

000

001 3D Reconstruction of Clothes using a Human

002 Body Model and its Application to Image-based

003 VTON

004

005

006

007

008

009

010

000

001 Anonymous ECCV submission

002

003

004

005

006

007

008

009

010

000

001 Paper ID ...

002

003

004

005

006

007

008

009

010

011 **Abstract.** Image-based virtual try-on (VTON) has drawn increasing

012 attraction for on-line apparel shopping mainly because not requiring 3D

013 information of try-on cloths and target humans. However, the existing

014 2D algorithms, even utilizing advanced non-rigid deformation algorithm,

015 could not handle the 3D shape change for the posture of target human.

016 In this study, we propose the 3D cloth reconstruction method using 3D

017 human body model. The 3D model of try-on cloth can be more easily

018 when applied to the rest posed standards human model. Thereafter the

019 pose and shape of cloth can be transferred to the ones of the target hu-

020 mans estimated from an 2D image. Finally the deformed cloth model

021 can be rendered and blended together with unchanged cloth and human

022 parts. The experimental results with a open dataset shows the recon-

023 structed cloth shapes are significantly more natural compared with the

024 2D imaged based deformation result, when the human pose and shape

025 are estimated accurately.

026

027

028

029 **Keywords:** We would like to encourage you to list your keywords within

030 the abstract section

031

032

033

034

035

036

037

038

039

040

041

042

043

044

1 Introduction

031 Online fashion market has been growing rapidly every year. Unlike electronics,

032 which makes it easy to standardize functions and performances, fashion apparel

033 are infinite variations in style, forms, colors, texture, and materials. Also the

034 difference between personal preferences is huge. As a result, clothing purchasing

035 decisions are very difficult to make with current un-customized information, like

036 the cloth and models' try fit images. Therefore, virtual try-on (VTON) is a

037 highly demanding technology for the on-line shopping.

038

039

040

041

042

043

044

The early VTON technologies were based on 3D computer graphics technol-

ogy that uses 3D models for a target human and clothing, which are usually

expensive and difficult to obtain. Therefore, recently 2D image-based VTON

technology are studied in academia and industry, fuelled by recent advance in

computer vision technology based on deep learning (DL). There have been many

assumption in problem settings from the general conditional human image gener-

ation related to VTON application. We consider the one with a try-on cloth and

target human image is a practical condition which is assumed in many papers VITON[3], CP-VTON[4], and the the following [ICCV19, ICIP19]. Therefore 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. In this paper we limit our application to top cloth only due to the restricted data set but consider the bottom, e.g. pants cases would be easier than top cloth cases.

The existing image based algorithms seemingly generate high quality VTON images, but our classified analysis on the cloth style, and human pose and shape reveals significant problems[]. One reason of the seemingly high quality in the existing algorithms are mainly due to the low complexity of dataset, i.e., most cloth are short-sleeved, and mono-colored, and the pose of human are mild. Specifically the results with the long-sleeved cloth arm and body posed shows far low quality of the presented result in their papers. We identified 5 issues in CP-VTON algorithms, some of which are tackled in the following papers. Firstly, the target try-on area is dependent upon current cloth shape. Especially, the neck area pixels are labelled as background and some body area 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 hands are missing in CP-VTON 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.

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

Fourthly, because no label the area of warped cloth in the same color as background, i.e., white are confused and improperly processed in the blending stage (Fig. 3 (c)) Finally, GMM module using Spatial Transform Network with TPS (Thin Plate Spline) deformation cannot handle strong 3-D deformation due to the target pose and also generates artifacts because of the person representation inputs. For examples, hands-up and folded arms. Note that many errors in the warping stage are often hidden in the blending stage when the cloth are single-colored, which can be expected in practical conditions (Fig. 3 (d)).

In this paper, we focus on the last but most difficult problems that can be solved in pure 2-D image based algorithm. The 3D cloth deformation is inherently difficult for 2D warping method, including non-rigid one, like TPS algorithm, we propose to first reconstruct 3D model of try-on cloth, then apply the the pose and shape transfer for the target human, and finally blending with unchanged image contents like the face, bottom cloth, and background. Therefore the one of main task now is to reconstruct 3D cloth model from 2-D try-on cloth image. The 3D cloth model reconstruction have been studied in previous studies [] but still needs significant improvement for general condition. Our key idea in this step is that once we can control the human pose and shape to become similar to

090 the try-on cloth's, the 3D reconstruction process can be made much easier and
091 the reconstruction quality would be much higher than general pose and shape
092 condition.

093 So in the Section 3, we describe the 3D cloth reconstruction algorithm. we
094 divide the reconstruction step into 2D matching of cloth to the standard body
095 silhouette and 3D reconstruction of cloth. The later 3D reconstruction step is
096 done through the SMPLify algorithm for the SMPL 3D body model. In Section
097 4, the blending method described, where the 3D cloth model are transferred to
098 the target human images, through SMPL body parameters of shapes and poses.
099 Then the transferred 3D is rendered and blended to the target human image. In
100 this step we reused the 2D VTON blending algorithm with the modification for
101 the condition. The sampled results from dataset are presented in Section 5 and
102 the paper is concluded in Section 6. In addition to our main study, we added the
103 classified quality evaluation of the previous 2D image based VTON algorithms
104 for the completeness of the paper.

105 2 Classified Image Based Performance Evaluation

106 2.1 Image-based VTON

107 In this Section, we started with evaluating the 2D image based VTON algo-
108 rithms. We considered CP-VTON published in 2018 as the benchmarking algo-
109 rithm. The previous and following in 2019 share same input image and informa-
110 tion conditions with CP-VTON and compare the results with it. Here we include
111 the SCM based-VTON, VITON, and CP-VTON, but believe the performance
112 strength and weakness are similar in the other algorithms too.

113 The Image based VTON algorithms are mostly composed of two stages: (1)
114 cloth warping step that warps the try-on cloth to align with the pose and shape of
115 the target model (called GMM in CP-VTON: geometric Manipulation Module),
116 and (2) blending step that blends the warped cloth onto the target human image
117 (called TON in CP-VTON: Try-On Network). CP-VTON assumes the target
118 human image is pre-processed for a cloth agnostic human representation by a
119 human pose estimation like OpenPose [] and human parsing like LIP[]. The
120 human representation is composed of 1) heat maps for each joints 2) silhouette of
121 human body, and 3) face and skin pixels patches (non-cloth and human identity
122 area). We use the same dataset collected by Han et al. used in VITON and
123 CP-TON papers.

124 2.2 Classified quality

125 Even though the success and failure cases are presented and compared with other
126 algorithms' results, the failure case analysis is not enough for understanding the
127 origin of failure cases and therefore difficult to find the solution for them. A
128 classified evaluation would be better for this understanding. Here we summarized
129 the classified results from our another study. We classify input try-on cloth and

target human images according to the posture and body type of the person, the degree of occlusion of the clothes, and the characteristics of the clothes. Quality 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 not include here for limitation of space, and the same cloth cases are enough to explain the tendency of performance. Though in general CP-VTON generates the best quality image, the relative comparison is not the main purpose of the analysis.

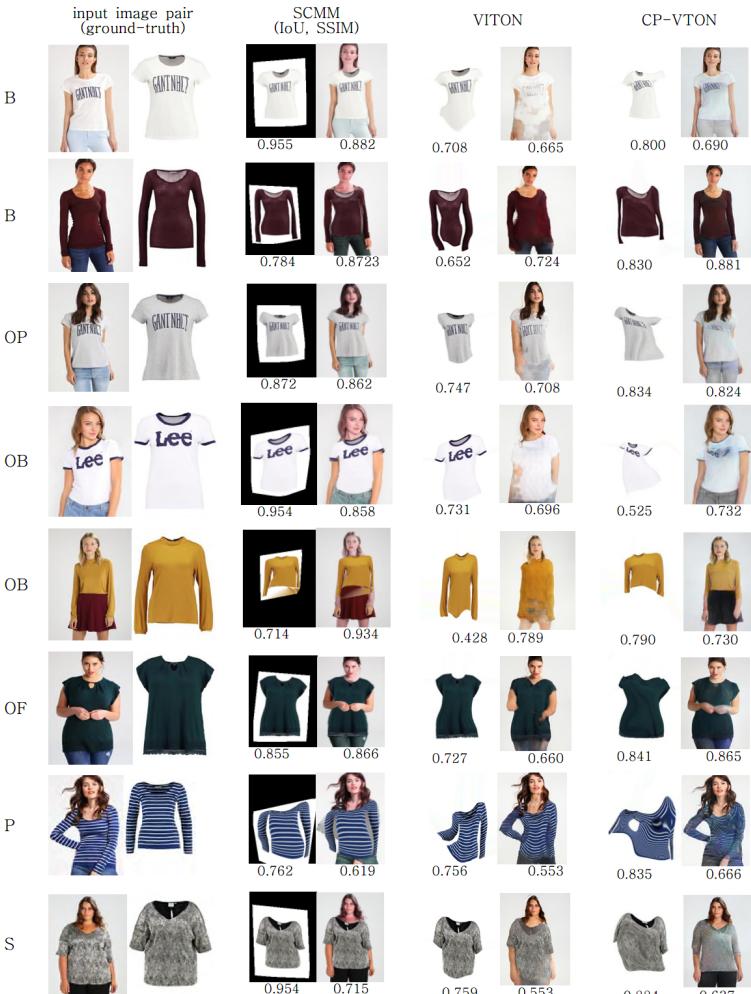


Fig. 1. Classified VTON result: same clothes

180 ////

181 Here I will describe the evaluation discussion, finally commenting that the
182 GMM has serious problem.

183 ////

184 Especially note that the warped cloth are often too much different for desired
185 shape. It is originated two facts. First the 3D deformation that any 2D defor-
186 mation including non-rigid transform such as TPS is quite limited, especially
187 any 2D deformation cannot handle when the two area in the original image are
188 overlapped in the destination images. There for when the arms of long sleeved
189 cloth occlude the main body, 2D warping cannot approximate the 3D defor-
190 mation properly. Second, the deformation needs corresponding points between
191 the source nd target image. The cloth are extremely difficult object to find the
192 corresponding points. The STN (spatial transform network) and SCM (shape
193 context matching) cannot find the corresponding points when the target cloth
194 and original cloth has different shapes. In conclusion, the 2D image based al-
195 gorithm has serious limitation in the range of applications. It can apply to the
196 mild posed target human only and simple short sleeved cloth, mainly because
197 the inherent limitation of 2D deformation method including non-rigid ones, and
198 the poor performance of matching algorithm. To overcome this limitation, we
199 consider to model the try-on cloth into 3D model and apply the 3D deformation

200 201 3 3D model reconstruction of cloth

202 3.1 Overview

203 For 3D human body model, we use Skinned Multi-Person Linear model (SMPL),
204 because SMPL has well defined control variable for shape and pose and also
205 well defined parameter estimation algorithms. For similar reasons, SMPL have
206 been utilized in many research works. Furthermore because it is based on blend
207 skinning, SMPL is compatible with existing rendering engines and we make
208 it available for research purposes. SMPL is a skinned vertex-based model that
209 accurately represents a wide variety of body shapes in natural human poses. The
210 parameters of the model are learned from data including the rest pose template,
211 blend weights, pose-dependent blend shapes, identity-dependent blend shapes,
212 and a regressor from vertices to joint locations. Unlike previous models, the
213 pose-dependent blend shapes are a linear function of the elements of the pose
214 rotation matrices. This simple formulation enables training the entire model from
215 a relatively large number of aligned 3D meshes of different people in different
216 poses.

217 For estimating the SMPL parameters, we use SMPLify method in this study.
218 However any other methods can be used because we assume nothing on the pro-
219 cedure and use estimated parameters only. SMPLify use 2D human body joint
220 information often obtained from deep learning based method like DeepCut or
221 OpenPose, and minimize the projected joint locations and the given (considered
222 true) 2D joint locations. The cost function can include other priors and silhouette

information. We made minor optimization for half body dataset, such as joint location mapping between the joints of used fashion data set and SMPLify joint definition, and conditional inclusion of invisible joints and initialization step. From our experiments with all 2032 test images, we found that the SMPLify quality should be much improved for fully automatic application to VTON application. So the result included in this paper excluded the bad matching cases which is around 30% of all test images.

Clothed human reconstruction using SMPL have been studied in several previous works [PhotoWake up and Rummians....]. Even we are successful in modelling human body, there are further difficulty to recover the clothed human model from body model. It is because the cloth vertices are not directly corresponds to the human body's, and even though it has it is still difficult to estimate the difference between two. Also the texture of cloth can be occluded by other part of cloth and human body parts. The previous work try to solve the problem in the given image condition. Therefore the results are strongly dependent upon the input image. In this paper, we make this step easy using simple standard human pose, where all frontal part of cloth is well separated and visible. This setup cannot handle all problem in the clothed human model reconstruction but can greatly make it easy. The following subsection describe the procedure in details.

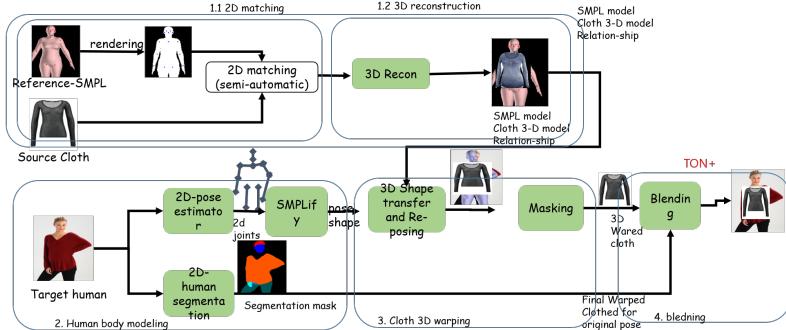


Fig. 2. Pipeline

3.2 2D Standard Cloth matching

To align the try-on cloth image with 3D SMPL body model, first their dimension spaces should be matched. Natural way would be first rendering the SMPL body model into 2D image space. However, again the matching with cloth image and body silhouette is not a simple task, for simplicity we assume we can segment the silhouette so that the the remaining area can be easily matched by SCM algorithms. We argue that this step can be monitored by service provider which

is practically acceptable; the manual operation from the customer in the try-on step would be not acceptable in general service environment.

$$(I_{c,warped}, M_{c,warped}) = T_{SMPL}((I_c, M_c)) \quad (2)$$

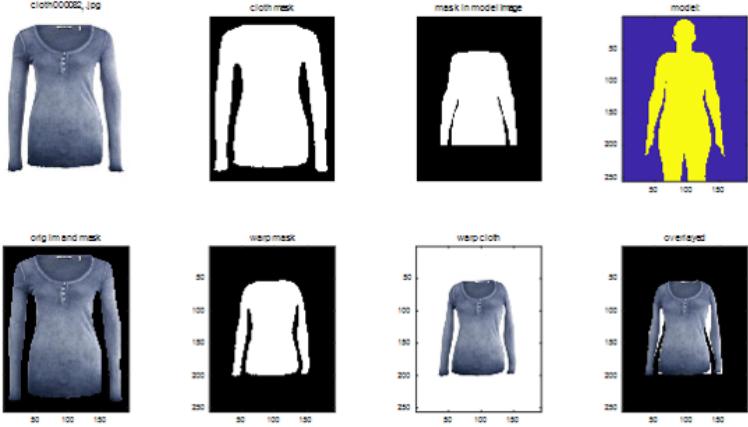


Fig. 3. 2D Matching

3.3 3D cloth model reconstruction

The 3D reconstruction process from aligned cloth image and projected silhouette consists of 2 steps. First, the vertices of 3D body mesh are projected into 2D image space, the boundary vertices in 2D spaces and the cloth boundaries are used for corresponding points. The corresponding points in the cloth boundary defined the closest points from the projected vertices. This step works well in our cases differently from PhotoWakeUp study, because the part of body and cloth are not self-overlapped. This is a implementation benefits of our approach. From the corresponding point pairs, a TPS parameter are estimated and applied to the mesh points. The new mesh points are considered the vertices projected from 3D mesh of cloth.

From the 2D points to 3D points are done with inverse projection with depth obtained from the body with a small constant gap. In reality the gap between the cloth and body cannot be constant but it works with tight or simple clothes. Further research should be needed for accurate depth estimation.

$$V_{clothed} = Pjt^{-1}(T((Pjt(V_{body})), depth(V_{body}))) \quad (3)$$

The try-on cloth images are used for the texture for the 3D cloth mesh. We can filter the vertices corresponding to cloth and get the cloth 3D mesh model. Figure xxxxx shows the reconstructed cloth examples.

Discussions needed

Fig. 4. 3D reconstructed cloth

4 Transfer of 3D cloth model to the target Human and Virtual Try-On

4.1 Transfer of 3D cloth model to the target Human

The 3D model and texture information obtained above are for the standard shape and posed person. To apply this information to the target human image, we have to apply the shape and pose parameters of estimated from SMPLify step. In stead of apply the shape and pose parameters to the obtained clothed 3D model, we transferred the displacement of cloth vertices to the target human body model, because the application of new parameters to the Body model provide much natural results.

Multiple option can be considered for the transferring. We could transfer the physical size of cloth or keep the fit, i.e., keep the displacement from the body to cloth vertex as before. We simply decide the Fit-preserving option for showing more natural results for final fitting.

Technically the displacement should be calculated locally. First we calculated the local coordinate at each vertices. We defines the local coordinates: surface normal vector as z -axis, and the vector to smallest indexed edge as x-axis, and their cross product vector as y-axis as the following equations.

$$u_z = \text{normal}(V_{body}) \quad (4)$$

$$u_x = u''_x / |u''_x|, u''_x = u'_x - u'_x \cdot u_z, u'_x = (V_{\text{argmin}(N_V)} - V_{body}) \quad (5)$$

$$u_y = u_z \otimes u_x, \quad (6)$$

where N_V is the neighbor vertex of V .

The displacement is expressed in the local coordinates and then used the same way in the new target body surfaces for location transfer

$$\vec{d} = (d_x, d_y, d_z) = V_{clothed} - V_{body} \text{ in } (u_x, u_y, u_z | V_{body}) \quad (7)$$

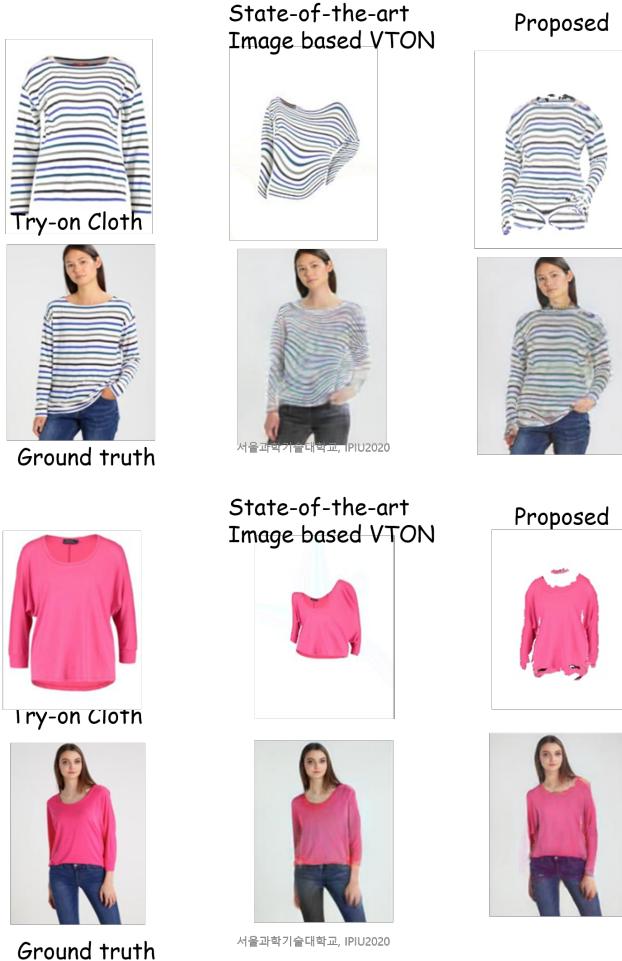
$$V'_{clothed} = V'_{body} + \vec{d} \text{ in } (u_x, u_y, u_z | V'_{body}) \quad (8)$$

360 4.2 Blending of warped cloth with target human image

361
362
363 This part is under implementation.

364 We first tried to use TOM. But we found when the reconstruction is not
365 perfect the blending is not natural

366 Other option is first reconstruct all the clothed information from the target
367 user and overlay the transferred cloth.



402 **Fig. 5.** VTON results

403
404

405 5 Conclusions

406
407 In this paper, we proposed 3D cloth model reconstruction method using single
408 cloth image. Leveraging the 3D body model, we can make it easy to reconstruct
409 3D shape information. The 3D cloth model is used for transferring the cloth
410 to target human model. The transferred clothed can be integrated with the
411 human image contents for realizing the pose and shape changes which can not
412 be realizable by existing image based VTON methods.

413 However, the algorithms in each step of the pipeline are not perfect and
414 has many thing to improve at present. Especially the in-accuracy in estimating
415 human pose and shape make the integrated VTON results is not natural enough.
416 Therefore we can consider to improve the SMPLify algorithm or use different
417 blending step that suit for the 3D model input.

418 419 References

- 420 1. M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black. *SMPL: A skinned multi-person linear model*. ACM transactions on graphics (TOG), Vol. 34, No. 6, 248, 2016
- 421 2. F. Bogo, A. Kanazawa, C. Lassner, P. Gehler, J. Romero, and M. J. Black. *Keep it SMPL: Automatic estimation of 3D human pose and shape from a single image*. In European Conference on Computer Vision (ECCV 2016) pp. 561-578, Oct. 2016.
- 422 3. B. Wang, H. Zheng, X. Liang, Y. Chen, L. Lin, and M. Yang, M. *Toward characteristic-preserving image-based virtual try-on network*. Proc. of the European Conference on Computer Vision, pp. 589-604, 2018
- 423 4. C. Y. Weng,, B. Curless, and I. Kemelmacher-Shlizerman *Photo wake-up: 3d character animation from a single photo*. In Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) pp. 5908-5917, July, 2019.
- 424 5. Belongie, S., Malik, J., Puzicha, J. *Shape matching and object recognition using shape contexts*. IEEE Transactions on PAMI, 25(4), 509-522. 2002
- 425 6. Cao, Z., Simon, T., Wei, S. E., Sheikh, Y. *Realtime multi-person 2d pose estimation using part affinity fields* Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 7291-7299.
- 426 7. Han, X., Wu, Z., Wu, Z., Yu, R., Davis, L. S. *Viton: An image-based virtual try-on network*. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 7543-7552, 2018
- 427 8. Liang, X., Gong, K., Shen, X., Lin, L. *Look into person: Joint body parsing and pose estimation network and a new benchmark*. IEEE transactions on PAMI, 41(4), 871-885, 2018
- 428 9. Wang, B., Zheng, H., Liang, X., Chen, Y., Lin, L., and Yang, M. *Toward characteristic-preserving image-based virtual try-on network*. Proceedings of the European Conference on Computer Vision, pp. 589-604. (2018)
- 429 10. Zanfir, M., Popa, A. I., Zanfir, A., and Sminchisescu, C. *Human appearance transfer* Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 5391-5399, 2018

450 Page 11 of the manuscript.

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495 Page 12 of the manuscript.

496 This is the last page of the manuscript.

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

Now we have reached the maximum size of the ECCV 2020 submission (excluding references). References should start immediately after the main text, but can continue on p.15 if needed.

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539