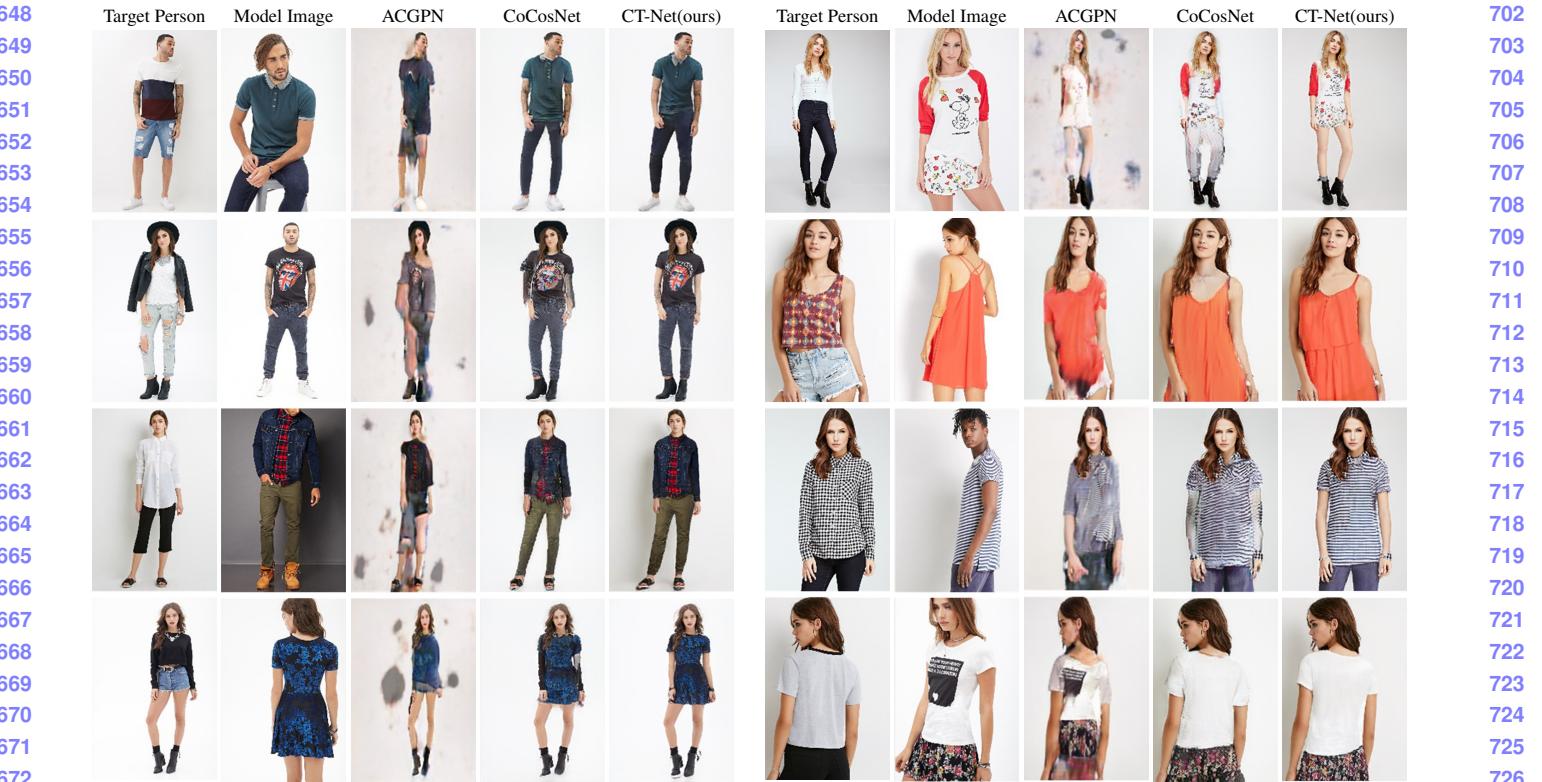
601 **Architecture.** Let  $C_k$  denote a Convolution layer with  
 602 kernel size of 4, a stride 2, and  $k$  filters, followed by  
 603 InstanceNorm2d Normalization [35] and ReLu activation  
 604 function. Let  $R_k$  denotes the a Convolution layer with  
 605 kernel size of 3, a stride 1, and  $k$  filters. Let  $L_k$  de-  
 606 notes a Linear function output  $k$  dimension. Let  $ResBlock$   
 607 denotes the original Residual Block proposed in [14], in  
 608 which the BatchNorm2d Normalization is replaced by In-  
 609 stanceNorm2d Normalizaiton and the filters of the convolu-  
 610 tion layers are 256. The two separate feature extractors  $\mathcal{F}$   
 611 in CWM share the same structure: R64, C128, R256, C256,  
 612 ResBlock \* 6. We adopt the same regression net as [6] for  
 613 TPS warping, which consists of C512, C256, C128, C64,  
 614 L32. Structure of generators in MRM and FM are the same  
 615 as U-Net [31] except the normalization is replaced by In-  
 616 stanceNorm2d Normalizaiton. All the discriminators are  
 617 from pix2pixHD [37].

618 **Implementation.** We adopt Adam [20] with  $\beta_1 =$   
 619 0.5,  $\beta_2 = 0.999$  as the optimizer in our all experiments.  
 620 Our model is jointly trained for 100 epoches. Learning rate  
 621 is fixed at 0.0002 for the first 50 epoches and then decays to  
 622 zero linearly in the remaining steps. To balance the scales  
 623 of losses in Eqn. 12, we set  $\alpha_{1,2,3,7,8} = 10$  and  $\alpha_{4,5,6} = 1$ .

### 4.3. Baselines

624 **ACGPN.** ACGPN is a state-of-art virtual try-on network  
 625 proposed by Yang *et al.* [40], which aims to transfer a stand-  
 626 alone clothing image onto a reference person. In compari-  
 627 on to previous methods [13, 36], ACGPN first predicts  
 628 the transferred clothing segmentation in a two-stage strat-  
 629 egy, then estimates the TPS warping by feeding the trans-  
 630 ferred and original clothing segmentation into a STN [17].  
 631 ACGPN shows state-of-art performance on VITON [13]  
 632 dataset with naturally transferred clothes and well-preserved  
 633 untarget body parts.

634 **CoCosNet.** CoCosNet stands for the cross-domain cor-  
 635 respondence network proposed by Zhang *et al.* [42], aiming  
 636 to synthesize realistic image given the exemplar im-  
 637 ages. Different from other methods, CoCosNet first align  
 638 images in different domain by establishing dense corres-  
 639 pondence matching and then use the aligned exemplar image  
 640 as guide to synthesize photo-realistic images, which shows  
 641 state-of-art performance in pose-guided image generation  
 642 task. However, CoCosnet lacks of the ability to generate



In Figure 6, we visualize warping results from different methods to make further comparisons. Benefited from the joint training of all modules, CT-Net shows superior performance in estimating both the dense correspondence warping and the Thin Plate Spline warping more precisely.

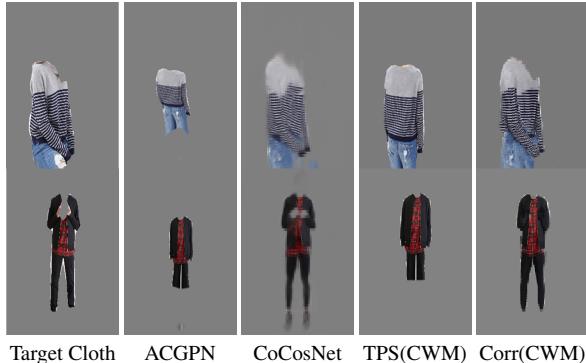


Figure 6. Visual comparisons of warping results. Corr(CWM) represents the warping results of dense correspondence warping estimated in ComplementaryWarp Module.

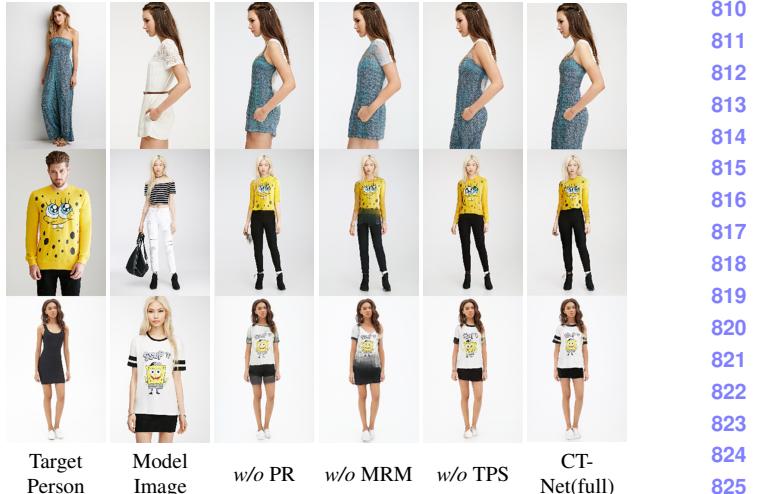


Figure 7. Visual comparisons of our ablation methods.

#### 4.6. Ablation Study

We conduct ablation experiments to explore the effectiveness of the main components in our model. In particular, *w/o PR* denotes removing the pose representation  $P^M$  inputting to the ComplementaryWarp Module. *w/o MRM* denotes removing MaskRefine Module. *w/o TPS* denotes removing the estimation of TPS warping.

Table 1 reports all the results of our ablation experiments. In specific, our full model outperforms *w/o PR* and *w/o MRM* by a margin in Warp-SSIM, which indicates that feeding  $p^M$  in the other side and adding structural constraints utilizing MaskRefine Module significantly facilitate the estimation of dense correspondence warping. We also find that our full model and *w/o TPS* have similar scores in all metrics, since SSIM only roughly measures the local similarity and IS only captures the realism of the images.

To further demonstrate the superiority of our full model, we visualize some examples to make qualitative comparisons in Figure 7. Since *w/o PR* and *w/o MRM* can not estimate the warping precisely, artifacts such as incorrect clothing shape (first row) and blurry boundaries (third row) can be observed. Although *w/o TPS* achieves the best scores in terms of H-SSIM and IS, visual results show that our full model synthesizes more photo-realistic images with clearer and more complete clothing patterns.

## 5. Conclusion



We propose Complementary Transferring Net (CTNet) for garment transfer with arbitrary geometric changes. In particular, our model addresses the issue of transferring outfits with high misalignments or severe occlusions, while preserving the characteristics of clothes and the identity of humans well. We introduce three novel modules: i) We

864 propose a Complementary Warping Module (CWM), which  
 865 estimates two complementary warpings to transfer the de-  
 866 sired clothes and segmentations in different granularities.  
 867 ii) A Mask Refine Module (MRM) is employed to refine the  
 868 coarse transferred segmentations. iii) A Fuse module is pro-  
 869 posed to adaptively integrate all the information to synthe-  
 870 size photo-realistic garment transfer results. Experiment re-  
 871 sults demonstrate that our model significantly outperforms  
 872 existing state-of-art methods both qualitatively and quanti-  
 873 tatively.

## References

- 875
- 876
- 877 [1] Badour AlBahar and Jia-Bin Huang. Guided image-to-image  
 878 translation with bi-directional feature transformation. In *Proceedings of the IEEE International Conference on Computer  
 879 Vision*, pages 9016–9025, 2019. 2
- 880
- 881 [2] Riza Alp Güler, Natalia Neverova, and Iasonas Kokkinos.  
 882 Densepose: Dense human pose estimation in the wild. In *Proceedings of the IEEE Conference on Computer Vision  
 883 and Pattern Recognition*, pages 7297–7306, 2018. 3, 5
- 884
- 885 [3] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh.  
 886 Realtime multi-person 2d pose estimation using part affinity  
 887 fields. In *Proceedings of the IEEE conference on computer  
 888 vision and pattern recognition*, pages 7291–7299, 2017. 4
- 889
- 890 [4] Wenzheng Chen, Huan Wang, Yangyan Li, Hao Su, Zhen-  
 891hua Wang, Changhe Tu, Dani Lischinski, Daniel Cohen-  
 892 Or, and Baoquan Chen. Synthesizing training images for  
 893 boosting human 3d pose estimation. In *2016 Fourth Interna-  
 894 tional Conference on 3D Vision (3DV)*, pages 479–488.  
 895 IEEE, 2016. 1, 3
- 896
- 897 [5] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha,  
 898 Sunghun Kim, and Jaegul Choo. Stargan: Unified gener-  
 899 ative adversarial networks for multi-domain image-to-image  
 900 translation. In *Proceedings of the IEEE conference on  
 901 computer vision and pattern recognition*, pages 8789–8797,  
 902 2018. 2
- 903
- 904 [6] Haoye Dong, Xiaodan Liang, Xiaohui Shen, Bowen Wu,  
 905 Bing-Cheng Chen, and Jian Yin. Fw-gan: Flow-navigated  
 906 warping gan for video virtual try-on. In *Proceedings of the  
 907 IEEE International Conference on Computer Vision*, pages  
 908 1161–1170, 2019. 1, 4, 6
- 909
- 910 [7] Patrick Esser, Ekaterina Sutter, and Björn Ommer. A vari-  
 911 ational u-net for conditional appearance and shape generation.  
 912 In *Proceedings of the IEEE Conference on Computer Vision  
 913 and Pattern Recognition*, pages 8857–8866, 2018. 2
- 914
- 915 [8] Leon Gatys, Alexander S Ecker, and Matthias Bethge. Texture  
 916 synthesis using convolutional neural networks. In *Advances  
 917 in neural information processing systems*, pages 262–270, 2015. 6
- 918
- 919 [9] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing  
 920 Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and  
 921 Yoshua Bengio. Generative adversarial nets. In *Advances  
 922 in neural information processing systems*, pages 2672–2680,  
 923 2014. 2
- 924
- 925 [10] Peng Guan, Loretta Reiss, David A Hirshberg, Alexander  
 926 Weiss, and Michael J Black. Drape: Dressing any person.  
 927
- 928
- 929
- 930
- 931
- 932
- 933
- 934
- 935
- 936
- 937
- 938
- 939
- 940
- 941
- 942
- 943
- 944
- 945
- 946
- 947
- 948
- 949
- 950
- 951
- 952
- 953
- 954
- 955
- 956
- 957
- 958
- 959
- 960
- 961
- 962
- 963
- 964
- 965
- 966
- 967
- 968
- 969
- 970
- 971

- 972       tion. In *Advances in neural information processing systems*,  
973       pages 406–416, 2017. 2
- 974       [25] Roey Mechrez, Itamar Talmi, and Lihi Zelnik-Manor. The  
975       contextual loss for image transformation with non-aligned  
976       data. In *Proceedings of the European Conference on Com-*  
977       *puter Vision (ECCV)*, pages 768–783, 2018. 6
- 978       [26] Mehdi Mirza and Simon Osindero. Conditional generative  
979       adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014. 2
- 980       [27] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan  
981       Zhu. Semantic image synthesis with spatially-adaptive nor-  
982       malization. In *Proceedings of the IEEE Conference on Com-*  
983       *puter Vision and Pattern Recognition*, pages 2337–2346,  
984       2019. 2
- 985       [28] Gerard Pons-Moll, Sergi Pujades, Sonny Hu, and Michael J  
986       Black. Clothcap: Seamless 4d clothing capture and retar-  
987       getting. *ACM Transactions on Graphics (TOG)*, 36(4):1–15,  
988       2017. 1, 3
- 989       [29] Amit Raj, Patsorn Sangkloy, Huiwen Chang, James Hays,  
990       Duygu Ceylan, and Jingwan Lu. Swapnet: Image based gar-  
991       ment transfer. In *European Conference on Computer Vision*,  
992       pages 679–695. Springer, 2018. 1, 2, 3
- 993       [30] Yurui Ren, Xiaoming Yu, Junming Chen, Thomas H Li, and  
994       Ge Li. Deep image spatial transformation for person image  
995       generation. In *Proceedings of the IEEE/CVF Conference on*  
996       *Computer Vision and Pattern Recognition*, pages 7690–  
997       7699, 2020. 2
- 998       [31] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-  
999       net: Convolutional networks for biomedical image segmen-  
1000       tation. In *International Conference on Medical image com-*  
1001       *puting and computer-assisted intervention*, pages 234–241.  
1002       Springer, 2015. 6
- 1003       [32] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki  
1004       Cheung, Alec Radford, and Xi Chen. Improved techniques  
1005       for training gans. In *Advances in neural information pro-*  
1006       *cessing systems*, pages 2234–2242, 2016. 8
- 1007       [33] Aliaksandr Siarohin, Enver Sangineto, Stéphane Lathuilière,  
1008       and Nicu Sebe. Deformable gans for pose-based human  
1009       image generation. In *Proceedings of the IEEE Conference on*  
1010       *Computer Vision and Pattern Recognition*, pages 3408–  
1011       3416, 2018. 2, 8
- 1012       [34] Karen Simonyan and Andrew Zisserman. Very deep convo-  
1013       lutional networks for large-scale image recognition. *arXiv*  
1014       *preprint arXiv:1409.1556*, 2014. 5
- 1015       [35] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. In-  
1016       stance normalization: The missing ingredient for fast stylization.  
1017       *arXiv preprint arXiv:1607.08022*, 2016. 6
- 1018       [36] Bochao Wang, Huabin Zheng, Xiaodan Liang, Yimin  
1019       Chen, Liang Lin, and Meng Yang. Toward characteristic-  
1020       preserving image-based virtual try-on network. In *Pro-*  
1021       *ceedings of the European Conference on Computer Vision*  
1022       *(ECCV)*, pages 589–604, 2018. 1, 2, 3, 6
- 1023       [37] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao,  
1024       Jan Kautz, and Bryan Catanzaro. High-resolution image syn-  
1025       thesis and semantic manipulation with conditional gans. In  
1026       *Proceedings of the IEEE conference on computer vision and*  
1027       *pattern recognition*, pages 8798–8807, 2018. 2, 6
- 1028       [38] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Si-  
1029       moncelli. Image quality assessment: from error visibility to  
1030       structural similarity. *IEEE transactions on image processing*,  
1031       13(4):600–612, 2004. 8
- 1032       [39] Zhonghua Wu, Guosheng Lin, Qingyi Tao, and Jianfei Cai.  
1033       M2e-try on net: Fashion from model to everyone. In *Pro-*  
1034       *ceedings of the 27th ACM International Conference on Mul-*  
1035       *timedia*, pages 293–301, 2019. 1, 2, 3
- 1036       [40] Han Yang, Ruimao Zhang, Xiaobao Guo, Wei Liu, Wang-  
1037       meng Zuo, and Ping Luo. Towards photo-realistic virtual  
1038       try-on by adaptively generating-preserving image content. In  
1039       *Proceedings of the IEEE/CVF Conference on Computer Vi-*  
1040       *sion and Pattern Recognition*, pages 7850–7859, 2020. 1, 2,  
1041       3, 4, 5, 6, 7, 8
- 1042       [41] Ruiyun Yu, Xiaoqi Wang, and Xiaohui Xie. Vtnfp: An  
1043       image-based virtual try-on network with body and clothing  
1044       feature preservation. In *Proceedings of the IEEE Interna-*  
1045       *tional Conference on Computer Vision*, pages 10511–10520,  
1046       2019. 1, 2, 3
- 1047       [42] Pan Zhang, Bo Zhang, Dong Chen, Lu Yuan, and Fang Wen.  
1048       Cross-domain correspondence learning for exemplar-based  
1049       image translation. In *Proceedings of the IEEE/CVF Con-*  
1050       *ference on Computer Vision and Pattern Recognition*, pages  
1051       5143–5153, 2020. 2, 4, 6, 7, 8
- 1052       [43] Jiapeng Zhu, Yujun Shen, Deli Zhao, and Bolei Zhou.  
1053       In-domain gan inversion for real image editing. *arXiv preprint*  
1054       *arXiv:2004.00049*, 2020. 2
- 1055       [44] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A  
1056       Efros. Unpaired image-to-image translation using cycle-  
1057       consistent adversarial networks. In *Proceedings of the IEEE*  
1058       *international conference on computer vision*, pages 2223–  
1059       2232, 2017. 2
- 1060       [45] Zhen Zhu, Tengteng Huang, Baoguang Shi, Miao Yu, Bofei  
1061       Wang, and Xiang Bai. Progressive pose attention transfer for  
1062       person image generation. In *Proceedings of the IEEE Con-*  
1063       *ference on Computer Vision and Pattern Recognition*, pages  
1064       2347–2356, 2019. 2