

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015

Abstract

Image-based virtual try-on aims to transfer a garment image onto a person’s image naturally. However, due to the lack of structure information, existing methods often fail to fit a garment seamlessly onto the human body while preserving its characteristics. To address this issue, we adopt new structural information, i.e. clothes parsing and DensePose, and further propose a novel framework powered by Structural Information and Cascade Warping (SICW) with three components: 1) A target parsing prediction network guided by rich structural information; 2) A cascade clothes warping module built on top of TPS transformation and pixel flow; 3) A decomposed neural rendering module dealing with skin inpainting and shading refinement separately. We further construct a new dataset RAV with Rich Annotations for Image-based Virtual try-on, which will be released to facilitate future research. Extensive experiments demonstrate the new structural information’s effectiveness and our model’s superiority over state-of-the-art methods.

034
035
036
037

1. Introduction

Virtual try-on [12, 34, 14, 40] has attracted lots of interest due to its massive potential in e-commerce and video production. Given a reference person image and a target clothes image, virtual try-on aims at generating a new image for that person wearing the given clothes. In this way, customers could see themselves in their desired garments, which would significantly improve the online shopping experience. Traditional methods achieve this goal via graphics approaches and generate try-on images by manipulating 3D models. Nevertheless, it is expensive and impractical to craft decent models for all clothes with diverse textures and designs on sale. To this end, image-based virtual try-on has attracted lots of attention because of the availability of abundant clothing and try-on images on the Internet.

However, image-based virtual try-on remains a challenging task due to the gap between 2D information it acquires

054
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

Powering Image-based Virtual Try-on with Structural Information and Cascade Warping

Anonymous CVPR 2021 submission

Paper ID 3405

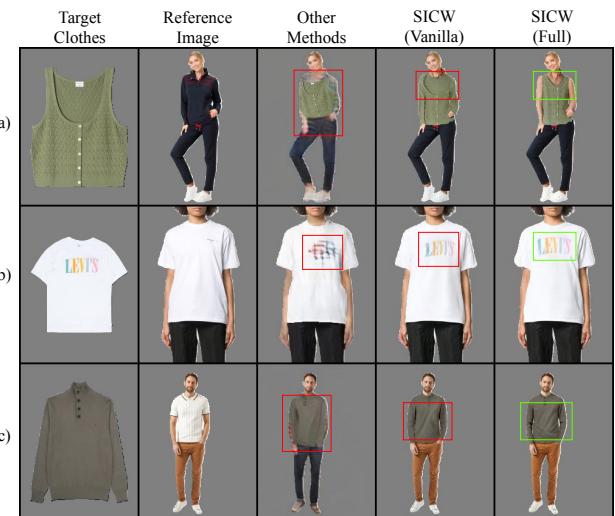


Figure 1. Representative comparisons of our model and other methods (a & c: CP-VTON [37], b: ACGPN [40]). SICW (Vanilla) indicates our method without using DensePose and clothes parsing, and SICW (full) is SICW with full structural information.

and non-rigid spatial manipulations it tries to simulate for clothes. As shown in Fig. 1, existing methods [37, 40] often fail to either (a) fit the clothing onto the human body, (b) preserve the garment’s characteristics, or (c) generate natural shading. We address these issues by introducing rich structural priors and decomposing the problem into several sub-tasks to free up network capacity.

Specifically, we improve the clothing-agnostic person representation by introducing DensePose [30], which provides semantic 3D body information while eliminates the effects of clothes. We further integrate the target clothing image and its clothes parsing into a clothing representation. Moreover, we propose a new framework powered by Structural Information and Cascade Warping, namely SICW. As shown in Fig. 2, it consists of three modules. Firstly, a target parsing prediction (TPP) network predicts the after-trying-on human parsing to guide subsequent

108 stages. Secondly, a cascade clothes warping (CCW) module warps the input clothing to fit the predicted shape. We
109 notice that TPS transformation [1] is good at general shape
110 matching but fails to conduct partial movements (e.g., lifting
111 arms), while appearance flow [42] handles minor move-
112 ment well but often causes distortion when dealing with
113 large displacements. Thus, we combine them by first per-
114 forming a TPS transformation for rough alignment and then
115 using appearance flow for refinement. Finally, a decom-
116 posed neural rendering (DNR) module is used to in-paint
117 the occluded skin area, refine the warped garment’s shad-
118 ing, and render the final try-on result.
119

120 However, the widely used VITON [14] dataset does not
121 meet our needs. It does not provide DensePose and clothes
122 parsing, and the annotations it has are of poor quality. We
123 also notice that most (1,575 out of 2,032) images in VI-
124 TON’s testing set also appear in its training set, which
125 would disturb the validation of a model’s generalizing abil-
126 ity. Thus, to facilitate our research and the community,
127 we collect a new dataset RAIIV with Rich Annotations for
128 Image-based Virtual try-on. RAIIV provides 16,052 high-
129 resolution image pairs ,as well as high-quality annotations.
130

131 Our contributions can be summarized as follows: (1) We
132 propose three improvements to image-based virtual try-
133 on, including new structural information (DensePose and
134 clothes parsing), cascade clothes warping for better shape
135 fitting and character preserving, and decomposed neural
136 rendering for skin inpainting and clothes shading refine-
137 ment. (2) We construct a high-resolution and high-quality
138 dataset RAIIV, which we believe could become a new bench-
139 mark for image-based virtual try-on. We will release the
140 dataset to facilitate the research community. (3) We con-
141 duct comprehensive module-by-module experiments, which
142 demonstrate the effectiveness of the new structural informa-
143 tion and the superiority of our SICW framework.
144

2. Related Work

2.1. Conditional Image Synthesis

145 Generative Adversarial Networks (GANs) [11] have
146 achieved great success in image generation, which revo-
147 lutionized many computer vision tasks, such as style
148 transfer [44, 38, 7], object detection [22], and image
149 inpainting [26]. Various techniques have been proposed
150 to optimize the vanilla GAN, among which conditional
151 GAN [23] (cGAN) becomes the most promising and widely
152 used method in image translation. cGAN adopt both the real
153 images and the random noises as the input of the generator,
154 making it possible to generate images while maintaining the
155 content of that source image. Built on cGAN, a series of
156 autoencoder-based image translation methods are proposed
157 to generate the realistic looking while content preserving
158 images [4, 5, 25]. Among them, PatchGAN [23] is a widely
159 used

160 used deblurring module to make the generated images more
161 photo-realistic, which is also adopted in our model to better
162 assist the skin in-painting process.
163

2.2. Virtual try-on

164 Virtual try-on allows people to see their overall appear-
165 ance without actually wearing the clothes. It is appearing to
166 both the customer and the company because of the enjoy-
167 ment and convenience when trying on numerous clothes.
168

169 Most traditional methods utilized 3D based models to
170 solve this problem. Some works [43, 6, 36, 15, 34, 2, 12, 27,
171 31] rely on the reconstruction and simulation of 3D cloth-
172 ing models to render the images, which produced reason-
173 able try-on results. However, such 3D-based methods are
174 usually built on the high-quality 3D reconstruction model,
175 which leads to slow rendering speed during the test period.
176 What’s more, such models require a large amount of high-
177 resolution 3D data, which is not accessible to ordinary cus-
178 tomers, making it further difficult to be employed in real-
179 world applications.
180

181 Recently, due to the superior representation ability of the
182 deep neural networks, a series of models are proposed to
183 solve this problem by only utilizing the 2D images. The
184 key challenge in such image-based methods is to model
185 the highly non-rigid transformation of the clothes. Most
186 of the existing methods [14, 37, 41, 20, 40, 28, 33] use
187 thin-plate spline (TPS) [1] to warp clothes. For exam-
188 ple, VITON [14] estimates TPS parameters by shape con-
189 text matching, which utilizes hand-crafted descriptors to
190 model the correspondences between the source and the tar-
191 get clothing shape. A composition mask is then utilized to
192 fuse the warped clothing with the coarse try-on result. CP-
193 VTON [37] instead proposes a two-stage pipeline which
194 includes clothing warping and coarse-to-fine try-on image
195 generation. Although this approach improves the preserva-
196 tion quality of clothing textures, it suffers from the distor-
197 tion of the body parts and clothing items that should not
198 be affected by the warped clothes’ layout. Later on, VT-
199 NFP [41] addressed this issue by estimating a human pars-
200 ing map after clothing warping, which helps preserve cloth-
201 ing and body details by indicating the region that should not
202 be affected by the warping stage. ACGPN [40] further im-
203 proved the preservation of the body detail by restricting the
204 estimated human parsing with the original image’s semantic
205 layout. Such models achieve reasonable results on a set of
206 basic transformations, however, the clothes would be heav-
207 ily distorted when the human pose becomes complicated.
208

209 On the other hand, ClothFlow [13] proposed a novel ap-
210 proach to tackle this problem, which predicts a dense pixel-
211 level transformation to address complex clothing deforma-
212 tions. However, flow-based warping still cannot handle
213 large deformations very well due to the lack of ground truth
214 flow. After systematically comparing and analyzing these
215

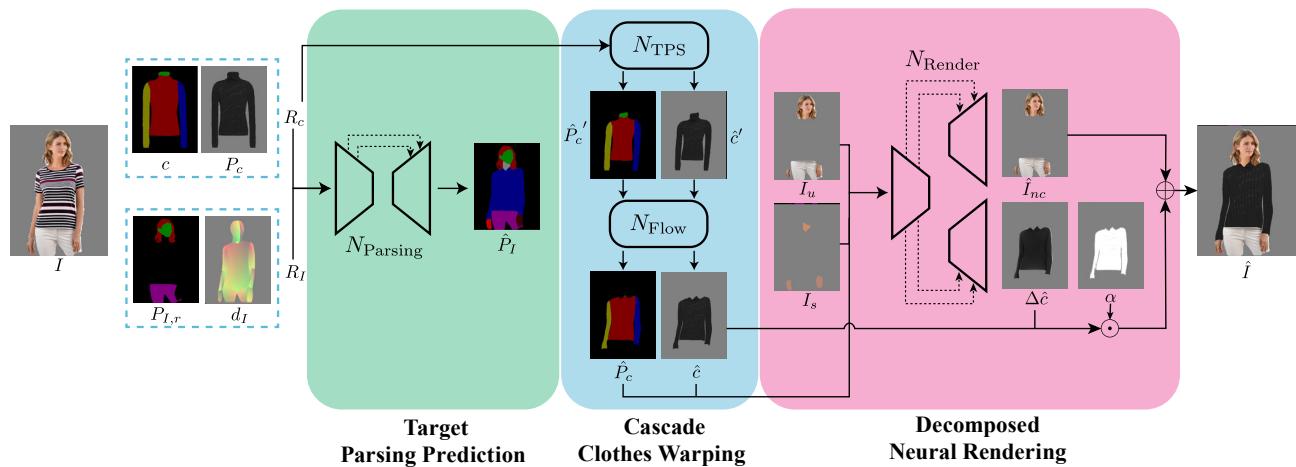


Figure 2. Model pipeline overview.

two clothing warping methods, we propose to combine the advantages of them, which leads to better shape fitting and texture preservation ability.

3. Our Approach

The task of image-based virtual try-on can be defined as follows: Given a clothing image (also known as target clothes) c and a person image (also known as reference image) I , we want to generate a photo-realistic image \hat{I} of that person wearing the c without changing the pose. However, the ground truth data for supervised learning is hard to get as it requires a person to change clothes while maintaining the pose. To address this dilemma, we use obtainable paired I and c to construct two integrated representations (See Sec. 3.1): a clothing-agnostic person representation [14] R_I extracted from I and a clothing representation R_c extracted from c . We then let our model learn to reconstruct I from R_I and R_c .

As shown in Fig. 2, the proposed SICW consists three modules. First, the target parsing prediction (TPP) network predicts the after-trying-on human parsing from rich structural information (See Sec. 3.2). Second, the cascade clothes warping (CCW) module deforms the input clothing c to the target shape on the predicted human parsing in a coarse-to-fine manner (See Sec. 3.3). Finally, the decomposed neural rendering (DNR) module is designed to generate a photo-realistic try-on result via skin inpainting and clothes shading refinement, as described in Sec. 3.4.

3.1. Person and Clothing Representations

The clothing-agnostic person representation [14] is designed for integrating descriptive and structural information for the person in I while eliminating the effects of the clothes. It is composed initially of a pose heatmap, a blurred body shape, and reserved regions (face and hair). We im-

prove this representation by adopting DensePose [30] and human parsing. Our enhanced clothing-agnostic representation R_I now consists of a 3-channel DensePose map and a 12-channel reserved human parsing that removes the tops and torso parsing label from the original human parsing of I .

Similarly, we construct a clothing representation R_c by integrating the 3-channel target clothing image c with its 5-channel clothes parsing P_c (explained later in Sec. 4).

3.2. Target Parsing Prediction

We first predict target human parsing, which is used to guide the clothes warping and body completion. Different from previous methods, our method takes cloting-agnosistic representation R_I and clothing representation R_c as inputs. The DensePose in R_I helps to exclude the influence of original shape of clothes in the reference image, and the clothes parsing helps to distinguish different parts of the clothes. Both annotations make the prediction on tops more correct and realistic.

As for network, an U-Net [32] N_{Parsing} is used to predict target human parsing, i.e. $\hat{P}_I = N_{\text{Parsing}}(R_I, R_c)$. In training procedure, we use paired data as training set. Note that the parsing of parts to be predicted is not available in inputs, the generation ability on unpaired data can be guaranteed. OHEM [35] loss and consistency loss are used in training procedure. OHEM loss calculates CE loss between P_I and \hat{P}_I on only the top-k most unconfident pixels, which improves the accuracy on details significantly. Consistency loss refers to the CE loss between results given transformed inputs and ground truth applied the same transformation, formulated as

$$\mathcal{L}_{\text{cons}} = \text{CE}(t(P_I), N_{\text{Parsing}}(t(R_I), t(R_c)))$$

where t refers to image transformation, and we use mirror transformation and affine transformation when training.

324

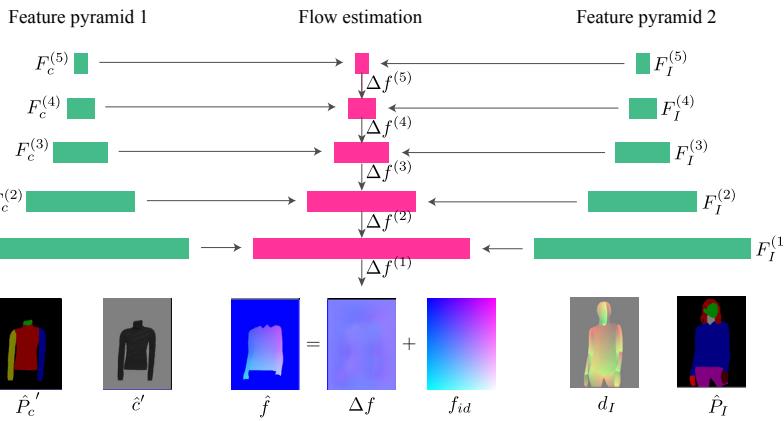


Figure 3. Flow model overview

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

340 The overall loss function can be formulated as

341
$$\mathcal{L} = \lambda_{ohem}\mathcal{L}_{ohem} + \lambda_{cons}\mathcal{L}_{cons} \quad (1)$$

344 where λ s are weights of each term to control the final performance.347

3.3. Cascade Clothes Warping

349 We then transform the clothing image c to fit the target shape $\hat{P}_{I,c}$, which is a 1-channel map extracted from the target human parsing \hat{P}_I . The clothes warping's ground truth I_c is cropped from the reference image I according to its human parsing.354 Since clothes deformations are highly non-rigid, we 355 adopt the cascade clothes warping (CCW) module to warp 356 clothes in a coarse-to-fine manner. Specifically, we first use 357 thin-plate spline (TPS) [1] transformation to roughly align 358 the clothes c with the target shape $\hat{P}_{I,c}$. TPS transformation 359 preserves the clothes details very well as it provides a 360 smooth interpolation between a set of control points. How- 361 ever, it has problems (Fig. 7) fitting the target shape due to 362 its limited degrees of freedom or handling occlusion since 363 it interpolates a continuous surface. Thus, we further use a 364 dense pixel flow (Appearance Flow [42]) to refine the warp- 365 ing result. Powered by DensePose and clothes parsing, the 366 flow-based warping can infer the body orientation and auto- 367 matically pick the non-occluded pixels. We do not use the 368 flow-based warping alone, as what has been done in [13], 369 because it produces blurred results when there is a large size 370 difference between the input clothes and the clothing on the 371 person.372 **Rough Alignment by TPS Transformation:** Formally, 373 given the clothing representation R_c , DensePose d_I , and the 374 target shape $\hat{P}_{I,c}$, we train a spatial transformation network 375 (STN) [19] N_{TPS} to estimate the coordinates of 25 con- 376 trol points $\hat{\theta} \in \mathbb{R}^{25 \times 2}$. We then apply the TPS transfor- 377 mation on the clothing representation, getting the correspond-

339 ing warped result by

380
$$\hat{P}_c' = \{\hat{c}', \hat{P}_c'\} = W_{\hat{\theta}}(R_c), \quad (2)$$

382 where $\mathcal{W}_{\hat{\theta}}(\cdot)$ is a bilinear sampling decided by $\hat{\theta}$, \hat{c}' is the 383 warped clothes, \hat{P}_c' is the warped clothes parsing.384 We adopt the second-order difference constraint \mathcal{L}_{reg} 385 from ACGPN [40], which makes the warping more sta- 386 ble by pushing the control points to be regularly arranged. 387 We also calculate the \mathcal{L}_1 Loss and perceptual loss [21] be- 388 tween \hat{c}' and the ground truth I_c , presented as $\mathcal{L}'_1(\hat{c}', I_c)$ 389 and $\mathcal{L}'_{perc}(\hat{c}', I_c)$, respectively.

390 Thus, the overall loss function for the STN is:

392
$$\mathcal{L}_{STN} = \lambda_{reg}\mathcal{L}_{reg} + \lambda'_1\mathcal{L}'_1 + \lambda'_{perc}\mathcal{L}'_{perc}, \quad (3)$$

394 where λ s are balancing weights.396 **Shape Refinement by Appearance Flow:** An appear- 397 ance flow [42] $f \in \mathbb{R}^{H \times W \times 2}$ specifies, for each pixel 398 $p^{(h,w)}$ in the output image, the coordinate at the input im- 399 age where the pixel value is sampled to reconstruct $p^{(h,w)}$. 400 Instead of estimating the flow directly [13], we let the net- 401 work output an offset Δf which denotes the pixel move- 402 ments and compute the final flow by $\hat{f} = f_{id} + \Delta f$. Here, 403 f_{id} denotes an identify flow that copies every input pixel 404 to the same place in the output image. Since adjacent pixels 405 tend to move together, i.e., they have the same offset, letting 406 the network learn to output the offset is easier than having 407 it output equally-spaced pixel coordinates.408 As shown in Fig. 3, we feed the warped clothing 409 representation \hat{P}_c' and human body information (d_I, \hat{P}_I) 410 into two separate encoders to get two feature pyramids 411 $\{F_c^{(l)}\}, \{F_I^{(l)}\}$ ($l = 1, 2, \dots, 5$). As for flow estima- 412 tion, we first feed the lowest-level features $F_c^{(5)}$ and $F_I^{(5)}$ 413 into a residual block [16] to predict the initial offset flow Δf^5 . 414 Furthermore, at each level l ($l = 1, 2, 3, 4$), the estimated 415 flow from the level above is upscaled by 2 and used for 416

432 estimating the current-level flow together with the features
 433 from pyramid 1 and features from pyramid 2:
 434

$$435 \Delta f^{(l)} = \text{ResBlock}^{(l)}(\mathcal{U}(\Delta f^{(l+1)}), F_c^{(l)}, F_I^{(l)}), \quad (4)$$

437 where $\mathcal{U}(\cdot)$ stands for the two-times bilinear upsampling.
 438

439 The final flow is computed by $\hat{f} = f_{id} + \mathcal{U}(\Delta f^{(1)})$, and
 440 is used for producing the final warping result, formulated as
 $\hat{c} = \mathcal{W}_{\hat{f}}(c')$.
 441

442 Similar to TPS transformation, we add the \mathcal{L}_1 Loss
 443 and perceptual loss on the final warping result, denoted as
 $\mathcal{L}_1(\hat{c}, I_c)$ and $\mathcal{L}_{perc}(\hat{c}, I_c)$, respectively. Besides, we adopt
 444 the flow regularization term and intermediate flow super-
 445 visions from [29] to make the training process stable and
 446 encourage smooth flow fields.
 447

448 3.4. Decomposed Neural Rendering

449 After getting target human parsing \hat{P}_I and the warped
 450 clothing representation \hat{R}_c , we finally use the DNR mod-
 451 ule to generate the final try-on result by fusing the warped
 452 clothes with the reference person, where we further decom-
 453 pose the task into skin inpainting and clothes shading re-
 454 finement.
 455

456 Following ACGPN [40], we first compute the nonaf-
 457 fected body parts I_u by comparing the reference human
 458 parsing P_I with the predicted target parsing \hat{P}_I . Intuitively,
 459 I_u denotes the body parts on I that will not be affected
 460 by the clothes' change, thus should be preserved through
 461 the rendering. Moreover, we compute the average color for
 462 each skin region (i.e. left arm, right arm, and torso) and fills
 463 the corresponding area \hat{P}_I with it. The computed region-
 464 specific skin color map I_s provides crude guidance for the
 465 skin inpainting process.
 466

467 Finally, our rendering network N_{render} takes the target
 468 human parsing \hat{P}_I , the warped clothing representation \hat{R}_c ,
 469 and the additional guidance computed above and get three
 470 outputs, formulated as:
 471

$$472 \hat{I}_{nc}, \Delta\hat{c}, \alpha = N_{\text{render}}(\hat{P}_I, \hat{R}_c, I_u, I_s), \quad (5)$$

473 where \hat{I}_{nc} refers to generated non-clothing image, $\Delta\hat{c}$
 474 refers to shading to be added onto warped clothes, and α
 475 is an alpha mask for fusion. To free up network capacity,
 476 we use a shared encoder and two separate decoder, each
 477 predicting \hat{I}_{nc} and $\Delta\hat{c}$ with the alpha mask.
 478

The final try-on result \hat{I} is then calculated by:

$$479 \hat{I} = \hat{I}_{nc} + \Delta\hat{c} \odot \alpha + \hat{c} \odot (1 - \alpha) \quad (6)$$

480 The training objective consists of the \mathcal{L}_1 Loss, percep-
 481 tual loss, and adversarial loss. We use three PatchGAN [18]
 482 discriminators to improve the generation photo-realism, in-
 483 cluding D that works on the whole image, D_{Skin} that fo-
 484 cuses on the skin region, and D_{Clothes} that focuses on the
 485 clothes region.
 486

	VITON [14]	RAIV	486
Dataset Size	16,253	16,052	487
Resolution	256×192	≥ 512×384	488
Person	Person Image	✓	489
	Human Silhouette ²	✓	490
	Human Parsing ¹	✓	491
	OpenPose [3]	✓	492
	DensePose [30]	✗	493
Clothes	Clothing Image	✓	494
	Clothing Silhouette ²	✓	495
	Clothing Parsing	✗	496

497 Table 1. A comparison of VITON and our RAIVdataset. Note that:
 498 1) The human parsing annotation [9] in RAIVinvolves full-body
 499 parts, while the one [10] used in VITON lacks the torso region.
 500 2) Human silhouette and clothing silhouette can be extracted from
 501 human parsing and clothing parsing, respectively.
 502

503 4. RAIV Dataset

504 We collect a large-scale image-based virtual try-on
 505 dataset named RAIV, containing 16,052 pairs of high-
 506 resolution upper clothing and person images together with
 507 rich annotations. Table 1 presents an overview compari-
 508 son between RAIVand the widely used VITON [14] dataset.
 509 Specifically, RAIVhas two appealing properties:
 510

511 • **Rich annotations.** RAIVprovides human parsing, Open-
 512 Pose [3], and DensePose [30] for person images, and
 513 clothing parsing for clothing images. We use a state-of-
 514 the-art human parser [8] to compute the human parsing,
 515 which follows the CIHP [9] annotation to have 19 se-
 516 mantic labels. OpenPose stands for 25 human body key-
 517 points. DensePose provides some 3D information for hu-
 518 man body and is clothes-agnostic. As for clothing pars-
 519 ing, we first manually annotate 6,000 clothes with five la-
 520 bels ('background', 'outer surface', 'inner surface', 'left
 521 sleeve', and 'right sleeve', shown in Fig. 4), and then train
 522 a parsing model [8] to get the rest.
 523

524 • **High quality.** As shown in Fig. 5, RAIVprovides high-
 525 quality full body human parsing comapred to VITON.
 526 Moreover, we adopt a strict reviewing process to elimi-
 527 nate flawed images and annotations, such as incomplete
 528 Densepose and wrong clothing parsing. Besides, the im-
 529 age resolution in RAIVis at least 512 × 384, which is two
 530 times better than that in VITON.
 531

532 5. Experiments

533 We conduct stage-by-stage ablation studies and com-
 534 parative experiments with the state-of-the-art models CP-
 535 VTON [37] and ACGPN [40]. First, we evaluate the impor-
 536 tance of introducing structural information (clothing pars-
 537 ing and DensePose) into the target human parsing predic-
 538 tion.
 539



540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
Figure 4. Clothing parsing labels (black: background, red: outer surface, green: inner surface, blue: left sleeve, yellow: right sleeve).



561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
Figure 5. A representative comparison of the human parsing quality between VITON [14] and RAI.

tion stage. Next, we evaluate how structural information and cascade warping can benefit clothes warping. Finally, we link three stages as a whole model and evaluate its ability to generate photo-realistic try-on images. While all the above methods use paired data for training, their performances are evaluated under the both paired and unpaired testing settings. We also comprehensively compare and analyze our models under both settings. All the experiments are conduct on the widely used VITON [14] dataset and our proposed RAI dataset. More details of the experiments are shown in the supplementary materials.

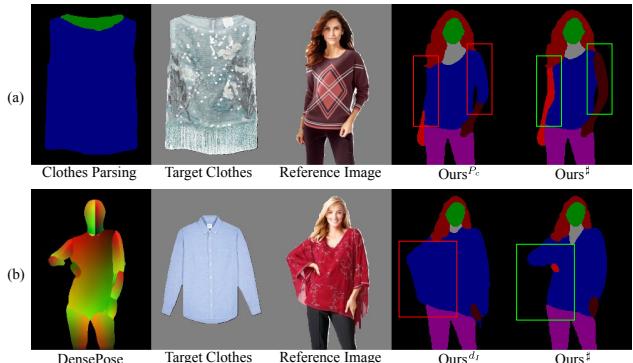
5.1. Evaluating Target Parsing Prediction

We compare the target human parsing prediction performance with ACGPN. Both methods are trained on our RAI dataset for 20 epochs. Table 2 lists the IoU measurements on paired data. We take the prediction of four semantic parts ('Tops', 'Torso', 'Left Arm', and 'Right Arm') that have the most influence on the final result. It shows that our full model (F) outperforms ACGPN (A) on all these parts.

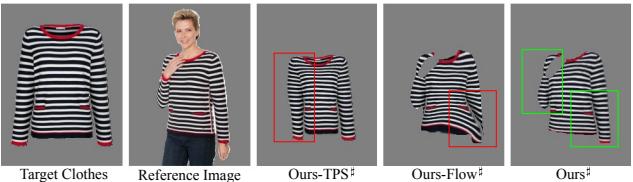
Besides, an ablation study (C-F) is conducted to prove the structural information's effectiveness. We can conclude that adding clothing parsing (C to D) improves all selected parts' prediction accuracy. The unpaired testing results in Fig. 6 (a) verifies this improvement. We also notice that adding DensePose (C to E) improves all but the 'Tops' part, which is caused by the paired testing setting, where human silhouettes used in (C) indicate the tops' boundary while DensePose in (E) is clothing agnostic. However, when test-

IDX	Method	IoU ↑				594
		Tops	Torso	Left Arm	Right Arm	
(A)	ACGPN [40]	0.799	0.664	0.731	0.767	595
(B)	ACGPN [40] [#]	0.919	0.767	0.799	0.824	596
(C)	Ours	0.948	0.672	0.627	0.670	597
(D)	Ours ^{P_c}	0.956	0.746	0.753	0.781	598
(E)	Ours ^{d_I}	0.929	0.725	0.792	0.825	599
(F)	Ours [#]	0.933	0.768	0.797	0.832	600

598
599
600
601
602
603
604
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
Table 2. Evaluation of Stage1 on RAI dataset. ↑ means bigger value is better; P_c denotes using clothing parsing; d_I denotes using DensePose; [#] denotes using both information.



630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
Figure 6. Ablation study on the use of clothes parsing and DensePose as inputs. (a) Clothes parsing helps improve the prediction performance when different clothes parts are hard to distinguish. (b) DensePose helps to exclude the influence of the origin shape of the reference image.



643
644
645
646
647
Figure 7. Ablation study on the effect of TPS and Flow. TPS transformation fails to fit the pose, and appearance flow causes distortion on texture. The combination of these two methods avoids both drawbacks.

ing under the real-world scenario, i.e., testing with unpaired data, human silhouettes will mislead the prediction while DensePose will not, as shown in Fig. 6 (b). Moreover, adding the new information to ACGPN (B) leads to significant improvements to its original implementation (A) and shows comparable results with ours (F), which further validates the importance of the rich guiding information.

5.2. Evaluating Cascade Clothes Warping

We evaluate different clothing warping methods with the paired testing setting. For the sake of fairness, we test each clothes warping methods independently by fetching inputs

	IDX	Method	SSIM ↑	MS-SSIM ↑
VITON	(A)	CP-VTON [37]	0.825	0.830
	(B)	ACGPN [40]	0.801	0.797
	(C)	Ours	0.873	0.890
RAIV	(D)	CP-VTON [37]	0.828	0.832
	(E)	CP-VTON [37]‡	0.831	0.837
	(F)	ACGPN [40]	0.824	0.815
	(G)	ACGPN [40]‡	0.828	0.821
	(H)	Ours	0.861	0.876
	(I)	Ours-TPS‡	0.821	0.825
	(J)	Ours-Flow‡	0.871	0.870
	(K)	Ours‡	0.880	0.889
	(L)	Ours-512‡	0.892	0.887

Table 3. Quantitative comparisons and ablation study of the clothing warping module. **Bolds** are the best results; ↑ means bigger value is better; ‡ denotes using our structural information; Ours-TPS and Ours-Flow are for the ablation study, which use only TPS or pixel flow for warping; Ours-512 denotes our model trained on 512×384 resolution. Note that all inputs are taken from the corresponding dataset, so warping results are not affected by the previous stage, if exists.

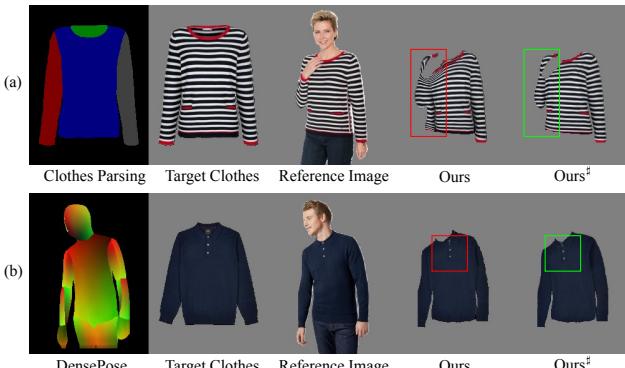


Figure 8. Ablation study of DensePose and clothes parsing. The clothes parsing prevents CCM from mixing sleeves and body together, and the DensePose helps CCM to predict the correct direction of reference person.

from the dataset.

For quantitative evaluation, we adopt SSIM and MS-SSIM [39] to measure the similarity between the warped clothing and the cropped clothing from the reference image). Table 3 lists the results of CP-VTON [37], ACGPN [40], and our SICW on both VITON [24] and RAIv datasets. It shows that our method (C, L) substantially outperforms the other methods on both datasets.

Besides, the ablation study (I-K) shows that cascade warping indeed improves the clothes warping quality. The qualitative comparison in Fig. 7 further proves this conclusion, where we see that TPS fails to fit the target shape while flow leads to distortions on the bottom of the warped garment.

Moreover, we study the effects of the new structural information. We see that all methods gain improvements (D

	IDX	Method	SSIM ↑	MS-SSIM ↑	FID ↓	
					Paired	Unaired
VITON	(A)	CP-VTON [37]	0.819	0.899	18.75	39.76
	(B)	ACGPN [40]	0.867	0.916	11.63	16.29
	(C)	Ours	0.892	0.949	10.25	15.80
RAIV	(D)	CP-VTON [37]	0.613	0.802	35.77	40.23
	(E)	CP-VTON [37]‡	0.661	0.805	35.42	39.26
	(F)	ACGPN [40]	0.843	0.893	18.70	22.41
	(G)	ACGPN [40]‡	0.846	0.894	18.03	21.55
	(H)	Ours	0.865	0.902	15.61	21.11
	(I)	Ours‡	0.864	0.915	14.89	19.63
	(J)	Ours-512‡	0.879	0.906	13.54	18.27

Table 4. Quantitative comparisons on final try-on results. **Bolds** are the best results; ↑ means bigger value is better; ↓ means smaller value is better; ‡ denotes using our structural information; Ours-512 denotes our model trained on 512×384 resolution.

to E, F to G, and H to K) by using the information, while our method is the most beneficiary one. The representative comparison in Fig. 8 further validates our claim.

We also train our model on the 512×384 resolution (L) and get better results to the one trained on the 256×192 resolution, which shows our model’s ability to generalize to higher resolutions.

5.3. Evaluating Final Try-on Results

We compare the final try-on results with VITON and ACGPN models on both VITON and RAIv datasets. For testing on paired data, we adopt SSIM and MS-SSIM as evaluation metrics. We also use FID [17] to measure the visual quality of synthesized images on both paired-setting (denoted as ‘Paired’) and unpaired-setting (denoted as ‘Unpaired’). Table 4 lists the quantitative comparisons. Our model clearly outperforms the other methods and becomes the state-of-the-art one in both datasets. We also see that each model gains some improvements by adopting the new structural information (D to E, F to G, and H to J). These improvements are not as significant as what we get in previous stages (our model even drops 0.001 on SSIM). Like what we mentioned in Sec. 5.1, this phenomenon is due to the limitation of the paired testing setting. When testing on unpaired data, as shown in Fig. 10, using the new structural information improves the try-on results.

We compare the final try-on results with VITON and ACGPN on both VITON and RAIv datasets. For testing on paired data, we adopt SSIM and MS-SSIM as evaluation metrics. We also use FID [17] to measure the visual quality of synthesized images on both paired-setting (denoted as ‘Paired’) and unpaired-setting (denoted as ‘Unpaired’). Table 4 lists the quantitative comparisons. Our model clearly outperforms the other methods and sets a new state of the art for both datasets. We also see that each model gains some improvements by adopting the new structural information (D to E, F to G, and H to J). These improvements



Figure 9. Comparison between three virtual try-on methods on clothes with simple texture (left side) and complex texture (right side). Our method outperforms the others on three aspects: (1) clothes pose matching (*e.g.* top right) (2) texture preserving (*e.g.* bottom right) (3) body parts completing (*e.g.* middle left). All these advantages make the final result more realistic and natural.

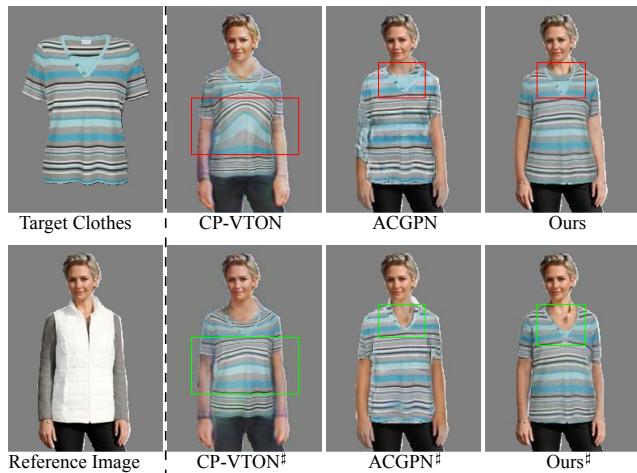


Figure 10. Comparison of three methods with and without clothes parsing and DensePose as inputs. Results show that all methods benefit from the new inputs that the performances on neckline and pose fitting are improved significantly.

are not as significant as what we get in previous stages (our model even drops 0.001 on SSIM). Like what we mentioned in Sec. 5.1, this phenomenon is due to the limitation of the paired testing setting. When testing on unpaired data, as shown in Fig. 10, using the new structural information improves the try-on results.

Moreover, We show qualitative comparisons in Fig. 9, including clothes with simple textures (on the left side) and complex textures (on the right side). Since CP-VTON does not predict the target human parsing, it often fail to fit the

garment to the human body correctly, as shown in the top right results of Fig. 9. Besides, the other two methods use TPS transformation only to warp the clothes, where the preserves parts could be occluded in the final results, while our method uses appearance flow to abandon occluded parts. This technique makes our method perform better on neckline (*e.g.* bottom left result). Additionally, in CP-VTON, the single TPS transformation causes distortion. ACGPN tries to handle this problem by introducing image generation, which makes the result more blurry. In contrast, our method performs better on texture preserving by cascade cloth warping and minor necessary shading, which not only avoids distortion (*e.g.* top right result) but also provides clearer texture (*e.g.* bottom right result).

6. Conclusion

We presented RAIv, a new image-based virtual try-on dataset with 16,052 image pairs and rich annotations. It surpasses the widely used VITON dataset in terms of image resolution, annotation quality, and annotation kind. Taking advantage of RAIv, we proposed a novel model that consists of three stages. Extensive experiments show that our SICWmodel outperforms state of the art on both VITON and RAIvdatasets.

864
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

References

- [1] Fred L. Bookstein. Principal warps: Thin-plate splines and the decomposition of deformations. volume 11, pages 567–585. IEEE, 1989. 2, 4
- [2] Remi Brouet, Alla Sheffer, Laurence Boissieux, and Marie-Paule Cani. Design preserving garment transfer. 2012. 2
- [3] Z. Cao, G. Hidalgo Martinez, T. Simon, S. Wei, and Y. A. Sheikh. Openpose: Realtime multi-person 2d pose estimation using part affinity fields. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019. 5
- [4] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *CVPR*, 2018. 2
- [5] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. Stargan v2: Diverse image synthesis for multiple domains. In *CVPR*, 2020. 2
- [6] Jun EHara and Hideo Saito. Texture overlay for virtual clothing based on pca of silhouettes. In *2006 IEEE/ACM International Symposium on Mixed and Augmented Reality*, pages 139–142. IEEE, 2006. 2
- [7] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *CVPR*, 2016. 2
- [8] Ke Gong, Yiming Gao, Xiaodan Liang, Xiaohui Shen, Meng Wang, and Liang Lin. Graphonomy: Universal human parsing via graph transfer learning. In *CVPR*, 2019. 5
- [9] Ke Gong, Xiaodan Liang, Yicheng Li, Yimin Chen, Ming Yang, and Liang Lin. Instance-level human parsing via part grouping network. In *ECCV*, 2018. 5
- [10] Ke Gong, Xiaodan Liang, Dongyu Zhang, Xiaohui Shen, and Liang Lin. Look into person: Self-supervised structure-sensitive learning and a new benchmark for human parsing. In *CVPR*, 2017. 5
- [11] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NeurIPS*, 2014. 2
- [12] Peng Guan, Loretta Reiss, David A Hirshberg, Alexander Weiss, and Michael J Black. Drape: Dressing any person. *ACM Transactions on Graphics (TOG)*, 31(4):1–10, 2012. 1, 2
- [13] Xintong Han, Xiaojun Hu, Weilin Huang, and Matthew R Scott. Clothflow: A flow-based model for clothed person generation. In *ICCV*, 2019. 2, 4
- [14] Xintong Han, Zuxuan Wu, Zhe Wu, Ruichi Yu, and Larry S Davis. Viton: An image-based virtual try-on network. In *CVPR*, 2018. 1, 2, 3, 5, 6
- [15] Stefan Hauswiesner, Matthias Straka, and Gerhard Reitmayr. Virtual try-on through image-based rendering. *IEEE transactions on visualization and computer graphics*, 19(9):1552–1565, 2013. 2
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 4
- [17] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *NeurIPS*, 2017. 7
- [18] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, 2017. 5
- [19] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. In *NeurIPS*, 2015. 4
- [20] Surgan Jandial, Ayush Chopra, Kumar Ayush, Mayur Heman, Balaji Krishnamurthy, and Abhijeet Halwai. Sievenet: A unified framework for robust image-based virtual try-on. In *WACV*, 2020. 2
- [21] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, 2016. 4
- [22] Jianan Li, Xiaodan Liang, Yunchao Wei, Tingfa Xu, Jiashi Feng, and Shuicheng Yan. Perceptual generative adversarial networks for small object detection. In *CVPR*, 2017. 2
- [23] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014. 2
- [24] Assaf Neuberger, Eran Borenstein, Bar Hilleli, Eduard Oks, and Sharon Alpert. Image based virtual try-on network from unpaired data. In *CVPR*, 2020. 7
- [25] Taesung Park, Jun-Yan Zhu, Oliver Wang, Jingwan Lu, Eli Shechtman, Alexei A. Efros, and Richard Zhang. Swapping autoencoder for deep image manipulation. *arXiv preprint arXiv:2007.00653*, 2020. 2
- [26] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders: Feature learning by inpainting. In *CVPR*, 2016. 2
- [27] Gerard Pons-Moll, Sergi Pujades, Sonny Hu, and Michael J Black. Clothcap: Seamless 4d clothing capture and retargeting. *ACM Transactions on Graphics (TOG)*, 36(4):1–15, 2017. 2
- [28] Amir Hossein Raffiee and Michael Sollami. Garmentgan: Photo-realistic adversarial fashion transfer. *arXiv preprint arXiv:2003.01894*, 2020. 2
- [29] Yurui Ren, Xiaoming Yu, Junming Chen, Thomas H Li, and Ge Li. Deep image spatial transformation for person image generation. In *CVPR*, 2020. 5
- [30] Iasonas Kokkinos, Riza Alp Güler, Natalia Neverova. Densepose: Dense human pose estimation in the wild. In *CVPR*, 2018. 1, 3, 5
- [31] Damien Rohmer, Tiberiu Popa, Marie-Paule Cani, Stefanie Hahmann, and Alla Sheffer. Animation wrinkling: augmenting coarse cloth simulations with realistic-looking wrinkles. *ACM Transactions on Graphics (TOG)*, 29(6):1–8, 2010. 2
- [32] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 3
- [33] Debapriya Roy, Sanchayan Santra, and Bhabatosh Chanda. Lgton: A landmark guided approach to virtual try-on. *arXiv preprint arXiv:2004.00562*, 2020. 2
- [34] Igor Santesteban, Miguel A Otaduy, and Dan Casas. Learning-based animation of clothing for virtual try-on. In *Computer Graphics Forum*, volume 38, pages 355–366. Wiley Online Library, 2019. 1, 2

- 972 [35] Abhinav Shrivastava, Abhinav Gupta, and Ross Girshick.
973 Training region-based object detectors with online hard ex-
974 ample mining. In *CVPR*, 2016. 3
975 [36] Hiroshi Tanaka and Hideo Saito. Texture overlay onto flex-
976 ible object with pca of silhouettes and k-means method for
977 search into database. In *MVA*, pages 5–8, 2009. 2
978 [37] Bochao Wang, Huabin Zheng, Xiaodan Liang, Yimin
979 Chen, Liang Lin, and Meng Yang. Toward characteristic-
980 preserving image-based virtual try-on network. In *ECCV*,
981 2018. 1, 2, 5, 7
982 [38] Xiaolong Wang and Abhinav Gupta. Generative image mod-
983 eling using style and structure adversarial networks. In
984 *ECCV*, 2016. 2
985 [39] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Si-
986 moncelli. Image quality assessment: from error visibility to
987 structural similarity. *IEEE transactions on image processing*,
988 13(4):600–612, 2004. 7
989 [40] Han Yang, Ruimao Zhang, Xiaobao Guo, Wei Liu, Wang-
990 meng Zuo, and Ping Luo. Towards photo-realistic virtual
991 try-on by adaptively generating-preserving image content. In
992 *CVPR*, 2020. 1, 2, 4, 5, 6, 7
993 [41] Ruiyun Yu, Xiaoqi Wang, and Xiaohui Xie. Vtnfp: An
994 image-based virtual try-on network with body and clothing
995 feature preservation. In *ICCV*, 2019. 2
996 [42] Tinghui Zhou, Shubham Tulsiani, Weilun Sun, Jitendra Ma-
997 lik, and Alexei A Efros. View synthesis by appearance flow.
998 In *ECCV*, 2016. 2, 4
999 [43] Zhenglong Zhou, Bo Shu, Shaojie Zhuo, Xiaoming Deng,
1000 Ping Tan, and Stephen Lin. Image-based clothes animation
1001 for virtual fitting. In *SIGGRAPH Asia 2012 Technical Briefs*,
1002 pages 1–4. 2012. 2
1003 [44] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A
1004 Efros. Unpaired image-to-image translation using cycle-
1005 consistent adversarial networks. In *ICCV*, 2017. 2
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025