

Supplementary Materials

1. User Study

The user study evaluates both the realism and accuracy of the generated image conditioned on correct layout. Because we want only aim to evaluate the image generator, we use the layout produced by the pre-trained human parser on the ground truth image (rather than the ones produced by the semantic layout generator). Because produced layout are often noisy, we choose examples where the layout is good, giving a fair representation of the top 20 percentile of our results.

We ran our user study with two populations of participants (vision researchers tested with Form B, and randomly selected people tested with Form A). Images are displayed using the highest possible resolution (512x512) and each participant is primed with two real and fake examples before the study. Each participant is then tested with 50 examples (25 real and 25 fake), without repeating products. During the study, the garment image and the model image are shown side-by-side, giving subject an easy way to determine whether the synthesized image accurately represent the garment – an important property for fashion e-commerce application.

For every model, we tested swapping a single garment. Each model in our dataset has a ground truth paired with only one garment (the other garments worn by the model are not matched with any product images). Both the real and fake images are shown with the same outfit. Form A and Form B mirror each other (i.e. if a garment is shown as the generated version in Form B, the real image will be shown in Form A).

The raw questions and responses are provided under the folder “user study”.

2. Semantic Layout and Pose Representation

2.1. Semantic Layout

We obtain our semantic layout using an off-the-shelf human parser [3]. Our Semantic Layout has 15 semantic labels to represent different human parts and garments. These are background, hair, face, neckline, right arm, left arm, right shoe, left shoe, right leg, left leg, bottoms, full-body, tops, outerwear, and bags. Among these semantic labels, bottoms, full-body, tops, and outerwear are garment labels. The semantic segmentation mask is $m \in R^{H \times W \times 15}$, where H and W correspond to the width and height of the image.

2.2. Incomplete Layout

During training, we create the incomplete layout m_i by hiding the target garment labels and relevant skin labels (by setting these labels to the background class). For tops and



Figure 1. A version of the Figure 1 without bounding boxes.

outerwear, we hide tops, outerwear, left arm, right arm and neckline; for bottoms, we hide bottoms, left leg and right leg; for full-body, we hide full-body, left leg, right leg, left arm, right arm and neckline. All original channels are still outputted as the incomplete layout ($m_i \in R^{H \times W \times 15}$).

2.3. Pose Representation

We first apply the pre-trained Openpose [1, 6] model on the model images to obtain 18 key points. Following prior work [5, 2, 7], we convert the key points into a heatmap with 18 channels. Our key point heatmap becomes $p \in R^{H \times W \times 18}$. In each channel, we draw a 5×5 square centered at each keypoint’s coordinate and set all values in the square to one. If a key point is missing, the channel will be set to all zeros.

3. Experiment Setups

3.1. Network Architecture

For both the semantic layout generator G_{layout} and the inpainting network of the multi-warp garment generator $G_{garment}$, we use a U-Net of 5 hidden layers. The channel sizes for the hidden layers are 64, 128, 256, 512, and 1,024 respectively. We downsample the image size by 2 at each layer, using bilinear interpolation.

For the warper module of the multi-warp garment generator $G_{garment}$, we use a ResNet-18 model with ImageNet pre-trained weights as the backbone.

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053
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

108	3.2. Training Procedure	162
109		163
110	We train our network using Adam Optimizer with a learning rate of $1e^{-4}$ for the semantic layout generator G_{layout} and $2e^{-4}$ for the multi-warp garment generator $G_{garment}$. Both networks are trained on a Quadro RTX 6000 GPUs (24GB). G_{layout} is trained for 50k steps with a batch size of 16. $G_{garment}$ is trained for 100k steps with a batch size of 8.	164
111		165
112		166
113		167
114		168
115		169
116		170
117	For training the G_{layout} , λ_1 and λ_2 are set to 1 and 0.2 respectively. For training the $G_{garment}$, γ_1 , γ_2 , γ_3 and γ_4 are set to 5, 5, 3 and 1, respectively.	171
118		172
119		173
120		174
121	4. More Qualitative Comparisons	175
122		176
123	We show more qualitative comparisons between our method and O-VITON [4]. Notice O-VITON [4] is the only prior work that supports multi-garment try-on, but they did not release their dataset or implementation. For a fair comparison, we found garment images that most closely resemble the garments chosen in [4] in terms of style, color, and texture. Image results for O-VITON are directly taken from their paper and supplementary materials. Notice the substantial improvement in generation quality in Figure 2.	177
124		178
125		179
126		180
127		181
128		182
129		183
130		184
131		185
132	References	186
133		187
134	[1] Z. Cao, G. Hidalgo Martinez, T. Simon, S. Wei, and Y. A. Sheikh. Openpose: Realtime multi-person 2d pose estimation using part affinity fields. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 2019. 1	188
135		189
136		190
137		191
138	[2] Xintong Han, Xiaojun Hu, Weilin Huang, and Matthew R. Scott. Clothflow: A flow-based model for clothed person generation. In <i>ICCV</i> , 2019. 1	192
139		193
140		194
141	[3] Peike Li, Yunqiu Xu, Yunchao Wei, and Yi Yang. Self-correction for human parsing. <i>arXiv preprint arXiv:1910.09777</i> , 2019. 1	195
142		196
143	[4] Assaf Neuberger, Eran Borenstein, Bar Hilleli, Eduard Oks, and Sharon Alpert. Image based virtual try-on network from unpaired data. In <i>CVPR</i> , 2020. 2, 3	197
144		198
145		199
146	[5] Bochao Wang, Huabin Zheng, Xiaodan Liang, Yimin Chen, and Liang Lin. Toward characteristic-preserving image-based virtual try-on network. In <i>Proceedings of the European Conference on Computer Vision (ECCV)</i> , 2018. 1	200
147		201
148		202
149		203
150	[6] Shih-En Wei, Varun Ramakrishna, Takeo Kanade, and Yaser Sheikh. Convolutional pose machines. In <i>CVPR</i> , 2016. 1	204
151		205
152	[7] Han Yang, Ruimao Zhang, Xiaobao Guo, Wei Liu, Wangmeng Zuo, and Ping Luo. Towards photo-realistic virtual try-on by adaptively generating↔preserving image content. In <i>CVPR</i> , 2020. 1	206
153		207
154		208
155		209
156		210
157		211
158		212
159		213
160		214
161		215



Figure 2. Qualitative comparison with O-VITON [4]. The top two rows in each cell show the garments in the outfit and the bottom row in each cell shows generated try-on results.

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

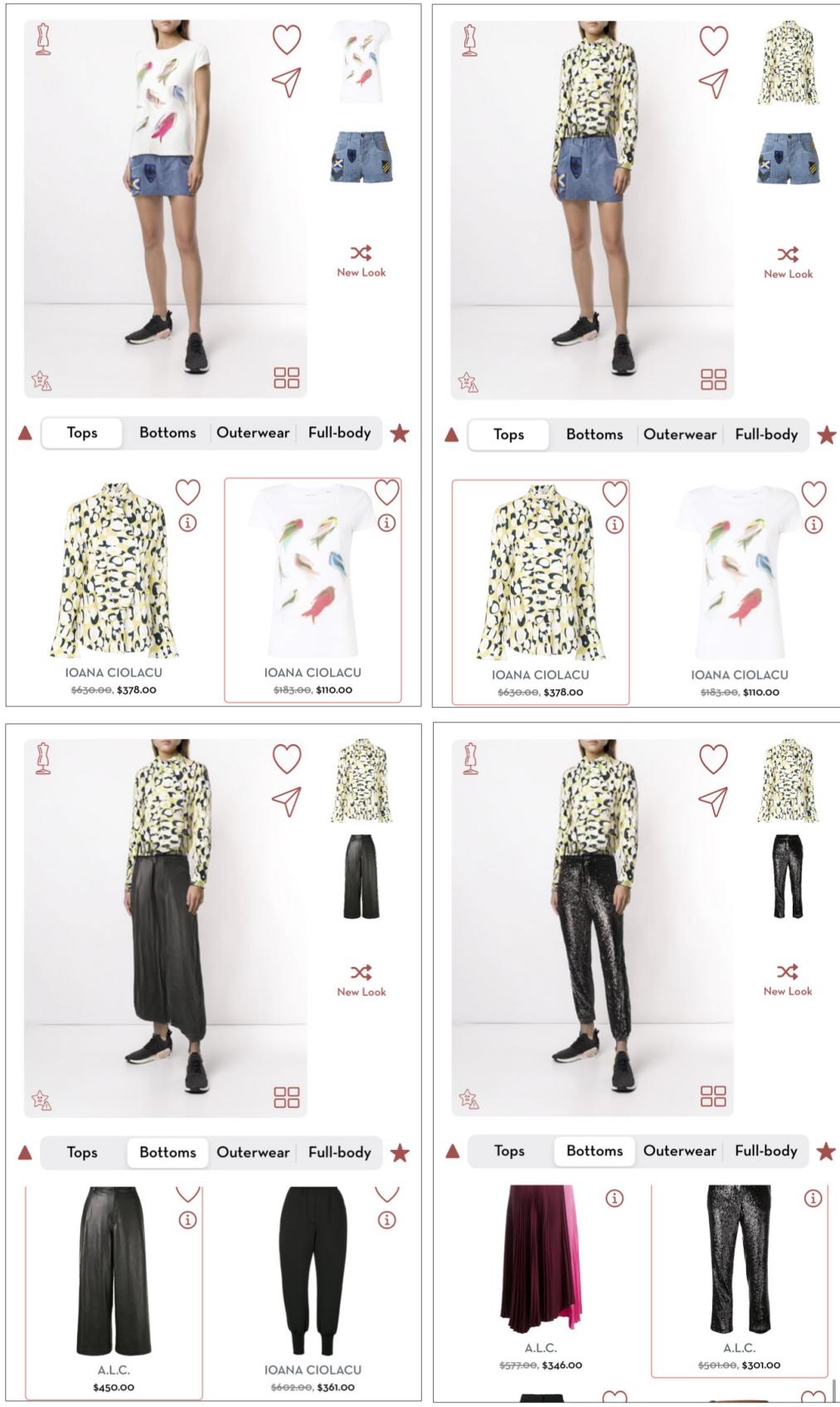


Figure 3. The figure shows examples of an interactive interface powered by our method. A user can select a garment listed below and our method will produce a visualization of the model wearing the selected garment in real-time. Garment images are products listed on an e-commerce site.

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

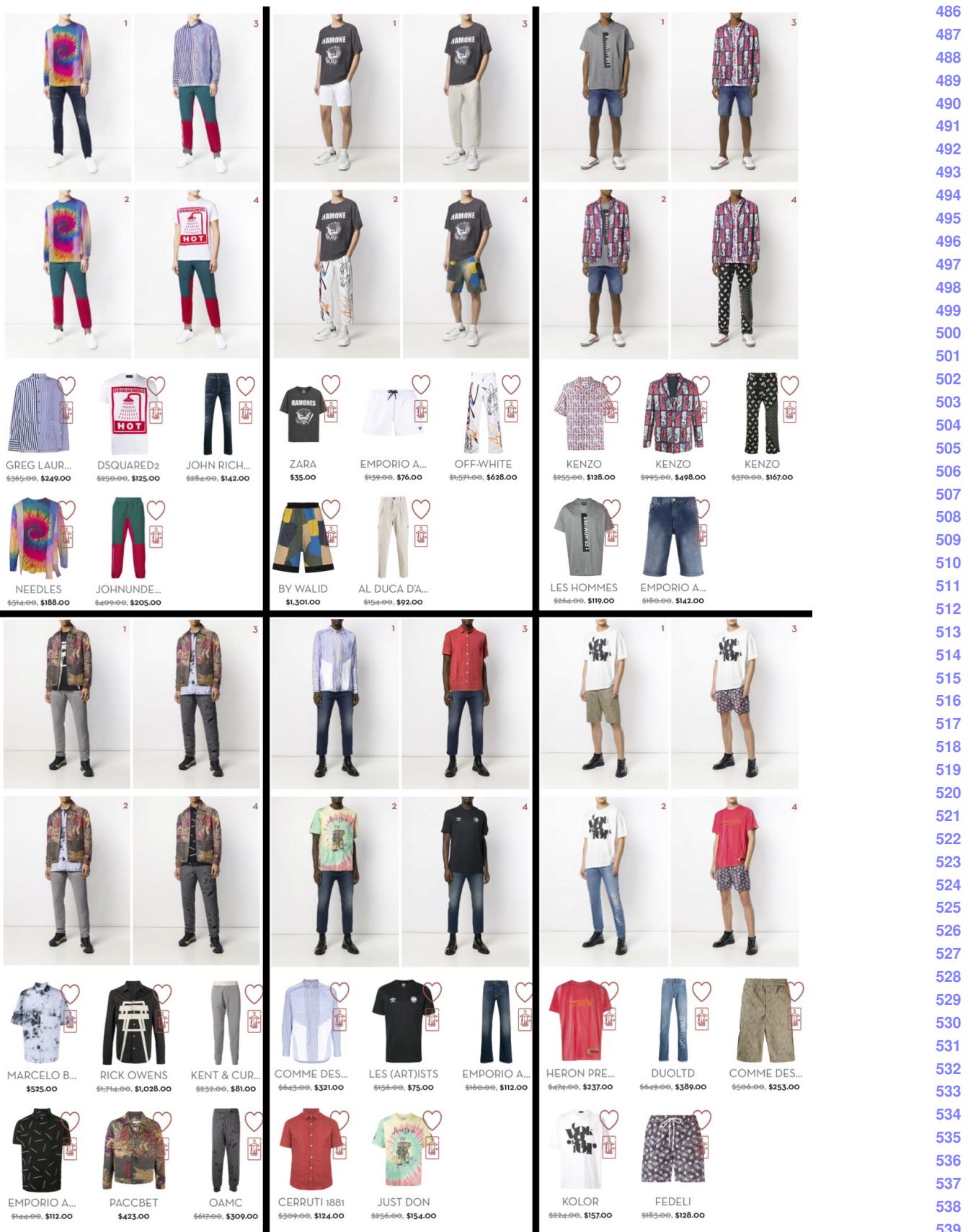


Figure 4. The figure shows outfits curated by users through the interactive interface (male models).

CVPR 2021 Submission #6278. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

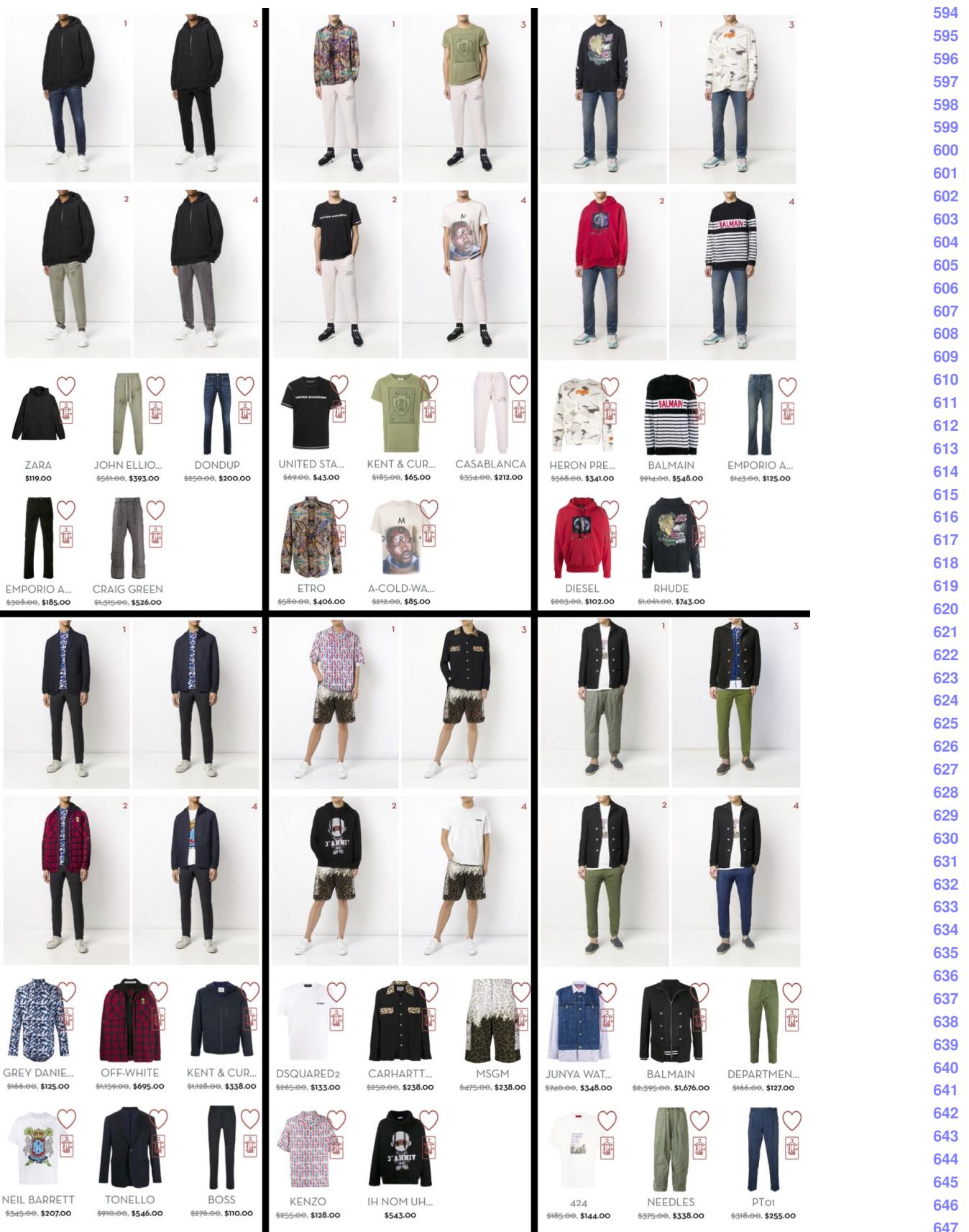


Figure 5. The figure shows outfits curated by users through the interactive interface (male models).

CVPR 2021 Submission #6278. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

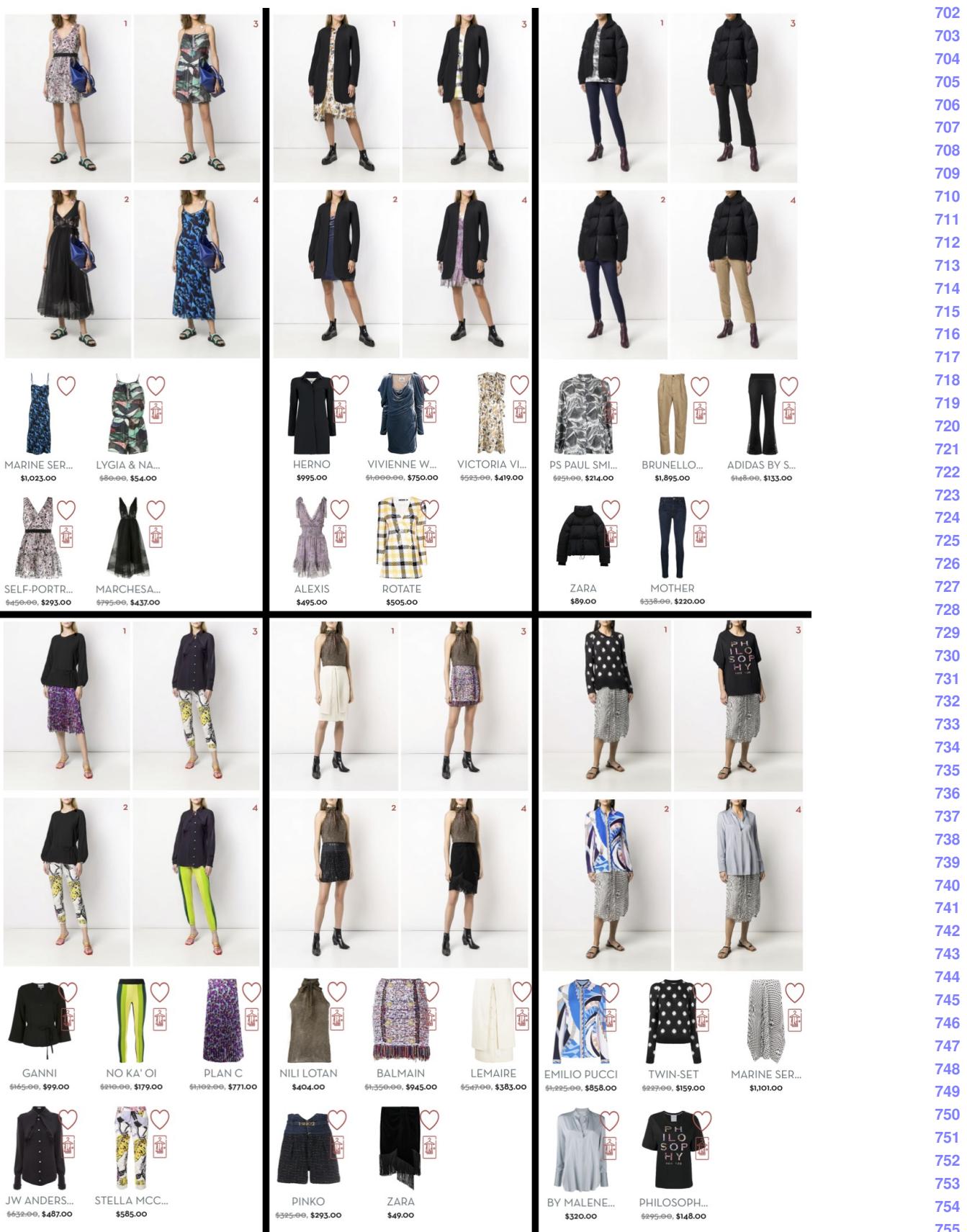


Figure 6. The figure shows outfits curated by users through the interactive interface (female models).

CVPR 2021 Submission #6278. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

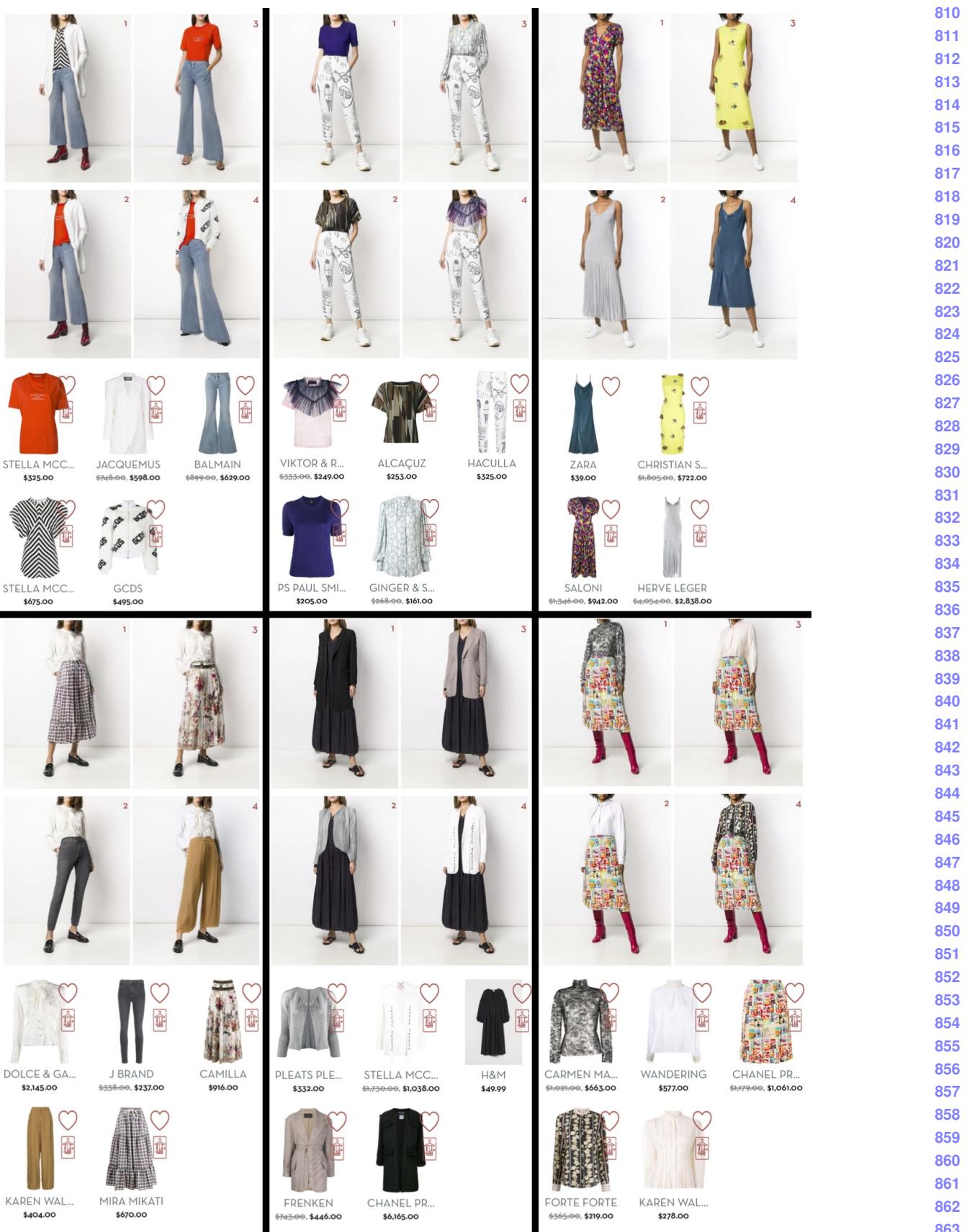


Figure 7. The figure shows outfits curated by users through the interactive interface (female models).