

APPENDIX

Anonymous authors

Paper under double-blind review

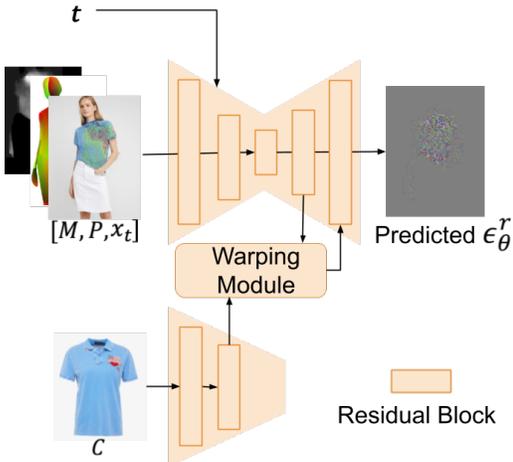


Figure 8: Architecture of our UNet in EARSB.

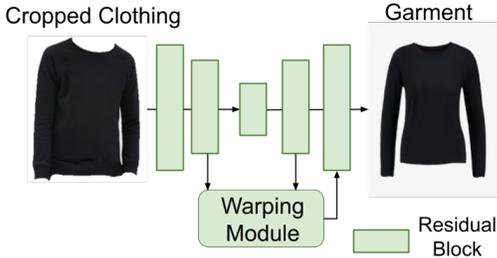


Figure 9: Architecture of our UNet in the human-to-garment model.

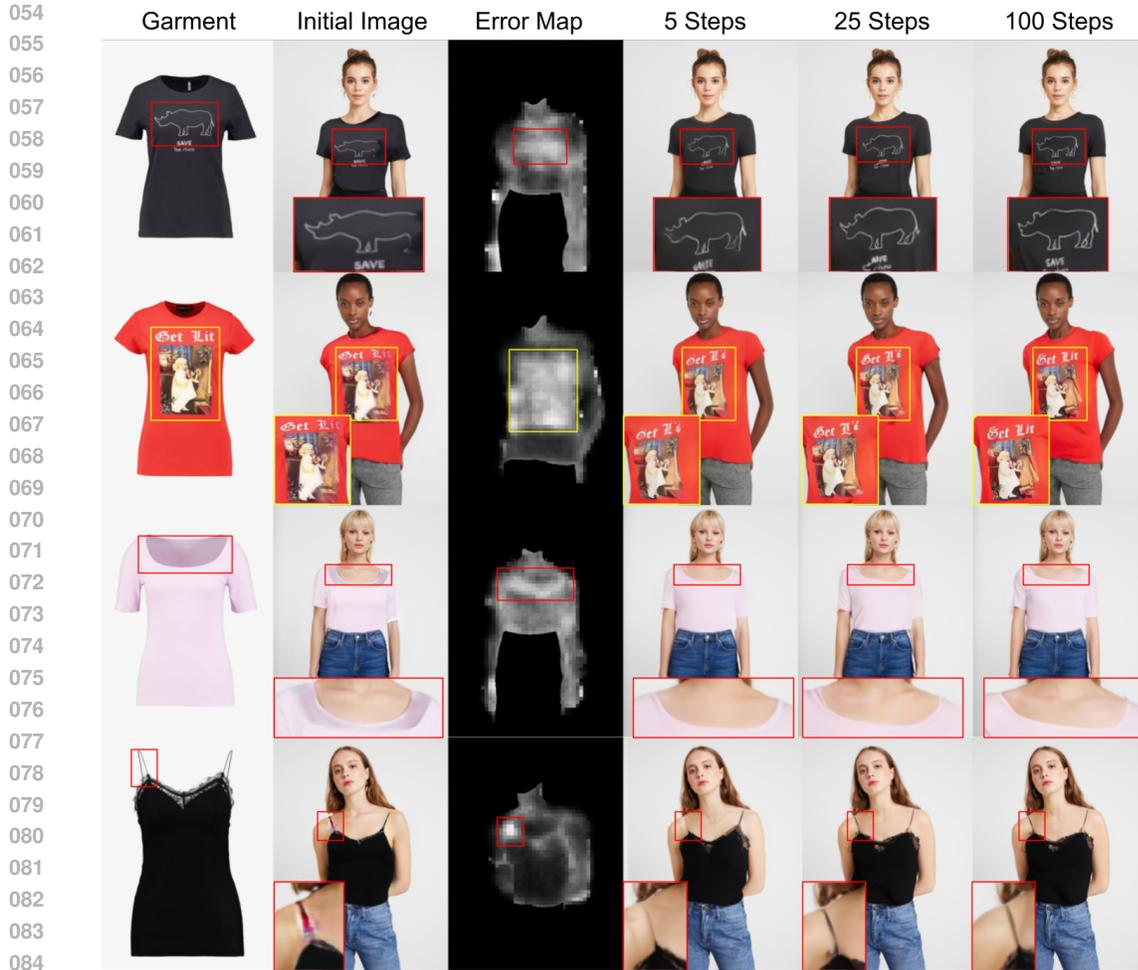
A IMPLEMENTATIONS DETAILS

For generating the initial image x_1 in our EARSB training, we employ two try-on GAN models: HR-VTON (Lee et al., 2022) and SD-VTON (Shim et al., 2024). All human images are processed to maintain their aspect ratio, with the longer side resized to 512 pixels and the shorter side padded with white pixels to reach 512. During training, images undergo random shifting and flipping with a 0.2 probability. The weakly-supervised classifier is trained for 100K iterations with a batch size of 8, while the human-to-garment GAN is trained for 90K iterations with a batch size of 16. EARSB+H2G-UH/FH is trained for 300K iterations with a batch size of 32, incorporating 15% synthetic pairs in each batch. The first 200K iterations are trained on $t \in [0, 1]$ while the following 100k iterations are finetuned on $t \in [0, 0.5]$ and $t \in [0.5, 1]$ respectively following (Balaji et al., 2022). All models utilize the AdamW optimizer with a learning rate of 10^{-4} .

For inference, we select the GAN model that demonstrates better performance on each dataset to generate the initial image. Specifically, we employ GP-VTON (Xie et al., 2023) for VITON-HD and SD-VTON (Shim et al., 2024) for DressCode-Upper. During the sampling process, the guidance score in Eq. (10) is scaled by a factor of 6 and clamped to the range $[-0.3, 0.3]$.

B UNET ARCHITECTURE

EARSB UNet. The UNet architecture in EARSB consists of residual blocks and garment warping modules. It processes the concatenation of the error map M , pose representation P , and noisy image x_t to predict the noise distribution ϵ_θ^r at time t . The UNet encoder has 21 residual blocks, with the number of channels doubling every three blocks to a maximum of 256. Similarly, the garment encoder has 21 residual blocks but reaches a maximum of 128 channels. The decoder mirrors the encoder’s structure, with extra garment warping modules. As shown in Fig. 8, each of the first 15 residual blocks in the UNet decoder is followed by a convolutional warping module. These modules concatenate encoded garment features and UNet-decoded features to predict a flow-like map for spatially warping the encoded garment features. The warped features are then injected into the



085
086
087
088

Figure 10: Results on different time steps. Our error map focuses on low-quality regions and maintains the quality of the sufficiently good regions.

089
090

subsequent decoder layer via input concatenation. Following (Rombach et al., 2022), all residual blocks and flow-learning modules incorporate timestep embeddings to renormalize latent features.

091
092
093
094
095
096

Human-to-Garment UNet. Our human-to-garment UNet architecture is adapted from the model proposed in (Han et al., 2019). As illustrated in Fig. 9, it shares similarities with the UNet in EARSB, but with two key distinctions: a) It is not timestep-dependent and takes cropped clothing as input to generate its product-view image. b) The garment warping module utilizes the i_{th} clothing features from both the encoder and decoder to learn a flow-like map, rather than using encoded features from the human.

097 098 099 C VISUALIZING ERROR MAPS

100
101
102
103
104
105
106
107

Our EARSB focuses on fixing specific errors and therefore can save the sampling cost when initial predictions are sufficiently good. For example, in the first row of Fig. 10, the error map highlights the graphics and text in the initial image. The highlighted low-quality part is being refined progressively as the number of sampling steps increases from 5 to 100. At the same time, other parts that our weakly-supervised classifier believes to be sufficiently good, which are mostly the solid-color areas, are kept well regardless of the number of sampling steps. Therefore, for an initial image whose error map has almost zero values, we can choose to use fewer steps in sampling. On the contrary, for an initial image whose error map has high confidence, we should assign more sampling steps to it to improve the image quality.

	HR-VTON (Lee et al., 2022)	SD-VTON (Shim et al., 2024)	GP-VTON (Xie et al., 2023)
Baseline	10.75	9.05	8.61
CAT-DM (Zeng et al., 2024)	10.03	8.76	8.55
EARSB	9.11	8.69	8.42

Table 4: FID scores of using different try-on GAN models to generate the initial image under the unpaired setting.

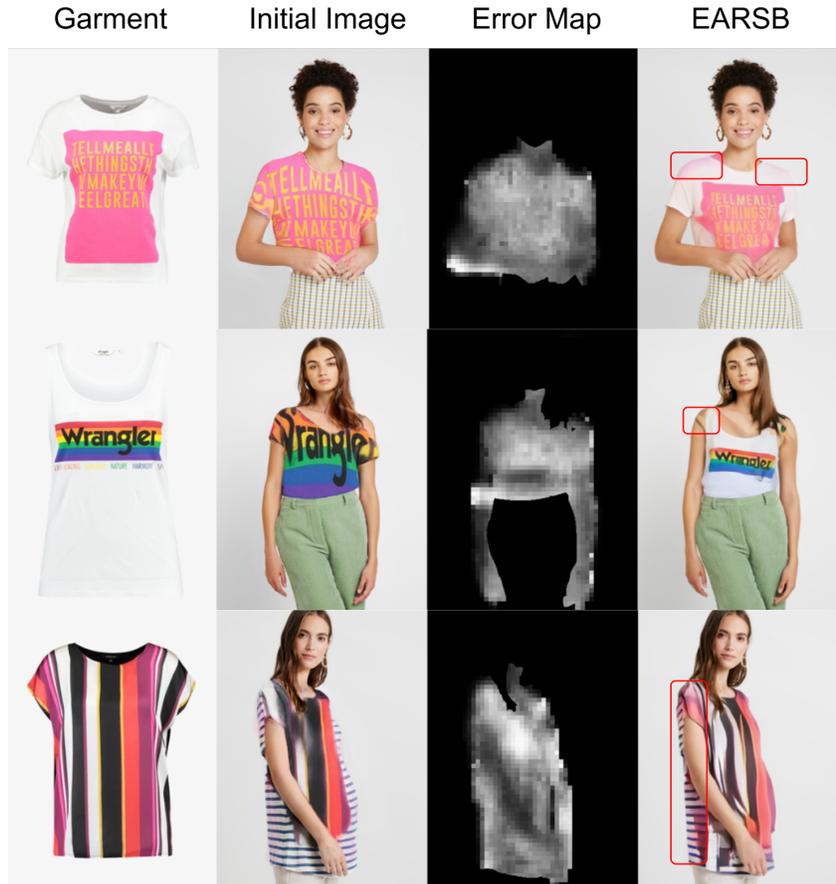


Figure 11: Failure cases on VITON-HD where the initial image has a poor-quality.

D ABLATIONS ON THE QUALITY OF THE INITIAL IMAGE x_1

In Tab. 4 we include the FID results of using different try-on GAN models to generate the initial image under the unpaired setting. Baseline means the GAN baseline. As previously stated in Sec. 4.1., we can draw three conclusions from the results: a) our EARSB can refine the GAN-generated image over the GAN baseline; b) the quality of the initial image x_1 is positively correlated with the quality of the sampled \hat{x}_0 ; c) our model achieves higher gains over CAT-DM, which also tries to refine the GAN-generated image but without error-aware noise schedule.

E LIMITATIONS

While our human-to-garment model can effectively generate synthetic paired data for try-on training augmentation, it has some imperfections. The overall quality of synthetic garments is regulated by our filtering criteria (Sec. 3.1.2), yet minor texture deformations occasionally occur. For instance, in Fig. 12, the second pair of the first row shows a misaligned shirt placket in the synthetic garment.

This limitation stems partly from the fact that our model is trained in the image domain which lacks 3D information. A potential solution is to utilize DensePose representations extracted from the garment as in (Cui et al., 2023).

A key constraint of our EARSB is its refinement-based nature, which makes the generated image dependent on the initial image. We assume that the initial image from a try-on GAN model is of reasonable quality, requiring only partial refinement. Consequently, if the initial image is of very poor quality, our refinement process cannot completely erase and regenerate an entirely new, unrelated image. Fig. 11 illustrates this limitation: in the first row, the initial image severely mismatches the white shirt with pink graphics. With EARSB refinement, while the shirt is correctly re-warped, color residuals from the initial image persist around the shoulder area.

F ADDITIONAL VISUALIZATIONS

Figures 12 and 13 showcase exemplars from our synthesized datasets H2G-UH and H2G-FH, respectively. The generated garment images closely mimic the product view of the clothing items, accurately capturing both the shape and texture of the original garments worn by the individuals. This approach to creating synthetic training data for the virtual try-on task is both cost-effective and data-efficient, highlighting the benefits of our proposed human-to-garment model.

Figures 14 and 15 give visualized results of the proposed EARSB and EARSB+H2G-UH. In contrast to previous approaches, EARSB specifically targets and enhances low-quality regions in GAN-generated images, which typically correspond to texture-rich areas. This targeted improvement is evident in the last row of Fig. 14, where EARSB more accurately reconstructs text *freinds*, and in the third row, where it successfully generates four side buttons. Furthermore, the incorporation of our synthetic dataset H2G-UH with EARSB leads to even more refined details in the generated images, demonstrating the synergistic effect of our combined approach.

REFERENCES

- Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Qinsheng Zhang, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, Tero Karras, and Ming-Yu Liu. ediff-i: Text-to-image diffusion models with ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022.
- Aiyu Cui, Sen He, Tao Xiang, and Antoine Toisoul. Learning garment densepose for robust warping in virtual try-on. *arXiv preprint arXiv:2303.17688*, 2023.
- Xintong Han, Xiaojun Hu, Weilin Huang, and Matthew R Scott. Clothflow: A flow-based model for clothed person generation. In *CVPR*, 2019.
- Sangyun Lee, Gyojung Gu, Sunghyun Park, Seunghwan Choi, and Jaegul Choo. High-resolution virtual try-on with misalignment and occlusion-handled conditions. In *ECCV*, 2022.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022.
- Sang-Heon Shim, Jiwoo Chung, and Jae-Pil Heo. Towards squeezing-averse virtual try-on via sequential deformation. In *AAAI*, 2024.
- Zhenyu Xie, Zaiyu Huang, Xin Dong, Fuwei Zhao, Haoye Dong, Xijin Zhang, Feida Zhu, and Xiaodan Liang. GP-VTON: Towards general purpose virtual try-on via collaborative local-flow global-parsing learning. In *CVPR*, 2023.
- Jianhao Zeng, Dan Song, Weizhi Nie, Hongshuo Tian, Tongtong Wang, and An-An Liu. Cat-dm: Controllable accelerated virtual try-on with diffusion model. In *CVPR*, 2024.

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

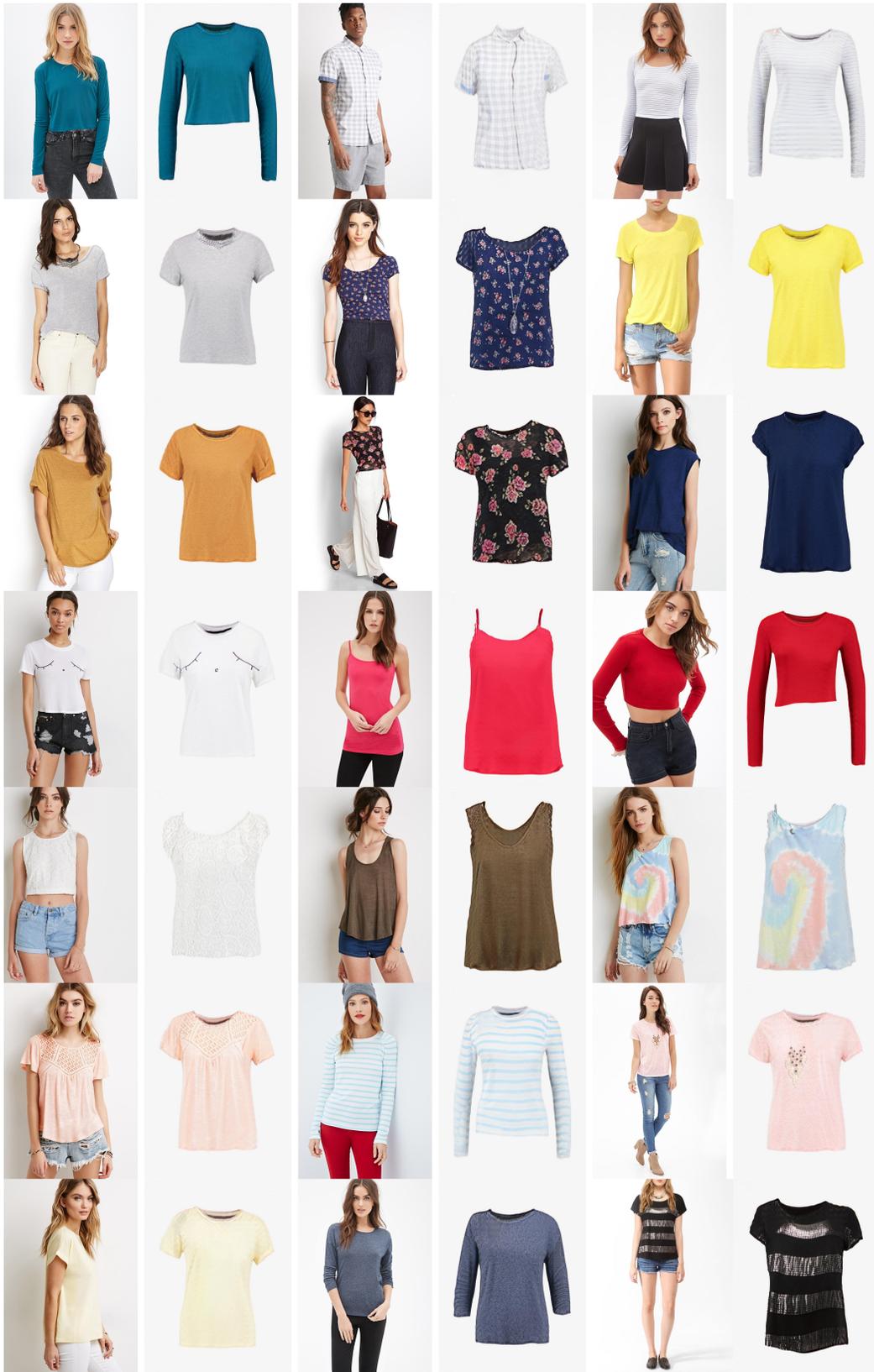


Figure 12: Visualized examples of the (human, synthetic garment) pairs on our proposed H2G-UH.

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

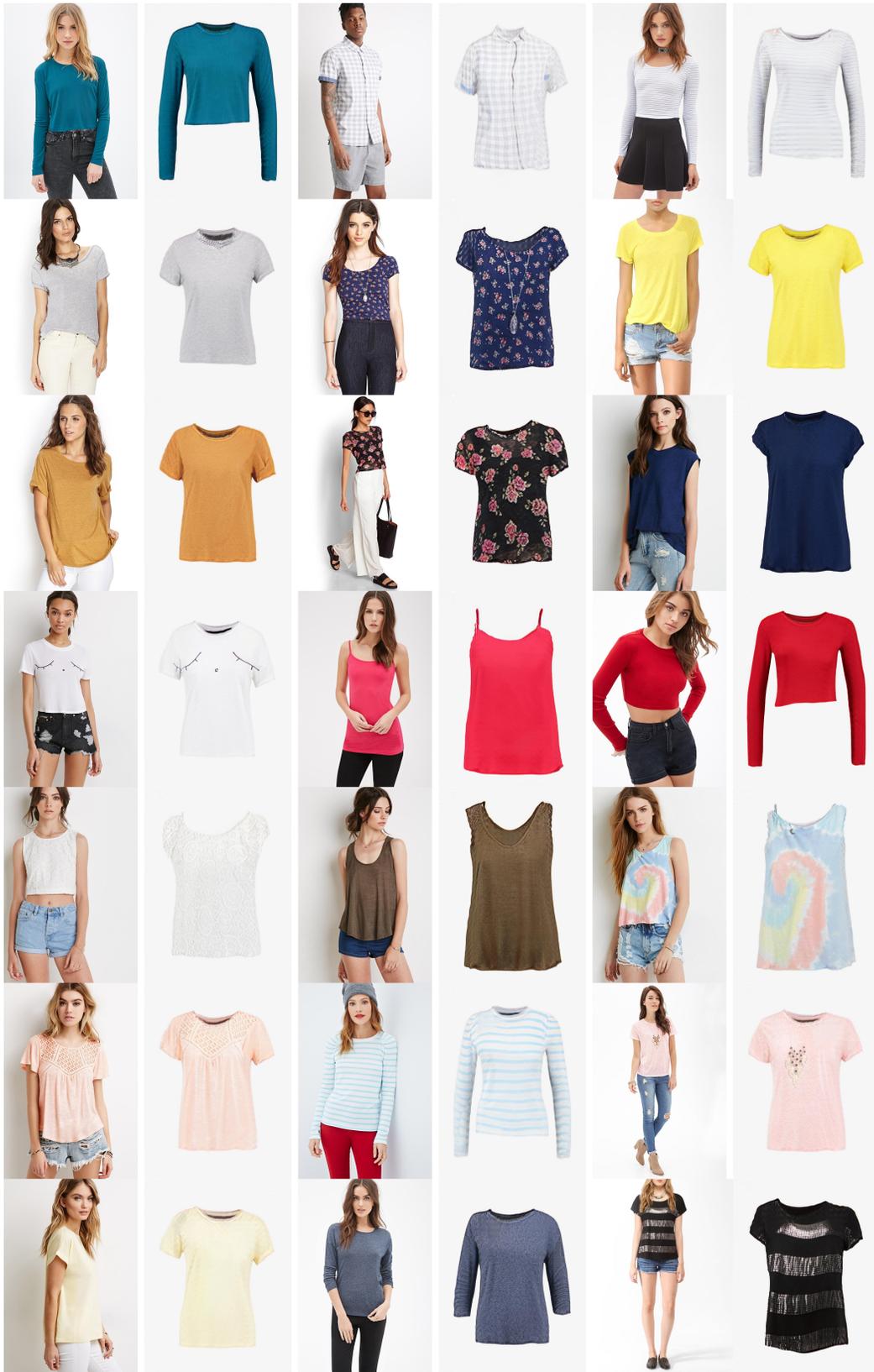


Figure 13: Visualized examples of the (human, synthetic garment) pairs on our proposed H2G-FH.

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

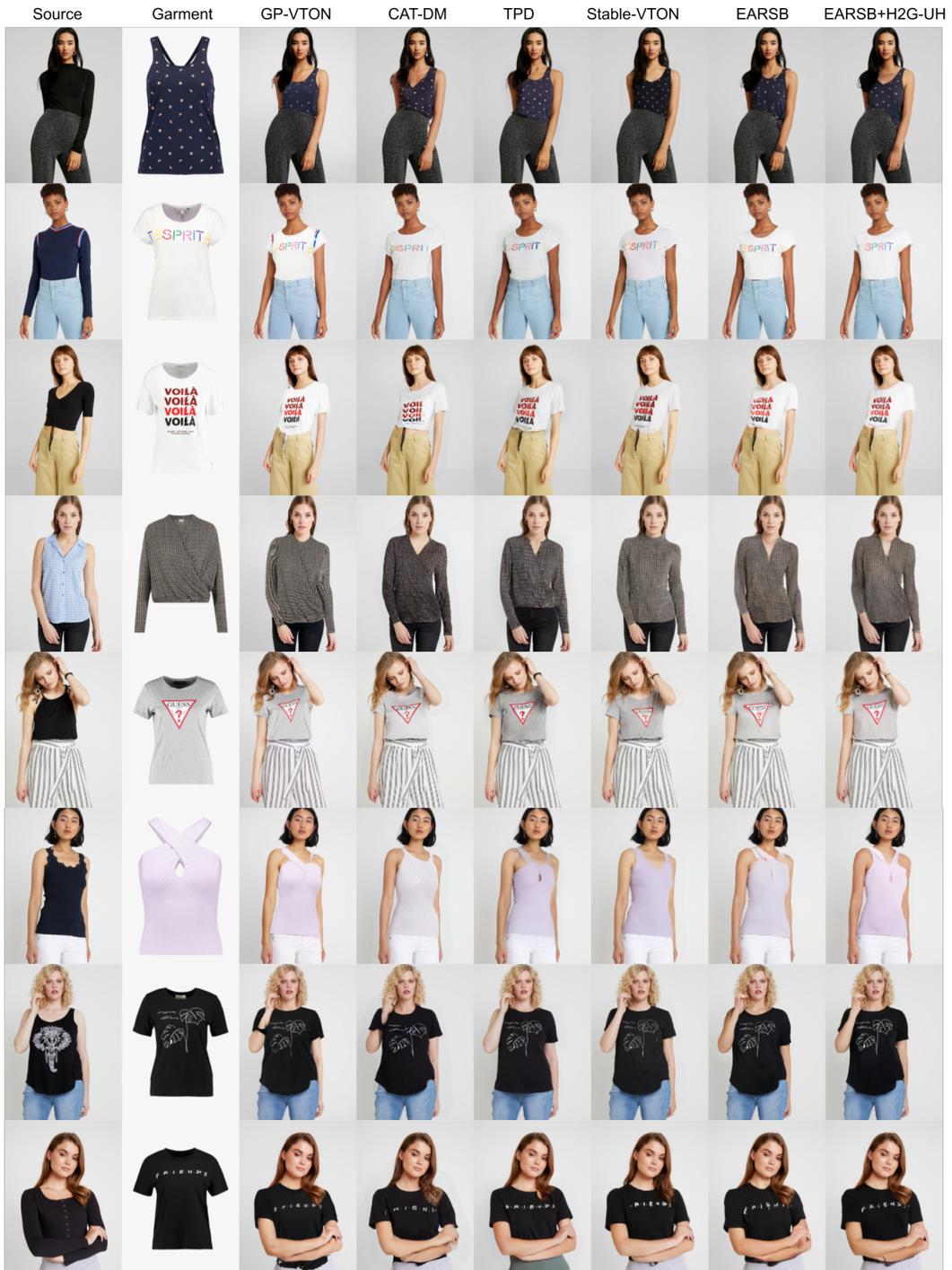


Figure 14: Visualized examples on VITON-HD. Our EARSB and EARSB+H2G-UH better recovers the intricate textures in the garment.

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

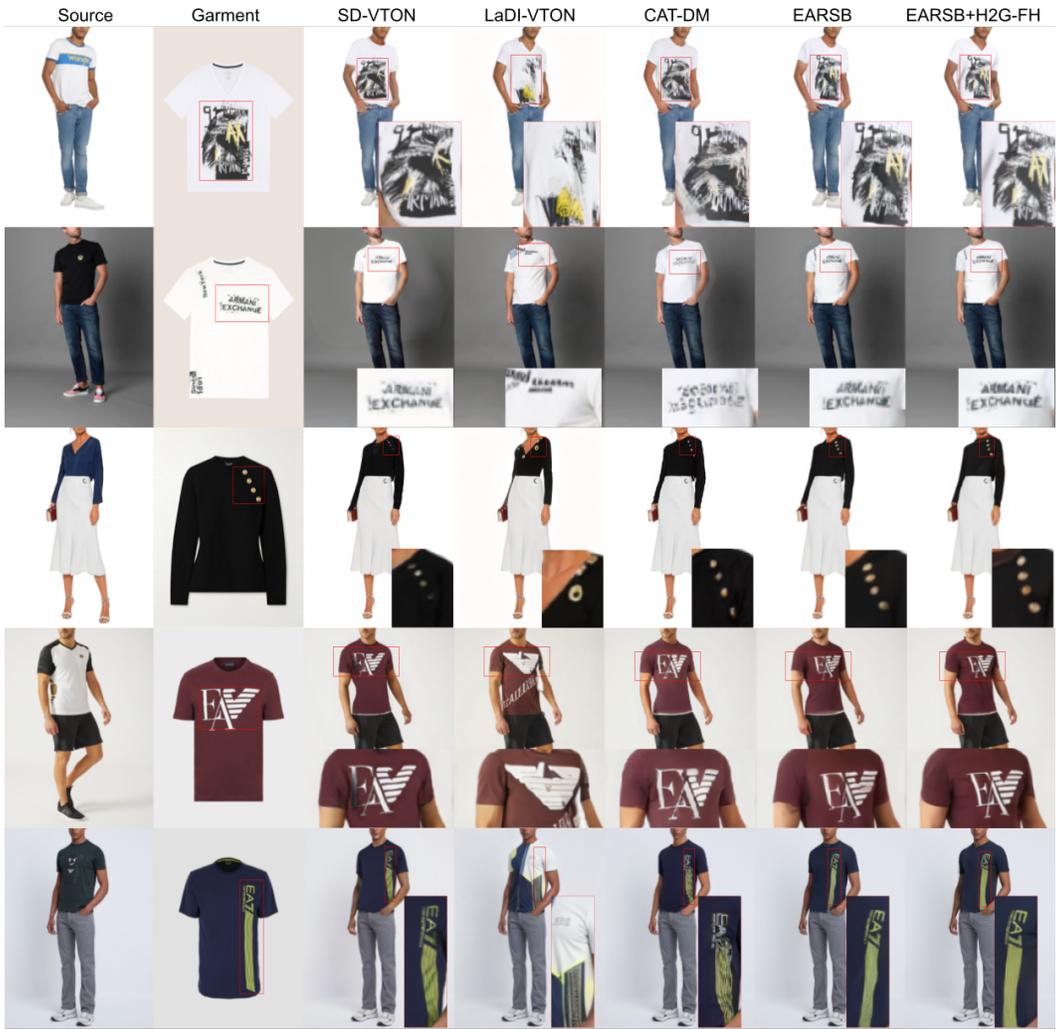


Figure 15: Visualized examples on DressCode-Upper. Our EARSB and EARSB+H2G-UH better reconstructs the texts and graphics in the garment.