

100  
101  
102  
103  
104  
105  
106  
107  
108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161  
162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199

## Supplementary Materials for

### “DCT-Net: Domain-Calibrated Translation for Portrait Stylization”

Paper ID 594

In this document we provide the following supplementary contents:

- Details of network architecture.
- Training details of content calibration network.
- Results of artistic portrait generation.
- Results of full-body image translation.
- Comparison with state-of-the-art methods.
- Limitations.

#### 1. Network architecture

For the content calibration network, the generator and discriminator follow the StyleGAN2 config-f models and the detailed architecture can be found in [9]. For the texture translation network, we provide network structure details in Table 1.

Table 1: Details of generator architecture.

Operation	Output Size
input	$256 \times 256 \times 3$
Conv+LReLU	$256 \times 256 \times 32$
Conv+LReLU+Conv+LReLU	$128 \times 128 \times 64$
Conv+LReLU+Conv+LReLU	$64 \times 64 \times 128$
Residual Block	$64 \times 64 \times 128$
Residual Block	$64 \times 64 \times 128$
Residual Block	$64 \times 64 \times 128$
Residual Block	$64 \times 64 \times 128$
Conv+LReLU+Conv transpose+LReLU	$128 \times 128 \times 64$
Conv+LReLU+Conv transpose+LReLU	$256 \times 256 \times 32$
Conv+LReLU	$256 \times 256 \times 3$

200  
201 Table 2: Details of discriminator architecture.  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
216  
217

Operation	Output Size
input	$256 \times 256 \times 3$
DeConv+SN+LReLU	$128 \times 128 \times 32$
Conv+SN+LReLU	$128 \times 128 \times 32$
DeConv+SN+LReLU	$64 \times 64 \times 64$
Conv+SN+LReLU	$64 \times 64 \times 64$
DeConv+SN+LReLU	$32 \times 32 \times 128$
Conv+SN+LReLU	$32 \times 32 \times 128$
Conv+SN	$32 \times 32 \times 1$

218  
219  
220  
221 Table 3: Details of regressor architecture.  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240

Operation	Output Size
input	$256 \times 256 \times 3$
DeConv+SN+LReLU	$128 \times 128 \times 32$
Conv+SN+LReLU	$128 \times 128 \times 32$
DeConv+SN+LReLU	$64 \times 64 \times 64$
Conv+SN+LReLU	$64 \times 64 \times 64$
DeConv+SN+LReLU	$32 \times 32 \times 128$
Conv+SN+LReLU	$32 \times 32 \times 128$
Mean pooling	128
Fully connected layer	512
Fully connected layer	3

300 2. Training details of content calibration network. 350

301 Given a style dataset contains approximate 100 images for a similar style. We start from the pretrained StyleGAN2 351  
302 model  $G_s$  trained on real faces (e.g., FFHQ dataset), and use a copy of  $G_s$  as our initialization model  $G_t$ . To adapt  $G_t$  352  
303 generating images in the target domain, we fine-tune  $G_t$  with full loss function consisting of the original adversarial 353  
304 loss and an identity loss: 354

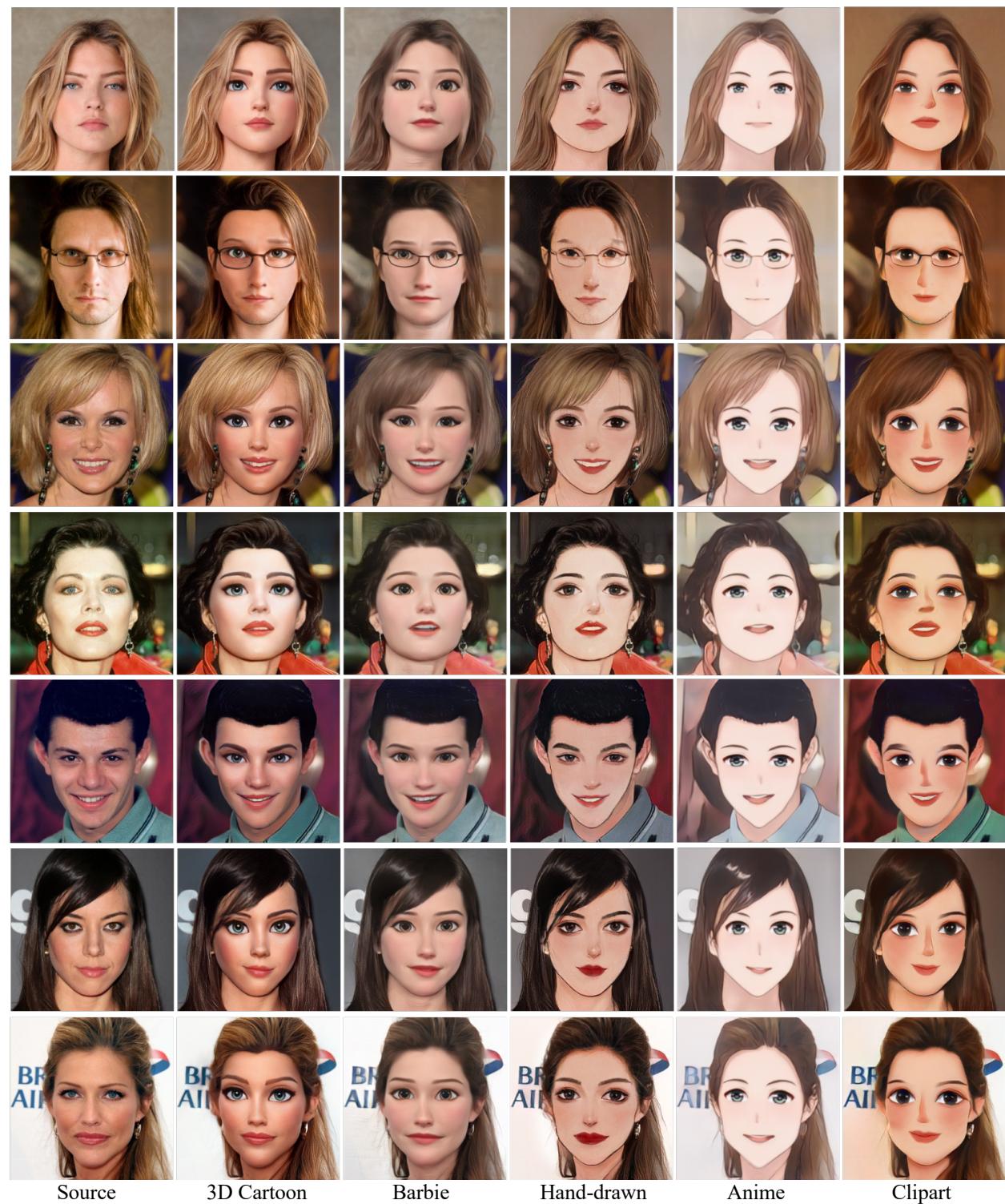
305 
$$\mathcal{L}_{ccn} = \mathcal{L}_{adv} + \lambda_{id}\mathcal{L}_{id},$$
 355

306 where  $\lambda_{id}$  denotes the weight of the identity loss.  $\mathcal{L}_{id}$  is formulated as: 356

308 
$$\mathcal{L}_{id} = 1 - \cos(z_{id}(\hat{x}_t), z_{id}(\hat{x}_s)),$$
 357

309 where  $\cos(\cdot, \cdot)$  represents the cosine similarity of two vectors and the id feature  $z_{id}$  is extracted from existing face 359  
310 recognition model [11].  $\hat{x}_t$  and  $\hat{x}_s$  are outputs of fixed generator  $G_s$  and learnable generator  $G_t$ , respectively. The 360  
311 define of  $\mathcal{L}_{adv}$  and more training parameters are the same with [9]. We set  $\lambda_{id}=0.1$  and train  $G_t$  for around 1000 361  
312 iterations. 362

313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323  
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

400 3. Results of artistic portrait generation.  
401444 Figure 1: Results of synthesized portraits in various styles. Source image credits: CelebA [2].  
445

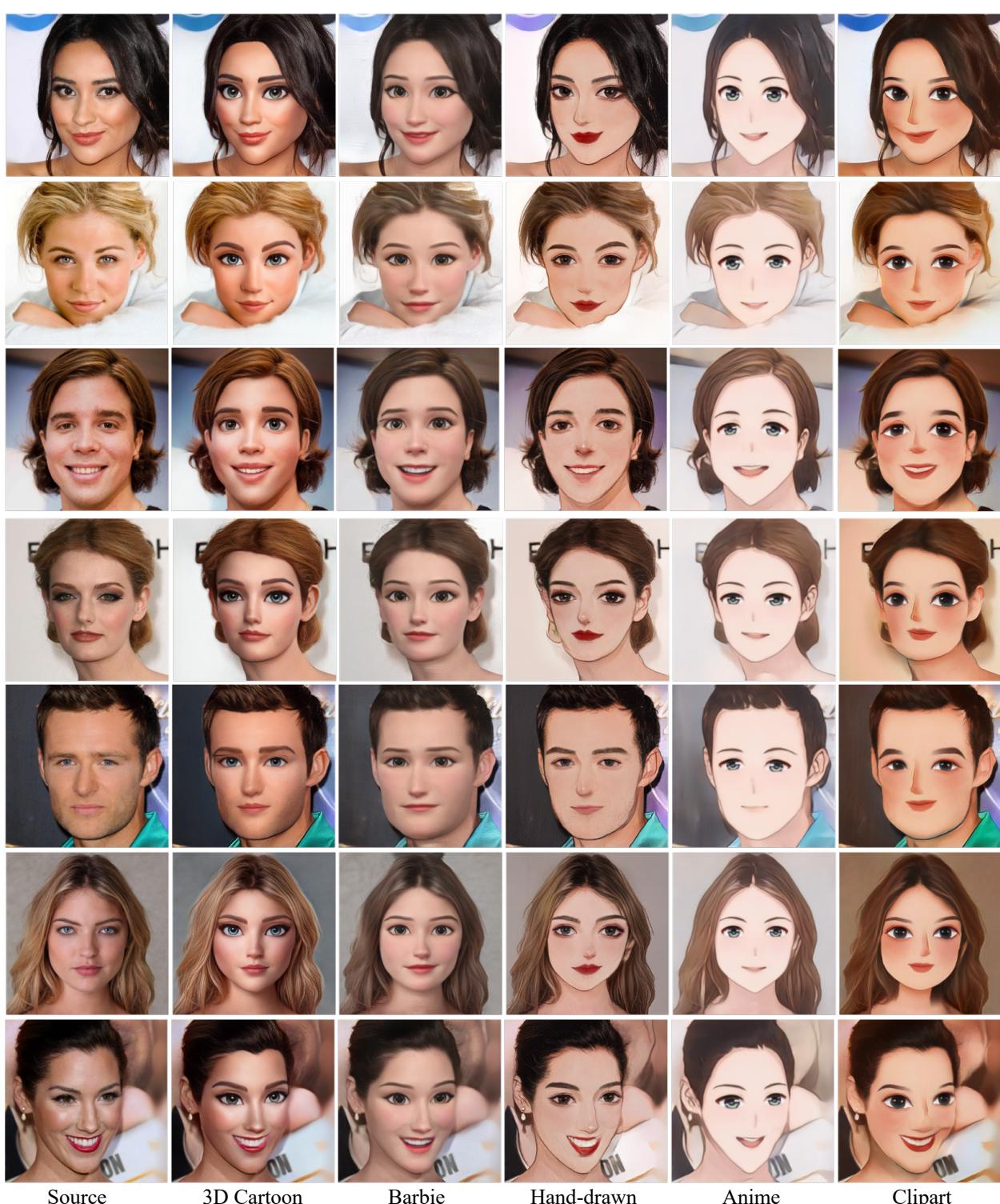
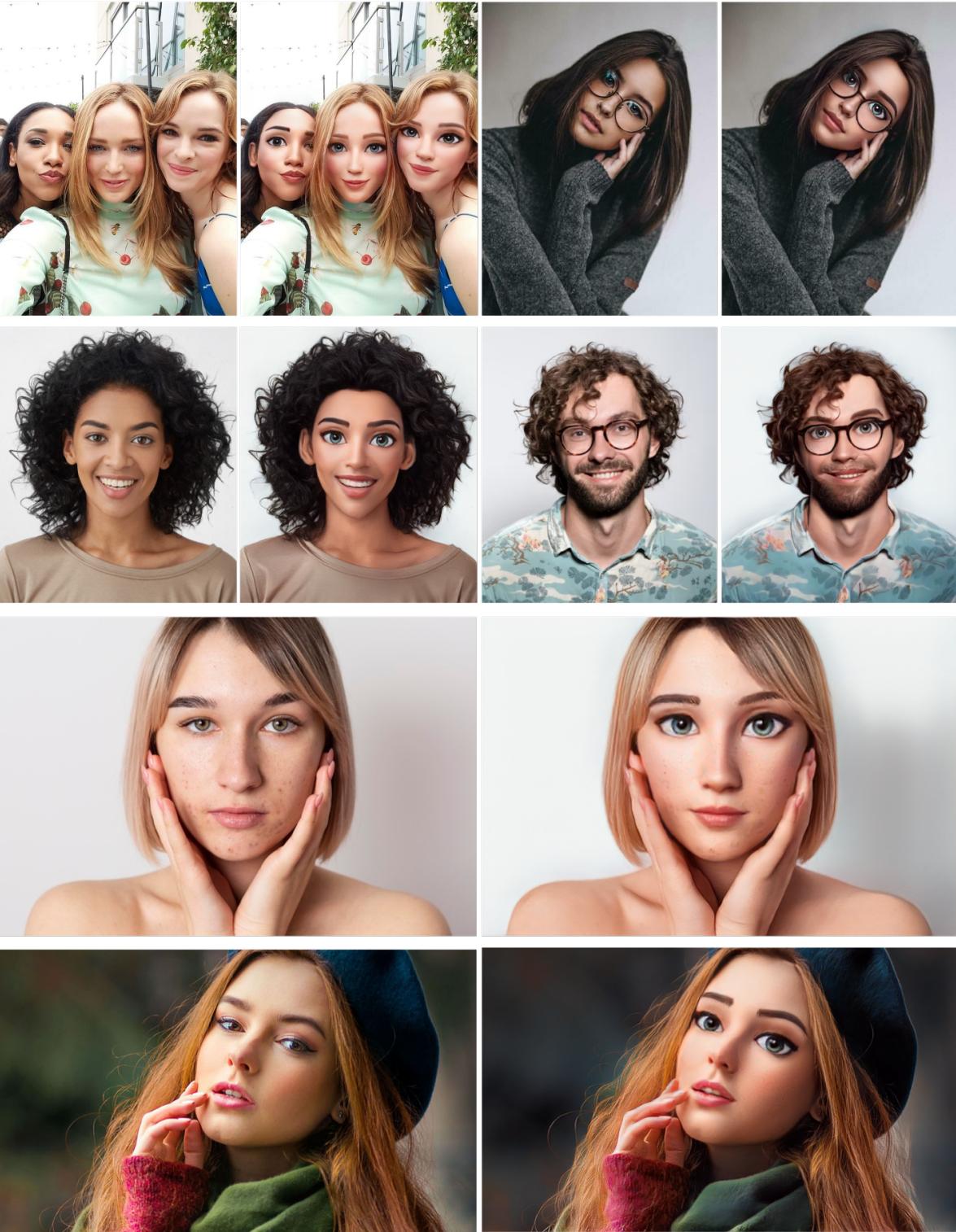


Figure 2: Results of synthesized portraits in various styles. Source image credits: CelebA [2].

600 4. Results of full-body image translation.

601 (a) 3D Cartoon

644 Figure 3: Results of stylized full images in 3D cartoon style. The source image in the left and the stylized result  
645 in the right. Source images: ©Unsplash[12], Google [1].

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

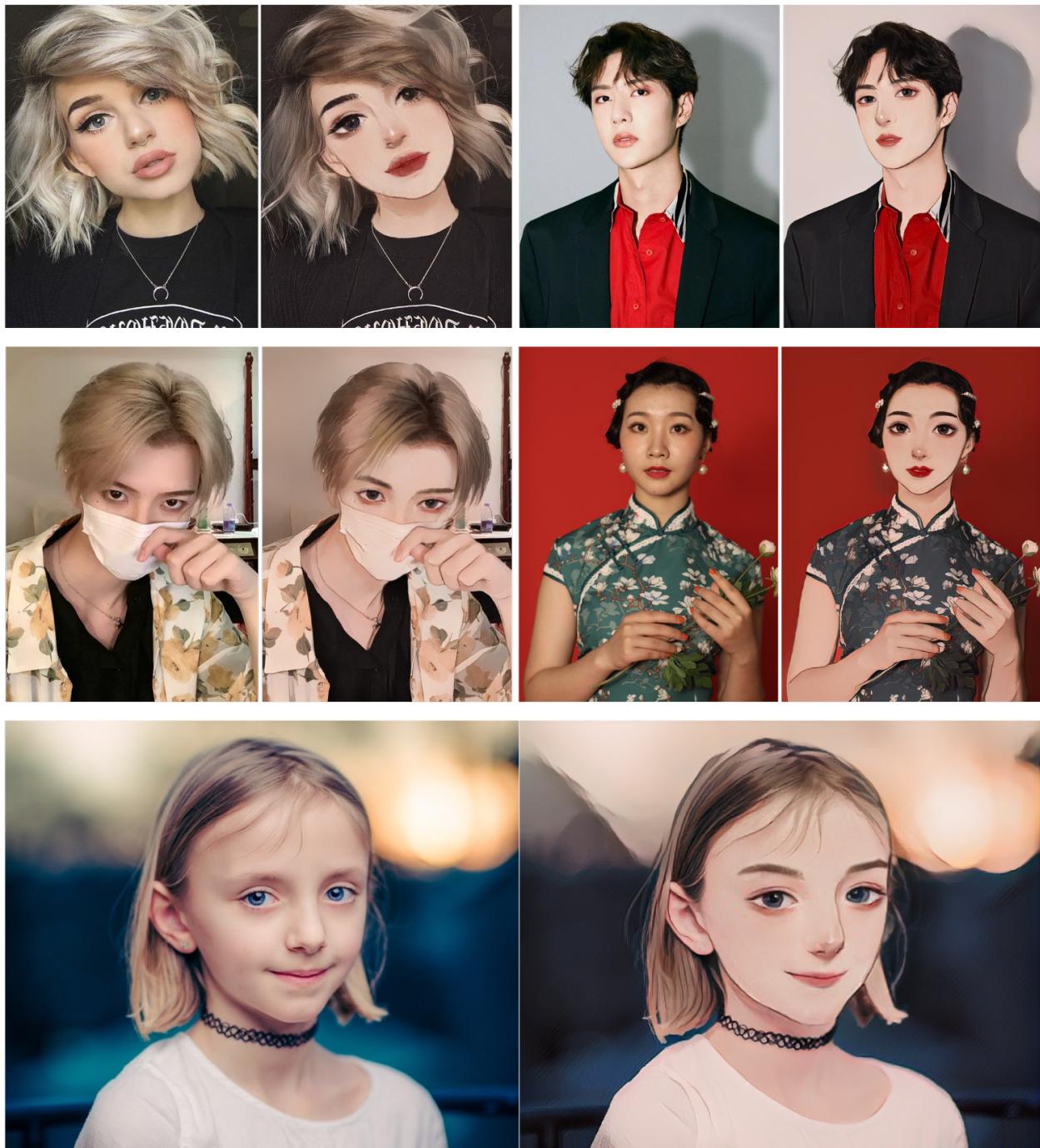
696

697

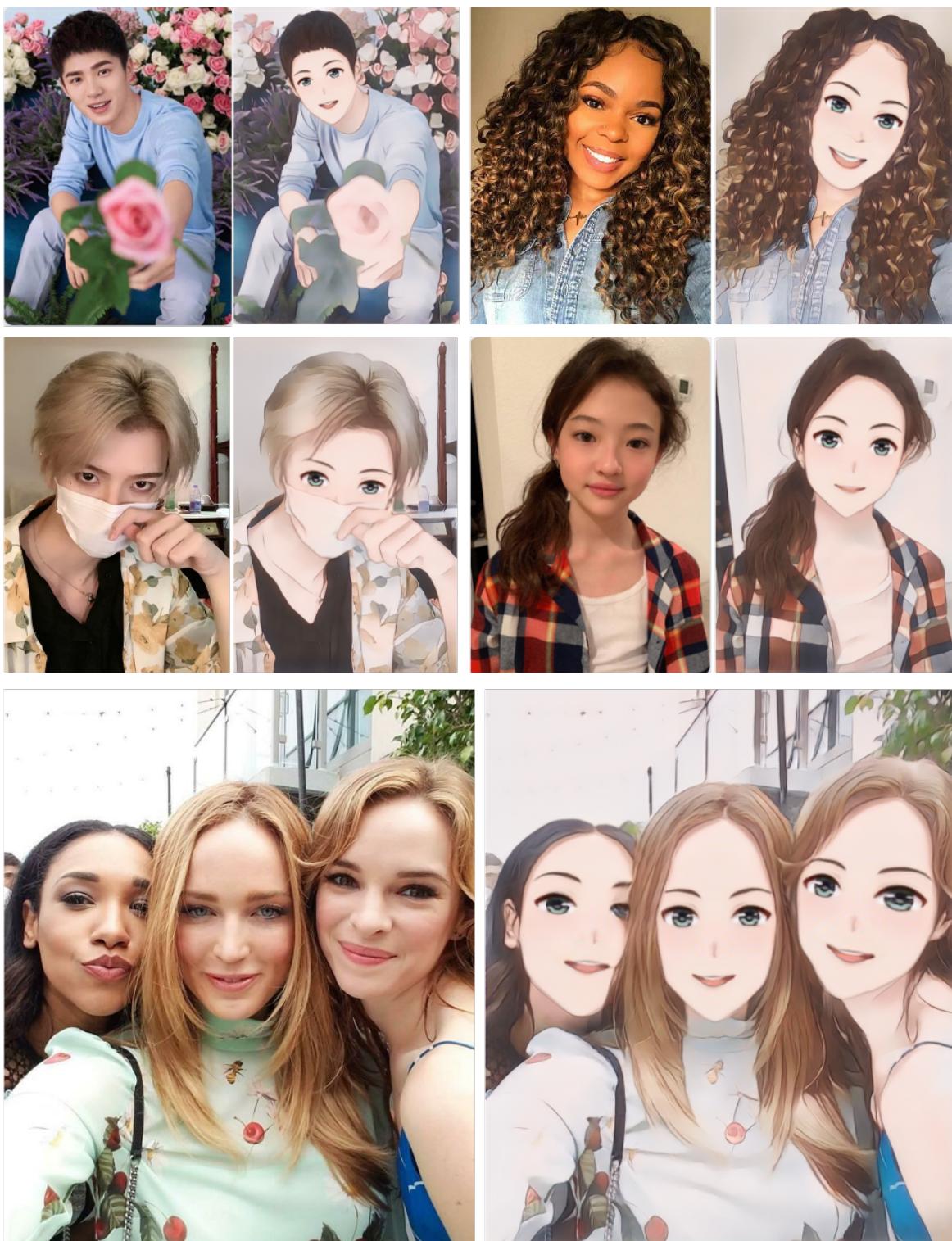
698

699

(b) Hand-drawn

742  
743  
744  
745  
746  
747  
748  
749  
Figure 4: Results of stylized full images in hand-drawn style. The source image in the left and the stylized  
result in the right. Source images: ©Unsplash[12], Google [1].

800 (c) Anime

844 Figure 5: Results of stylized full images in anime style. The source image in the left and the stylized result in  
845 the right. Source images: ©Google [1].

## 5. Comparison with state-of-the-art methods.

We provide more comparison results with four SOTA methods: CycleGAN [3], U-GAT-IT [4], Toonify [5], and PSP [6].

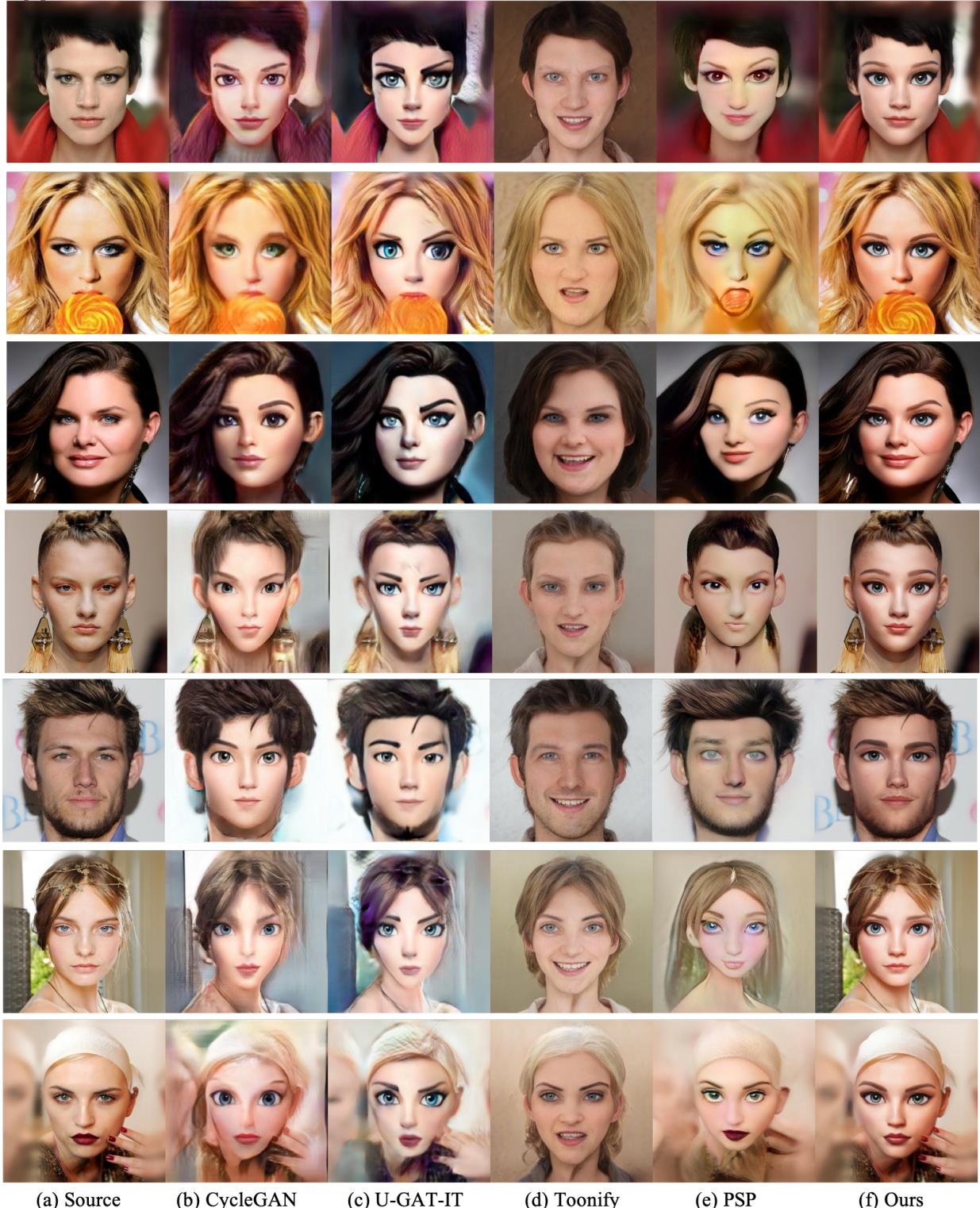


Figure 6: Qualitative comparison with state-of-the-art methods. Source images: ©CelebA [2].

We provide more comparison results with two SOTA methods: AgileGAN [7] and Few-shot-Ada [10].

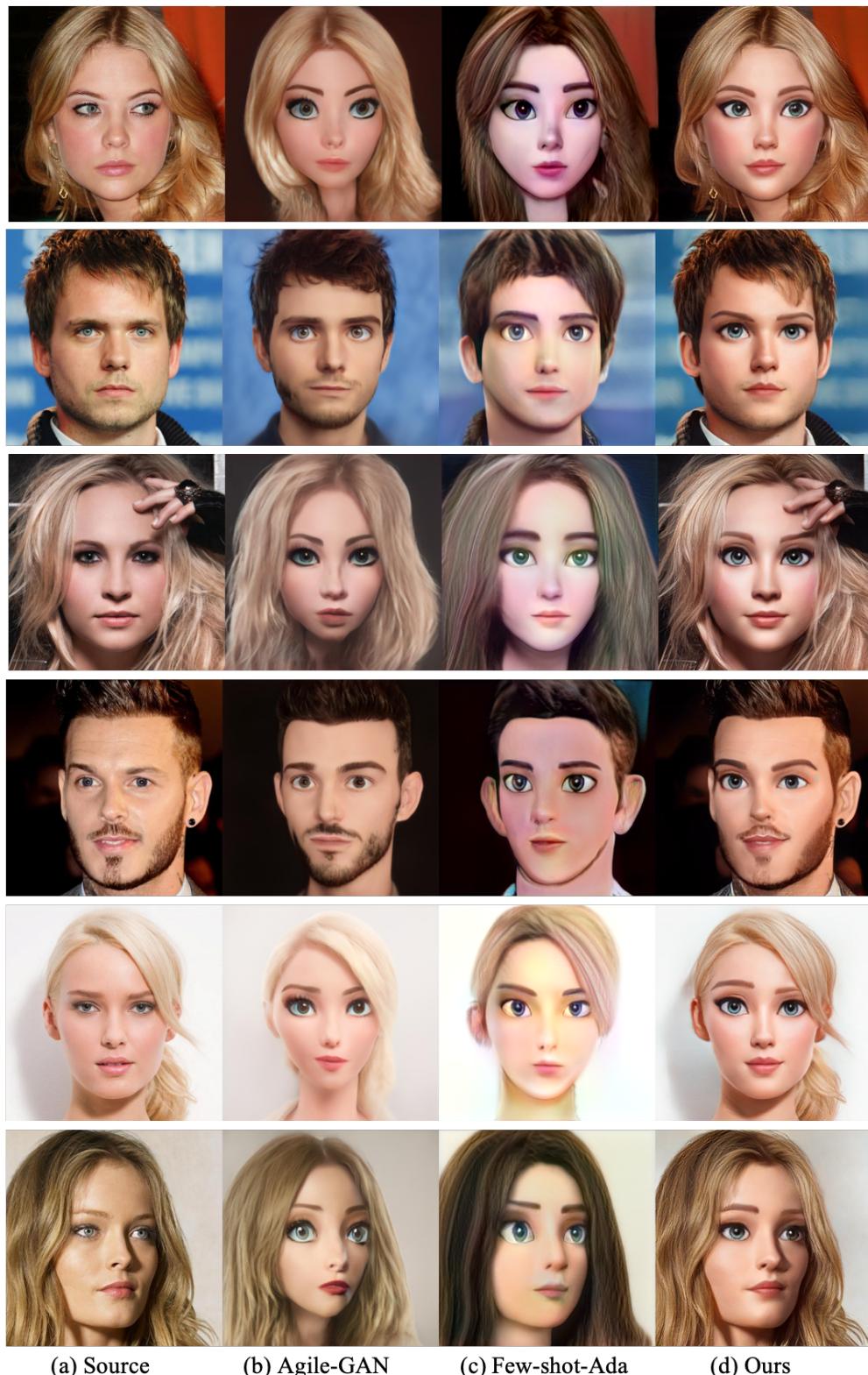


Figure 7: Qualitative comparison with state-of-the-art methods. Source images: ©AgileGAN [7].

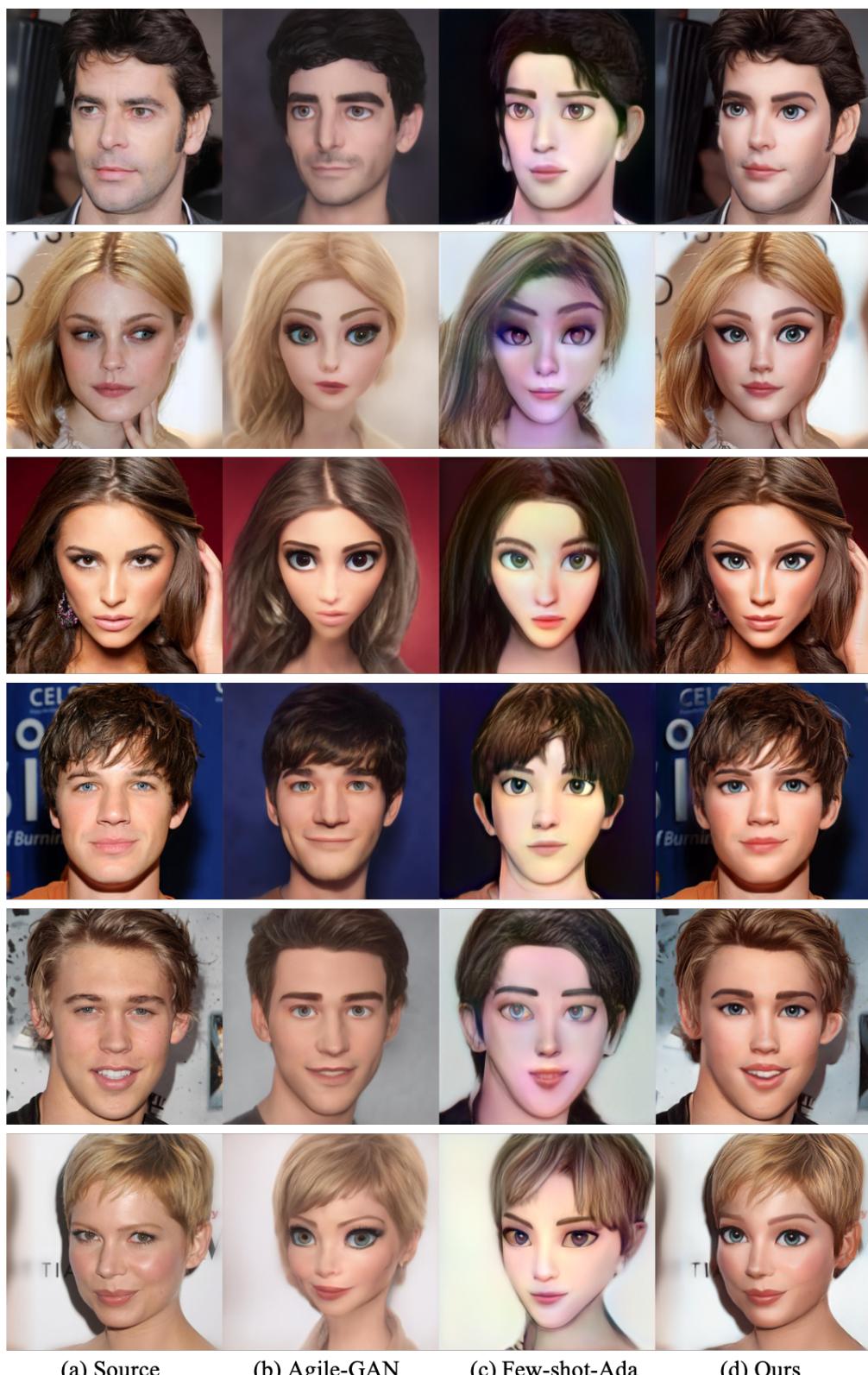
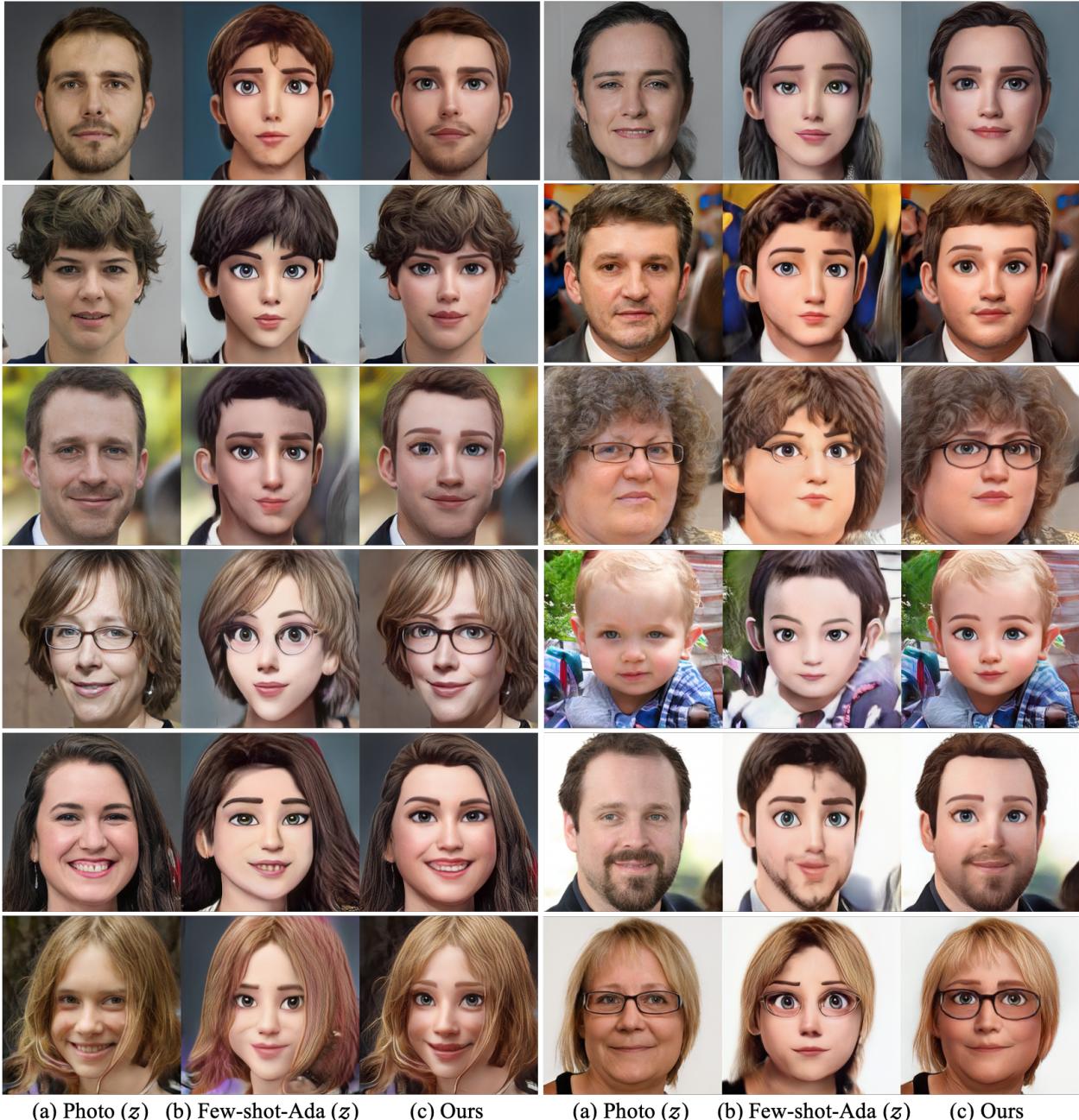


Figure 8: Qualitative comparison with state-of-the-art methods. Source images: ©AgileGAN [7].

1200 Due to the nature of unconditional generative model, Few-shot-Ada performed worse for arbitrary photo transfer.  
1201 So, we also evaluate this method in noise manner: random noises  $z$  are randomly sampled in its latent space, and we  
1202 input  $z$  into both the source generator and the adapted style generator to produce photo images  $I_p$  and stylized results  
1203  $I_r$ , respectively. These aligned pairs ( $I_p, I_r$ ) indicate the best capability of this adapted generation model for the  
1204 transfer task (no inversion error is introduced). We use the synthesized photo image  $I_p$  as the input of our network to  
1205 produce comparison results in Figure 9. As we can see, our method still outperforms this method with more content  
1206 details preserved.  
1207



1244 Figure 9: Comparison with Few-shot-Ada [8] in noise manner. Source images: ©Agile-GAN [7].  
1245  
1246  
1247  
1248  
1249

## 1300 6. Limitations.

1301 Because of the inherent characteristics of some styles (i.e., hand-drawn and anime), our synthesized results might  
1302 not be natural enough when there exist server lighting shadows in human faces, as shown in Figure 10. But some  
1303 styles (i.e., 3D cartoon) can still be well handled owing to its specific nature.



1304 (a) Source (b) 3D cartoon (c) Hand-drawn (d) Anime

1311 1312 1313 1314 1315 1316 1317 1318 1319 1320 1321 1322 1323 1324 1325 1326 1327 1328 1329 1330 1331 1332 1333 1334 1335 1336 1337 1338 1339 1340 1341 1342 1343 1344 1345 1346 1347 1348 1349 Figure 10: Failure cases due to the disturbed illumination.

## Reference

- [1] Google. [EB/OL]. <https://www.google.com/>.
- [2] Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and interactive facial image manipulation. In IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2020.
- [3] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In Proceedings of the IEEE 1064 international conference on computer vision, pages 2223–1065 2232, 2017.
- [4] Junho Kim, Minjae Kim, Hyeonwoo Kang, and Kwang Hee Lee. U-gat-it: Unsupervised generative attentional networks with adaptive layer-instance normalization for image-to-image translation. In International Conference on Learning Representations, 2020.
- [5] Justin NM Pinkney and Doron Adler. Resolution dependent gan interpolation for controllable image synthesis between domains. arXiv preprint arXiv:2010.05334, 2020.
- [6] Elad Richardson, Yuval Alaluf, Or Patashnik, Yotam Nitzan, Yaniv Azar, Stav Shapiro, and Daniel Cohen-Or. Encoding in style: a stylegan encoder for image-to-image translation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 2287–2296, 2021.
- [7] Guoxian Song, Linjie Luo, Jing Liu, Wan-Chun Ma, Chunpong Lai, Chuanxia Zheng, and Tat-Jen Cham. Agilegan: stylizing portraits by inversion-consistent transfer learning. ACM Transactions on Graphics (TOG), 40(4):1–13, 2021.
- [8] Tero Karras, Samuli Laine, and Timo Aila. A style-based 947 generator architecture for generative adversarial networks. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 4401–4410, 2019.
- [9] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 8110–8119, 2020.
- [10] Utkarsh Ojha, Yijun Li, Jingwan Lu, Alexei A. Efros, Yong Jae Lee, Eli Shechtman, and Richard Zhang. 2021. Few-shot Image Generation via Cross-domain Correspondence. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 10743–10752.
- [11] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. 2019. Arcface: Additive angular margin loss for deep face recognition. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 4690–4699.
- [12] Unsplash. [EB/OL]. <https://unsplash.com/>.