

000

# Image-based Virtual Try-On: its limitations and 001 an improvement

002

003

004 Anonymous ECCV submission

005

006 Paper ID ...

007

008

009 **Abstract.** Recently, a series of studies on virtual try-on (VTON) using  
010 a try-on cloth and human image have been published. These algorithms  
011 are composed of two stages: (1) warping the try-on cloth to align with  
012 the pose and shape of the target model, and (2) blending the warped  
013 cloth onto the target human image. Our classified/strategic comparison  
014 study shows that CP-VTON generates the best quality image among  
015 SCM-based non-deep learning method, and deep learning-based VITON  
016 and CP-VTON. However, we identified 5 key problems of CP-VTON,  
017 such as improper human segmentation labelling, the pixel generation of  
018 un-intended areas, missing warped cloth mask and the cost function used  
019 in the learning. Tacking the issues, a new refined pipeline, CP-VTON+ is  
020 proposed. CP-VTON+ shows consistent improvements in SSIM, LPIPS,  
021 and IS, and outperforms the previous ones significantly in qualitative  
022 evaluations.

023 **Keywords:** Virtual Try-On, Image-based, Deep-Learning, Quality Com-  
024 parison

025

## 026 1 Introduction

027

028 Online fashion market has been growing rapidly every year. Clothing purchasing  
029 decisions are very difficult to make with current non-customized information,  
030 like the cloth and models' try fit images. Unlike other products, such as elec-  
031 tronic devices, whose function, performance, and styles can be expected through  
032 few images and specification tables. Fashion apparels have infinite variations in  
033 style, forms, colors, texture, and materials. Also the difference between personal  
034 preferences is huge. Therefore, virtual try-on (VTON) is a highly demanding  
035 technology for the on-line shopping.

036 The early VTON technologies were based on 3D computer graphics technol-  
037 ogy that uses 3D models of target humans and clothing. The 3D models are  
038 usually expensive and difficult to obtain. Therefore, recently 2D image-based  
039 VTON technologies are being studied in academia and industry, powered by the  
040 recent advances in computer vision technologies based on deep learning.

041 There have been many studies with different problem settings related to  
042 image-based VTON, from clothed human pose transferring using conditional  
043 GAN [8], swapping two humans clothes [6], to VTONs with a try-on cloth and a  
044 target human image [4]. The last configuration with a try-on cloth and a target

human image has been considered practical in many papers, [4, 10], and more recently [9, ?,?]. In this paper, we also consider the VTON problem that use the try-on cloth and human images and generated a new virtual image that the target human replaced the current top or bottom cloth with the try-on cloth. Our implementation is also limited to top clothes due to the restricted dataset but the bottom clothes, e.g. pants or skirts would be easier than top cloth cases because they are simpler than upper clothes in style and shapes.

Although the previous image-based VTON studies shows high quality VTON results, our classified analysis on the cloth styles and human posture in the Section 2 reveals significant problems in the previous works. In this paper, we started with evaluating the quality of SCM[1] based-VTON, VITON [4], and CP-VTON [10], using a cloth and a human image, according to the posture and body type of the human, the degree of occlusion of the clothes, and the texture of clothes with IoU of warped clothes and SSIM of final re-try-on image result. The more recent works in 2019 could not included this work, the general pattern of image based algorithm would be similar to the three algorithms. Even though CP-VTON generates the best quality outputs among them, all algorithms shows similar problems, the mistakes in human representation, improper network cost function, and the inherent limitations of 2D-based approaches. We would emphasize here that one reason of the seemingly high quality in the existing algorithms are mainly due to the dataset with low-complexity bias, i.e., most clothes are short-sleeved and monochromatic, and the poses of humans are mostly in an up-right position. Specifically, as it will be shown in the following section, the results with the long-sleeved cloth arm and body posed shows very low quality. In Section 3, we point out 5 serious issues in previous works including CP-VTON [10].

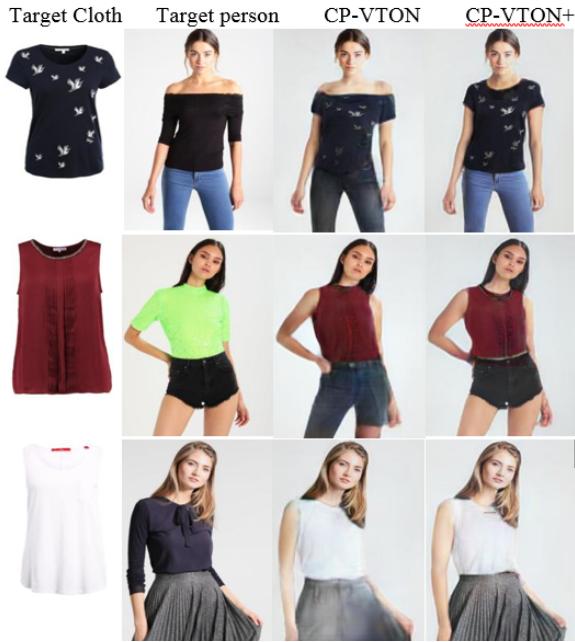
The key contribution of the papers are 3 ways: First, we provide the classified performance evaluations of existing Image-based algorithms. Second, the origins are the limitations and reason of seemingly well working of the existing algorithms are identified. Third, the direct solutions are proposed to tackles the identified problems. The proposed method might be not the optimal way to solve the problem but it can explicitly show the level of effects of identified issues in the final results.

## 2 Classified Performance Evaluation of Image-Based VTON

### 2.1 Image-based VTON

In this Section, we started with evaluating the 2D image based VTON algorithms. We considered CP-VTON [10] published in 2018 as the benchmark algorithm. The previous and following in 2019 share same input image and information conditions with CP-VTON and compare the results with it. Here we include the SCM based-VTON, VITON[4], and CP-VTON [10], however, we believe the performance strength and weakness are similar in the other algorithms too.

The Image based VTON algorithms are mostly composed of two stages: (1) cloth warping step that warps the try-on cloth to align with the pose and shape



**Fig. 1.** The proposed VTON results

of the target model (called GMM in CP-VTON: geometric Manipulation Module)[10], and (2) blending step that blends the warped cloth onto the target human image (called TON in CP-VTON: Try-On Network)[10]. CP-VTON assumes the target human image is pre-processed for a cloth agnostic human representation by a human pose estimation like OpenPose[3] and human parsing like LIP[7]. The human representation is composed of 1) heat maps for each joints 2) silhouette of human body, and 3) face and skin pixels patches (non-cloth and human identity area). We use the same dataset collected by Han et al. used in VITON[4] and CP-VTON[10] papers.

## 2.2 Classification Rule

First of all, as shown in Table 1, the criteria for classifying the experimental samples were divided into the degree of occlusion of the costume, the pose of the subject, and the complexity of the costume itself. The degree of obscuration is a factor that affects the accuracy of the object of deformation, the posture is the degree of deformation, and the complexity of the clothes means the processing complexity of the clothes themselves. However, it is included in the range of classification, but not included in the actual experiment is shown in parentheses. Excluded conditions are those that are not included in the test data or that the evaluation is considered to be complex in the current technology. Based on this, six cases were classified as follows.

- B: Little obscuration and posture (long sleeves, short arms)
- OP: Same as S, but partly covered by hair and arms
- OB: A large part of the clothing is covered by the bottoms.
- OF: If the front of the clothing is covered by the arms (long sleeves, short sleeves)
- P: If there is a large posture deformation (large movement of the arm or twisted or lateral posture)
- S: When there is a large body shape change (all or part of the body is thick or pregnant)

For reference, the costume of the data used is a T-shirt (without a collar) or the like, and the costume is generally formed in a simple form. It is also worth noting that most clothing is limited, such as monochromatic or monotone patterns.

Although several to tens of images were used for each type of experiment, they are not included in the paper due to the relationship between pages. 5 and Fig. 6, only representative images are presented for explanation.

#### – Wearing the Same Costume

In the same clothing experiment, performance was compared based on IoU, SSIM and visual results.

- B: All three algorithms showed high performance for short arm costumes. In particular, SCMM showed high performance. However, for long arms, the deformation of the arm was not good. This seems to be because the matching and deformation are mainly based on the whole silhouette. VITON and CP-VTON show that they are synthesized by complementing some of them in the synthesis process. Including the skeletal information can be supplemented, but at present, no algorithm has been developed for automatically extracting the skeletal information of the garment.
- OP: Partial obstruction caused by hair, etc., caused this part to be excluded from the scope of the purpose, and had some influence on the deformation of the clothes. VITON and CP-VTON also exhibited the same problem because they included head and skin in their representation. Therefore, it is expected to improve the GMM part by removing elements such as hair and using only body shapes.
- OB: Occlusion due to bottoms occurs. SCMM algorithms do not distinguish between deformation and occlusion. Therefore, IoU shows a deterioration, and especially the long arm is reduced and the skin is lifted up. Since VITON and CP-VTON use body information rather than box body, this effect seems to be reduced.
- OF: When the front part is covered by the arm, human expression can be used to distinguish the area of the clothes from the area of ??the arm. However, this part may not be clear during deep learning synthesis.
- P: If there is a large pose change, all three algorithms have a big error in the clothing deformation. This is considered to be a big limitation using the two-dimensional algorithm.

- 180     • S: In the case of the data of the same costume, the costume itself was  
181        largely prepared, so the error was not large.

182     – Experiment with wearing a new costume

183        The wearing of a new outfit is, in effect, the ultimate result of the application.  
184        As above, objective evaluation would be possible, but at present, one  
185        model could not present dataset that has more than two costumes. Although  
186        limited, it is possible to compare the relative differences of each algorithm  
187        based on visual results.

- 188        • B: The clothing deformation itself shows good results similar to the result  
189        of the same costume, but there are some differences according to the algo-  
190        rithms in the synthesis. When switching from long to short arms, SCMM  
191        does not restore the skin color of the arm and VITON and CP-VTON  
192        seem to produce it. In particular, CP-VTON shows excellent generating  
193        ability. This is an advantage of using pix2pix-based deep learning over  
194        non-deep learning.

- 195        • OP: The performance was similar to that of the same clothes.  
196        • OB: The same characteristics as the same image, that is, the SCMM  
197        showed a problem that can not distinguish between the mask and defor-  
198        mation.  
199        • OF: The same characteristics as the costume. Here too, in the case of  
200        SCMM, the synthesis algorithm needs to be improved when switching  
201        from short arm to long arm.  
202        • P: As in the case of the same costume, there was a large error in defor-  
203        mation.  
204        • S: It showed some adaptation to body shape change. In particular, VI-  
205        TON and CP-VTON have been shown to adapt very well to body shape  
206        changes.

207        In addition to the analysis of each condition in addition to the analysis of  
208        each condition, it can be seen that there are the following big features. First,  
209        it was confirmed that the shape change of clothes by GMM has an influence  
210        on the current wear clothes. The reason is that SCMM uses the area of the  
211        current costume, and deep learning methods use the body itself, but the area  
212        of the current costume is reflected in the correct answer mask used in the  
213        learning process.

## 214 215     2.3 Analysis Summary

216

217        Even though the success and failure cases are presented and compared with other  
218        algorithms' results, the failure case analysis is not enough for understanding the  
219        origin of failure cases and therefore difficult to find the solution for them. A  
220        classified evaluation would be better for this understanding. Here we summarized  
221        the classified results from our another study. We classify input try-on cloth and  
222        target human images according to the posture and body type of the person, the  
223        degree of occlusion of the clothes, and the characteristics of the clothes. Quality  
224        is compared in IoU for the warping step and in SSIM for the final blending step

for same cloth re-try-on cases. We also tested for the new cloth try-on cases but did not include here for limitation of spaces, and the same cloth cases are enough to explain the tendencies of the performances. Though in general CP-VTON generates the best quality image, the relative comparison is not the main purpose of the analysis. Please refer to Figure ?? for detailed results.



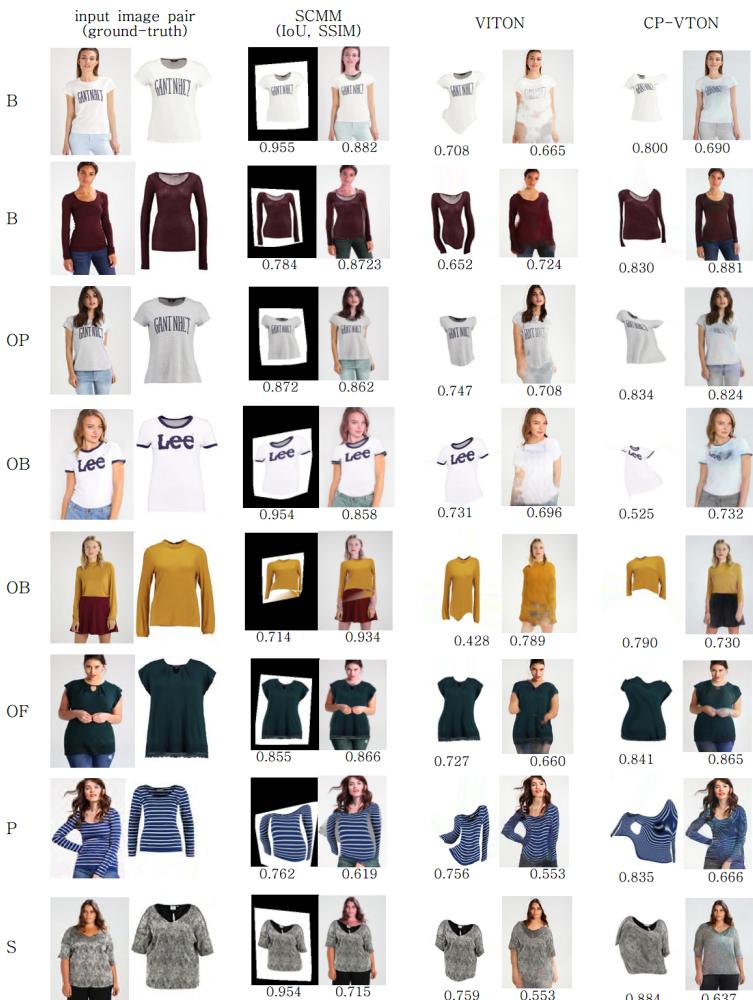
**Fig. 2.** The issues in CP-VTON

Firstly, the target try-on area is dependent upon current cloth shape. Especially, the neck area pixels are labeled as background and some body areas are occluded by hairs or accessories (Fig. 3 (a) left), which affects in cloth warping and blending. Secondly, all the unintended part, faces, bottom-clothes and legs have to be preserved in blending stage. But other parts except face and hair are missing in CP-VTON[10] human representation and generated at blending stage, which is all right for general synthesis application but not desirable in VTON application (Fig. 3 (b) left). Thirdly, the texture is often not vivid, which is due to the composition. Examining the original loss function of TON network, the term for the composition alpha mask are poorly formulated as simple regularization loss.

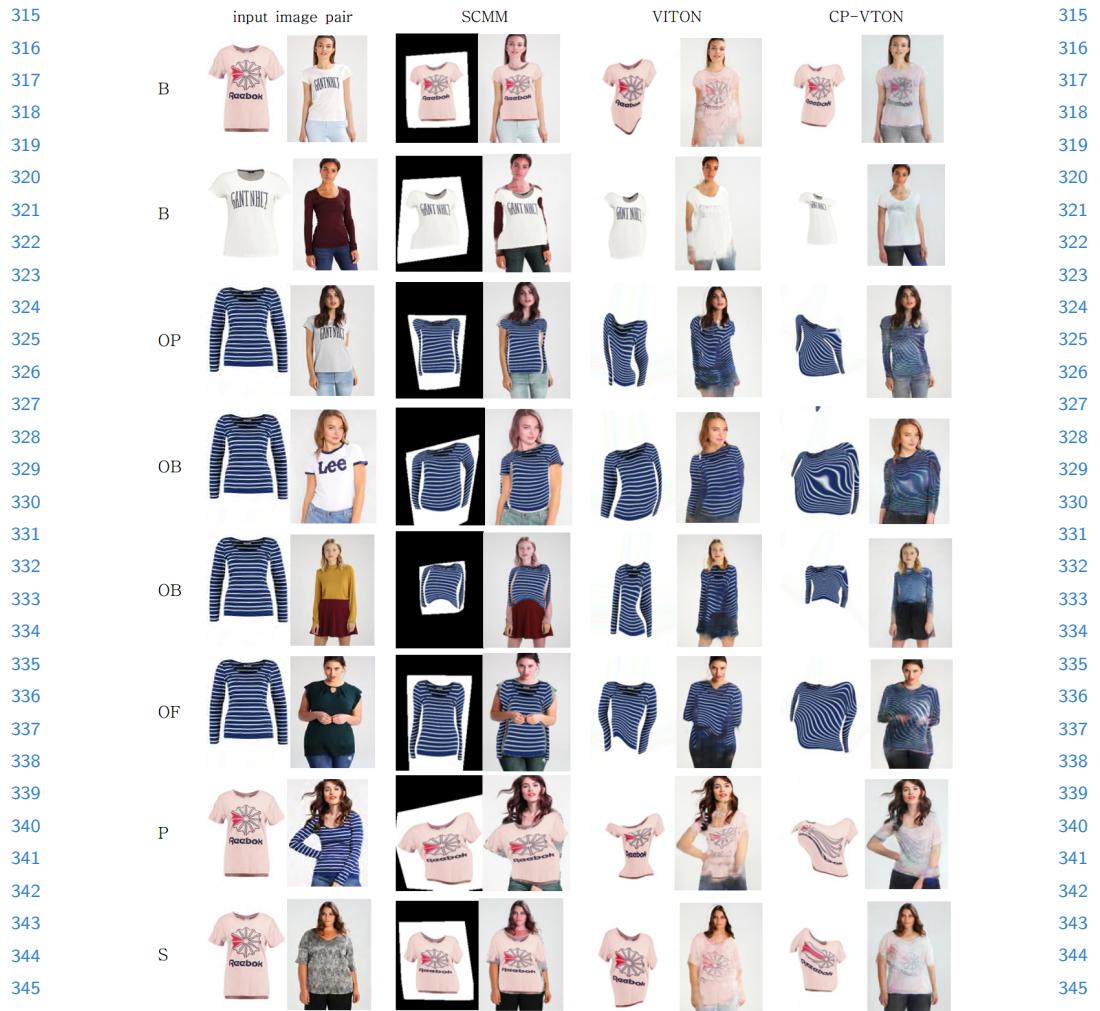
$$L = c_1|I_0 - I_{GT}| + c_2 L_{VGG} + c_3|1 - M_0| \quad (1)$$

270 Fourthly, since no label in the area of warped cloth is the same color as  
 271 background, white colored clothes are confused and improperly processed in the  
 272 blending stage (Fig. 3 (c))

273 Finally, GMM module using Spatial Transform Network[5] with TPS (Thin  
 274 Plate Spline)[2] deformation cannot handle strong 3D deformation due to the  
 275 target pose and also generates artifacts because of the person representation  
 276 inputs. For example, hands-up and folded arms. Note that many errors in the  
 277 warping stage are often hidden in the blending stage when the target clothes are  
 278 single-colored, which can be expected in practical conditions



314 **Fig. 3.** Classified Evaluation of Image-based VTONs (same cloth re-try-on)



**Fig. 4.** Classified Evaluation of Image-based VTONs (new cloth try-on)

### 3 CV-VTON+

#### 3.1 Overview

Fig. 4 illustrates the new pipeline emphasizing the modifications from CP-VTON, which tackles all 4 issues above mentioned except extreme 3D pose of target human.

- Correction on the cloth agnostic representation

- Modification 1: We added a new label ‘Skin’ to the human parsing data to represent the human body shape more accurately (Fig. 2. (a) right).
- Modification 2: We removed the hair label from the Reserved Regions of GMM’s Person Representation input, i.e. only face remains.

– Un-changed human area inclusion

- Modification 3: We added extra human components except the target cloth area, e.g. bottom clothes and legs in the Reserved Regions of Person Representation to TOM, along with face and hair (Fig. 2. (b) right).

– Improving Composition alpha-map

- Modification 4: In the mask loss term in TOM loss function, we replaced the Composition Mask with supervised ground truth mask for a strong alpha mask.

$$L = c_1|I_0 - I_{GT}| + c_2 L_{VGG} + c_1|M_{GT} - M_0| \quad (2)$$

- Modification 5: Lastly, we added the binary mask of warped cloth to TOM network input so that TOM can clearly differentiate the target cloth area regardless of cloth color.

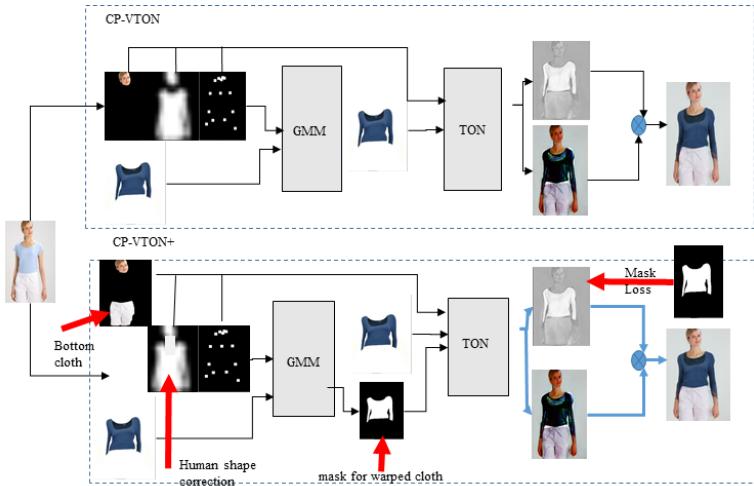
– GMM improvement Will be here: GIC loss and Mask ETC

TODO!!!

## 4 Experiment and Results

### 4.1 Implementation details

We used the same dataset used in CP-VTON, collected first for VITON. We used IoU and SSIM performance metrics for the same cloth retry-on cases for GMM and TOM. For final output quality measures, we used SSIM, Inception Score (IS)[9] and LPIPS[10] for different cloth try-on. The subjective qualities can be examined in Fig. 5/7. Special comments are required for the IoU values, where CP-VTON (0.78) is slightly higher than CP-VTON+ (0.75). The un-expected results originated due to CP-VTON generating as in the current cloth shape. However, similar clothes are not always applicable, furthermore, it generates wrong shaped results for different clothes. Fig. 6 illustrates this with two typical example, plugging and normal tops



**Fig. 5.** Comparison: CP-VTON and CP-VTON+

**Fig. 6.** VTON results

## 4.2 Comparative Results

Final, i.e., TOM results are evaluated with non-reference methods, LPIPS[10], IS[9] and SSIM. Our proposed method, CP-VTON+ outperforms CP-VTON in LPIPS with 0.1263 against 0.1397, in SSIM with 0.8076 over 0.7798, and in IS with 2.76 over 2.7417. The subjective evaluation shows significant visual improvements, especially in cloth textures such as logos and patterns (Fig. 7). We added the test results for all test cases for comparison between CP-VTON and CP-VTON+ in supplementary materials, together with the categorized comparison of SCM-based VTON, VITON and CP-VTON.

## 4.3 Ablation Study

Figure 7 we highlight the impact of the identified problems and improvement of the proposed method step-by-step through the ablation study of CP-VTON+. The first and second columns are target humans and try-on clothes, respectively. The third column is vanilla CVP-VTON results. The fourth column is when unchanged clothes and body parts are added to the reserved region inputs of TOM, retaining the original pants texture. The fifth column is when the mask loss function of TOM is updated with the target cloth area, making the texture and color of cloth sharp and vivid. Finally the sixth and last column is when the body masks are updated, replacing the skin area wrongly labelled as background and hair are removed from the reserved region input of GMM, making GMM can better cloth-and-hair-agnostic human representation.

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  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449

450 WHEN WE GET the GMM improvement, Where we can add this? may be  
 451 the first step?  
 452  
 453



470 **Fig. 7.** Ablation study of CP-VTON+. From left to right column. Target human, try-on  
 471 cloth, CP-VTON, w. human representation, warped cloth mask and mask loss function  
 472 updated, and CP-VTON+

#### 4.4 Known further issues

473  
 474  
 475  
 476 Even though our modification improve the VTON results a lot, and showing  
 477 highly natural results, note that the dataset has limited samples for difficult  
 478 cases, like the long sleeved, complicated shaped, or textured cloth and large  
 479 posture target human. We amplify the key problems identified not try to list all  
 480 the small problems.  
 481

482 As Figure 8 shows two typical failure cases due to the cloth warping. First  
 483 row shows when the arms heavily covers the body area. The warped cloth does  
 484 not match to the human body and TON failed in hiding the warping error. It is  
 485 due to the limitation of STN (Spatial Transform Network). STN is originally de-  
 486 veloped for invert the (augmented) input images for different camera views and  
 487 camera distortion. Non-rigid transforms, including TPS algorithms, cannot han-  
 488 dle the strong 3D deformations of cloth. Also the 3D poses induce self-occlusions.  
 489 The TON network should recognized the cloth area and skin areas, like naked  
 490 arms. One practical short-term solution would be to restrict the pose of target  
 491 human image from the customer. And the long-term solution would be developed  
 492 an 3D cloth deformation techniques for the GMM step, which is under studied  
 493 by the authors.  
 494

The second rows shows the another problems. Even without strong 3D posture of the target human, the warped cloth often shows un-realistic results. The accuracy for matching and warping of STN is not fully studied for VTON applications.

WE need to say something.

The output image quality of all image-based VTON algorithms including ours depend upon the quality of input human representation, i.e., estimated joint locations and parsed human segmentation. The poses of target humans are usually (or forced) rather simple so that the state-of-the-art pose estimation algorithms can provide fairly accurate positions. However the quality of parsed human are not always good enough, especially when the target human wear complicated clothes. Also one can restrict the complexity of human images in pose and cloth style, but still the accuracy at the segmentation boundary sometimes affects the blended image results as shown in third row case in Fig 7, where the pixels of the current top cloth, which is mislabelled as bottom cloth, remained in the blending result. There fore high quality human parsing algorithm especially around boundaries are required.



**Fig. 8.** Cloth Warping Failure of CP-VTON+. From left to right column. Target human, try-on cloth, CP-VTON, and CP-VTON+ (Why not same first and second ??)

## 5 Conclusions

Almost all real computer vision algorithms have a certain condition where they work successfully and not. Therefore it is more important than merely developing better algorithms to identify the working condition of algorithms and approaches.

540 By categorized cloth and human inputs and analysis not only final try-on results  
541 but also intermediate results of the pipeline, e.g. the warped cloth, we showed the  
542 key successful and unsuccessful conditions and origins of a typical image-based  
543 VTON algorithm, CP-VTON.

544 With these identified issues, a CP-VTON+, an improvement to CP-VTON  
545 was proposed, which produces significant quality improvements over existing  
546 state-of-the-art algorithm, CP-VTON.

547 However, there remains several areas that we could not yet solved. The auto-  
548 matic warping to the target human shape is still challenging in feature point  
549 search and matching and limitation of non-rigid 2-D transforms, and the accu-  
550 racy human parsing needs to be improved.

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

## References

1. Belongie, S.J., Malik, J., Puzicha, J.: Shape matching and object recognition using shape contexts. *IEEE Trans. Pattern Anal. Mach. Intell.* **24**(4), 509–522 (2002). <https://doi.org/10.1109/34.993558>, <https://doi.org/10.1109/34.993558>
2. Bookstein, F.L.: Principal warps: Thin-plate splines and the decomposition of deformations. *IEEE Trans. Pattern Anal. Mach. Intell.* **11**, 567–585 (1989)
3. Cao, Z., Martinez, G.H., Simon, T., Wei, S.E., Sheikh, Y.: Openpose: Realtime multi-person 2d pose estimation using part affinity fields. *IEEE transactions on pattern analysis and machine intelligence* (2018)
4. Han, X., Wu, Z., Wu, Z., Yu, R., Davis, L.S.: Viton: An image-based virtual try-on network. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition* pp. 7543–7552 (2017)
5. Jaderberg, M., Simonyan, K., Zisserman, A., Kavukcuoglu, K.: Spatial transformer networks. In: Cortes, C., Lawrence, N.D., Lee, D.D., Sugiyama, M., Garnett, R. (eds.) *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada.* pp. 2017–2025 (2015), <http://papers.nips.cc/paper/5854-spatial-transformer-networks>
6. Jetchev, N., Bergmann, U.: The conditional analogy gan: Swapping fashion articles on people images. In: *Proceedings of the IEEE International Conference on Computer Vision Workshops*. pp. 2287–2292 (2017)
7. Liang, X., Gong, K., Shen, X., Lin, L.: Look into person: Joint body parsing and pose estimation networks and a new benchmark. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **41**, 871–885 (2018)
8. Ma, L., Jia, X., Sun, Q., Schiele, B., Tuytelaars, T., Van Gool, L.: Pose guided person image generation. In: *Advances in Neural Information Processing Systems*. pp. 406–416 (2017)
9. Sun, F.D., Guo, J., Su, Z., Gao, C.Y.: Image-based virtual try-on network with structural coherence. *2019 IEEE International Conference on Image Processing (ICIP)* pp. 519–523 (2019)
10. Wang, B., Zhang, H., Liang, X., Chen, Y., Lin, L., Yang, M.: Toward characteristic-preserving image-based virtual try-on network. In: *ECCV* (2018)