

# DiffStyle360: Diffusion-Based 360° Head Stylization via Style Fusion Attention

## Supplementary Material



Figure 1. Dataset pairs used in our training.

001

### 1. Dataset Pair Examples

002 Fig. 1 presents an example triplet from the *Joker* style domain,  
 003 consisting of a content image generated by a base  
 004 3D-GAN generator [1], a cross-identity style exemplar,  
 005 and the corresponding ground-truth stylized output produced by  
 006 a domain-adapted 3D-GAN model [2]. Note that the style  
 007 exemplar itself is also generated using the same domain-  
 008 adapted model. Using a style image from a different identity  
 009 encourages a clearer separation between identity and  
 010 style, enabling more effective disentanglement of content  
 011 and stylistic attributes during training. For our dataset, we  
 012 generate 150 samples for each of the six style domains,  
 013 yielding a total of 900 training pairs. A representative sub-  
 014 set of these samples is shown in Fig. 2.

015

### 2. Implementation Details

016

#### 2.1. Training Setup

017 We fine-tune the Style Appearance and View Consistency  
 018 modules of our model on a synthetic 3D-GAN-generated  
 019 stylized dataset for 800 iterations. The dataset includes six  
 020 style categories: Joker, Pixar, Sketch, Statue, Werewolf, and  
 021 Zombie. We use the AdamW optimizer with a learning rate  
 022 of  $10^{-5}$  and apply Classifier-Free Guidance (CFG) with a  
 023 weight of 3.0. The batch size is 16, corresponding to the  
 024 number of generated views per generation. All experiments  
 025 are conducted on a single NVIDIA A100 GPU. We use  $\tau =$   
 026 1.05 as the default key scaling value.

027

#### 2.2. Baselines

028 For diffusion-based stylization methods (IP2P and Instan-  
 029 tID), we first stylize the input content image using a text-  
 030 driven style prompt. The resulting stylized image is then  
 031 passed through DiffPortrait360 [6] to synthesize multi-view  
 032 stylized outputs.

033 We train the latent mapper of StyleCLIP [8] using the  
 034 PanoHead generator. For StyleGAN-NADA [5] and Style-  
 035 GANFusion [9], we rely on the official implementation  
 036 of StyleGANFusion and adapt the generator backbone to  
 037 PanoHead. In the case of StyleGANFusion, we use their

EG3D configuration for PanoHead, and for StyleGAN-  
 NADA [5] we incorporate adaptive layer selection as de-  
 scribed in the original method. DiffusionGAN3D [7] is im-  
 plemented directly from the paper, as no official code is  
 available. We preserve each baseline’s original hyperpa-  
 rameters, such as denoiser settings, noise schedule, learn-  
 ing rate, and optimizer, unless training instability requires  
 adjustments.

### 2.3. Quantitative Scores

To evaluate our approach, we generate three reference image distributions using the Stable Diffusion pipeline. The first distribution is created by adding noise at timestep  $t = 25$  to the input images and denoising them for 50 steps with the corresponding style prompt using each baseline’s diffusion checkpoints. The resulting edited images serve as ground-truth targets for computing FID and CLIP similarity. The second distribution consists of the stylized outputs produced by the 3D head stylization methods, which we compare against the first distribution for FID and against matched pairs for CLIP scores. The third distribution contains the original, unedited images and is used to evaluate identity preservation (ID) and  $\Delta D$ . For ID preservation, we compute ArcFace-based [4] identity similarity between the stylized outputs and their corresponding input images. For depth consistency,  $\Delta D$  is measured using depth maps estimated by the DepthAnything model [11]. These metrics are computed by comparing images from the second and third distributions. Fig. 3 shows representative examples from all three distributions.

### 3. User Study

We conducted a user study with **15 participants**, each of whom completed **18 comparison questions**. Options were randomly permuted to ensure fairness. In each question, the participant was shown:

- the original content image,
- a style exemplar image
- stylized outputs produced by different methods.

Participants were asked to select the image that provides the **best balance of identity preservation and stylization**, i.e., the image that both retains recognizable identity and reflects the target style. We did not ask separate questions for identity and stylization; participants provided a single holistic judgment for each comparison.

The stylization results are provided in gif format.

As shown in Table 1, participants preferred our method in **78.2%** of all comparisons, indicating that our approach

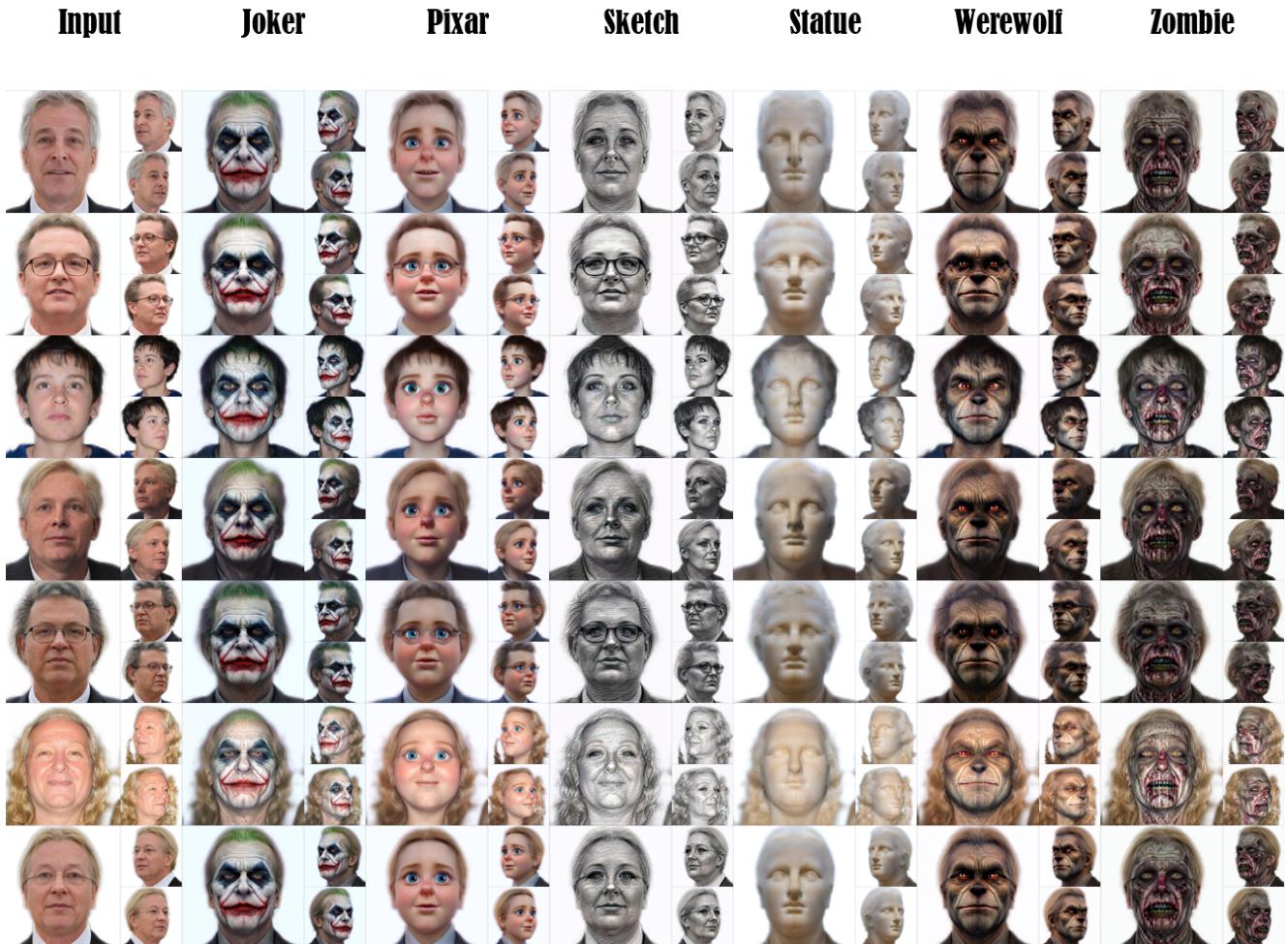


Figure 2. Dataset pair examples.

Method	User Study (%)
IP2P [3]	6.7%
StyleGAN-Fusion [9]	2.4%
Bilecen <i>et al.</i> [2]	12.7%
<b>Ours</b>	<b>78.2%</b>

Table 1. User study preference on the FFHQ dataset, aggregated over 270 responses (15 participants  $\times$  18 questions). Higher values indicate stronger user preference.

084 achieves the most favorable perceptual trade-off between  
085 identity retention and stylistic fidelity.

#### 086 4. Stylization-Identity Trade-Off Analysis

087 As shown in Tab. 2 and Fig. 6, the quantitative trends re-  
088 veal a clear stylization–identity trade-off among baseline

methods. InstantID [10] consistently reports high identity similarity scores across all styles; however, these metrics reflect its generally weak stylization capability rather than true identity preservation under meaningful edits. Because InstantID applies only subtle tonal changes, the output remains close to the input image, inflating ID scores while failing to capture the intended style.

089 IP2P [3] shows more inconsistent behavior across do-  
090 maines. For styles such as *Sketch*, *Statue*, and *Werewolf*,  
091 IP2P performs filter-like transformations that preserve most  
092 facial structure, leading to high identity scores in Tab. 2. In  
093 contrast, for semantically demanding styles such as *Joker*  
094 and *Pixar*, IP2P introduces stronger but less controlled  
095 changes that reduce identity similarity while still failing to  
096 deliver faithful stylization.

097 These quantitative observations are clearly reflected in  
098 the qualitative comparisons in Fig. 6. When integrated into  
099 the DiffPortrait360 pipeline for novel-view synthesis, both  
100

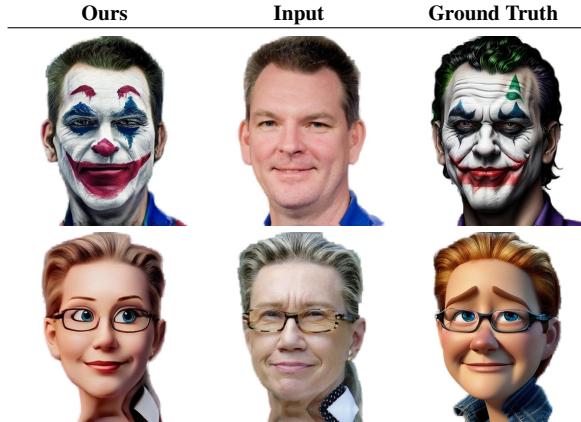


Figure 3. From left to right: a stylized image from the second distribution (outputs of our model), the corresponding unedited input image from the third distribution (used for evaluating identity preservation), and the edited result from the first distribution obtained by using the Stable Diffusion pipeline (used as the ground truth for FID and CLIP evaluations).

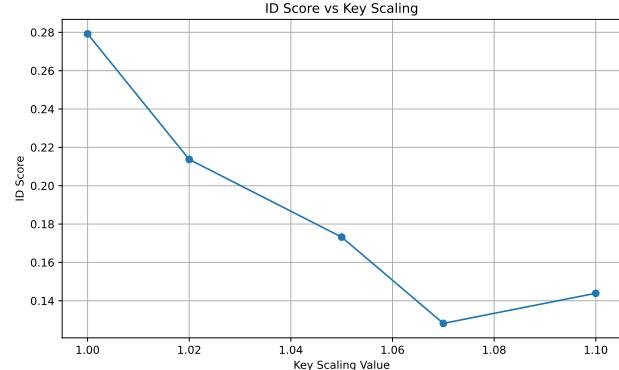
107 IP2P and InstantID struggle to maintain consistency across  
 108 viewpoints. Their limited or unstable stylization does not  
 109 generalize well to unseen views, frequently resulting in de-  
 110 graded back-view quality, texture collapse, and Janus-like  
 111 artifacts where facial features incorrectly appear on the back  
 112 of the head. These multi-view failures further indicate that  
 113 high ID scores in Tab. 2 do not correspond to robust styl-  
 114 ization, but rather to insufficient or inconsistent transforma-  
 115 tions.

116 In contrast, our approach delivers semantically aligned,  
 117 domain-consistent stylizations that remain coherent across  
 118 all viewpoints. As evidenced by both Tab. 2 and Fig. 6,  
 119 our method achieves a more favorable balance, produc-  
 120 ing strong stylization while maintaining significantly bet-  
 121 ter identity retention and multiview consistency than prior  
 122 methods.

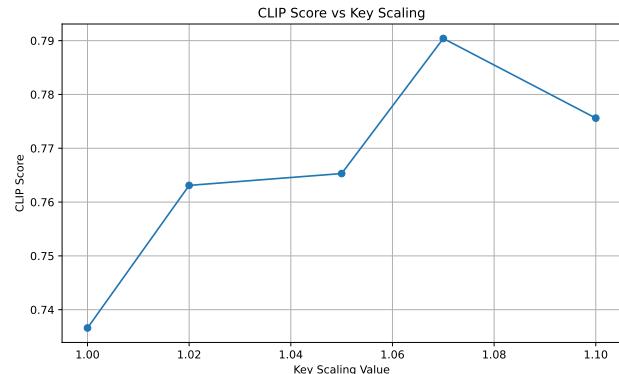
#### 123 4.1. Effect of Key Scaling on Stylization and Identity 124 Preservation

125 To analyze the impact of the proposed Key Scaling mecha-  
 126 nism, we evaluate identity similarity and CLIP-based style  
 127 alignment across different scaling values. As shown in  
 128 Fig. 4a, increasing the scaling factor consistently lowers  
 129 identity similarity: the ID score drops from 0.28 at  $\tau = 1.0$   
 130 to 0.13 at  $\tau = 1.07$ , indicating stronger deviation from  
 131 the input identity as stylistic cues are amplified in atten-  
 132 tion. Conversely, Fig. 4b shows that CLIP alignment im-  
 133 proves with larger scaling values, peaking at  $\tau = 1.07$  before  
 134 slightly declining at  $\tau = 1.10$ .

135 Overall, Key Scaling provides a controllable trade-  
 136 off between stylization strength and identity preserva-  
 137 tion. Smaller values retain identity but weaken stylization,  
 138 whereas larger values yield stronger style fidelity at the cost



(a) Identity score vs. Key Scaling.



(b) CLIP score vs. Key Scaling

Figure 4. Analysis of the effect of Key Scaling on stylization and identity preservation for Statue style on the RenderMe360 dataset. Increasing the scaling factor  $\tau$  enhances stylistic alignment (higher CLIP) but leads to reduced identity similarity (lower ID), demonstrating a clear stylization–identity trade-off.

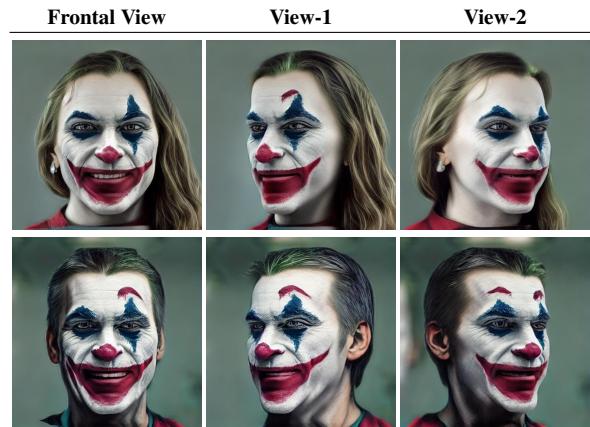


Figure 5. Limitations of our method. In rare cases, certain style elements may not remain fully consistent across views. As shown, the Joker eyebrow appearance exhibits slight inconsistency between angles.

Method	Joker		Pixar		Sketch		Statue		Werewolf	
	CLIP ↑	ID ↑	CLIP ↑	ID ↑						
IP2P [3]	0.8832	0.2237	0.8096	0.2046	0.6592	0.7334	0.7693	0.5829	0.7782	0.4554
InstantID [10]	0.6997	0.4662	0.6621	0.4857	0.6994	0.4946	0.7627	0.4706	0.6589	0.4915
StyleCLIP [8]	0.7821	0.1239	0.7481	0.1753	0.7413	0.1845	0.7123	0.2059	0.6634	0.1170
StyleGAN-Fusion [9]	0.8452	0.1105	0.7813	0.2211	0.7652	0.1636	0.8232	0.1385	0.8008	0.1477
StyleGANNADA [5]	0.8035	0.1728	0.7710	0.1873	0.7375	0.1666	0.7302	0.0934	0.7611	0.1329
DiffusionGAN3D [7]	0.7884	0.2298	0.8077	0.1545	0.6745	0.1707	0.8450	0.1552	0.7856	0.2036
Bilecen <i>et al.</i> [2]	0.8617	0.1877	0.8211	0.2336	0.7577	0.2752	0.8197	0.1952	0.8444	0.1755
<b>Ours</b>	0.8334	0.3132	0.7868	0.3500	0.7578	0.3428	0.8090	0.3529	0.7865	0.3186
<b>Ours + Key Scaling</b>	0.8649	0.1843	0.8116	0.2598	0.7466	0.1945	0.8182	0.2031	0.8545	0.1611

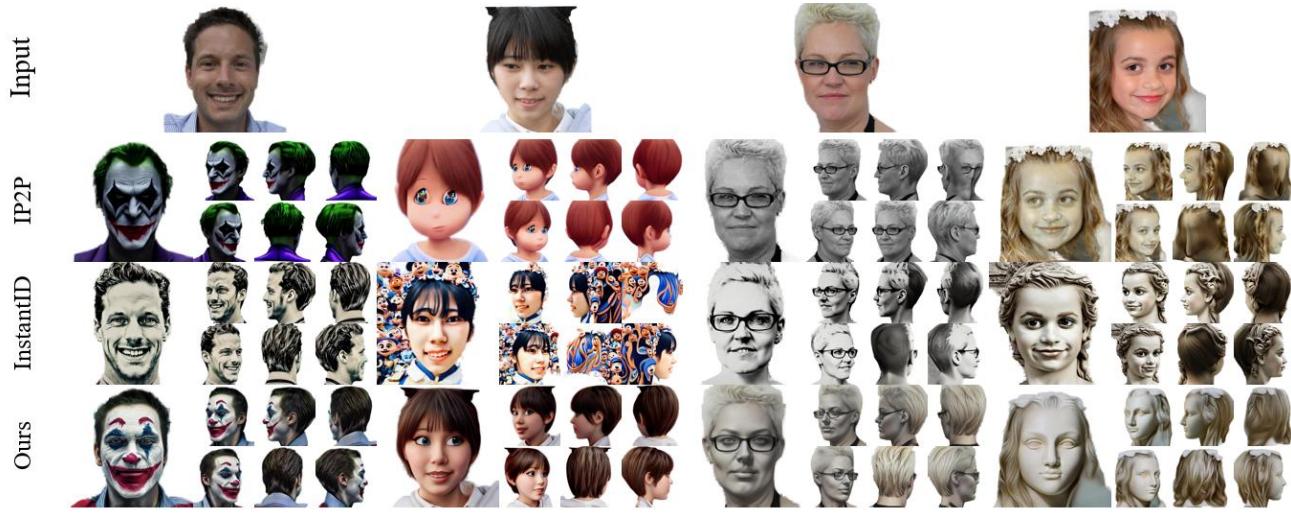
Table 2. CLIP and identity similarity scores across **Joker**, **Pixar**, **Sketch**, **Statue**, and **Werewolf** style domains on the FFHQ dataset.

Figure 6. Qualitative illustration of the stylization–identity trade-off and view-consistency behavior across methods for Joker, Pixar, Sketch, and Statue domains. IP2P applies filter-like adjustments for some domains such as Sketch and Statue (leading to inflated ID scores) but introduces uncontrolled changes for others (e.g., Joker and Pixar), causing identity drift. When combined with DiffPortrait360 for novel-view synthesis, both methods frequently fail to maintain consistent 3D structure, showing back-view degradation and Janus-like artifacts. InstantID produces only subtle tonal modifications, resulting in weak stylization but high identity scores as reflected in Tab. 2. Our method achieves strong, semantically aligned stylization with stable identity across all viewpoints and robust multi-view geometry.

139 of identity drift. We select  $\tau = 1.05$  as a balanced setting  
 140 that offers improved stylization with manageable identity  
 141 loss.

## 142 5. Limitations

143 Although our method achieves strong multi-view consistency  
 144 in the vast majority of cases, certain challenging styles  
 145 can still introduce localized inconsistencies across view-  
 146 points. As illustrated in Fig. 5, the stylized Joker eye-  
 147 brows appear coherent in some views but exhibit slight de-  
 148 formation or misalignment when rendered from other an-  
 149 gles. These artifacts typically arise when the target style  
 150 contains highly exaggerated or fine-grained facial elements

151 that place strong demands on cross-view geometric align-  
 152 ment.

153 Importantly, such failures occur only in rare cases and do  
 154 not reflect the general behavior of our system. Across most  
 155 styles and identities, our approach maintains stable geom-  
 156 etry, consistent textures, and coherent stylization across all  
 157 novel views.

## 158 6. Additional Results

159 Figs. 7 and 8 present the style reference images for a di-  
 160 verse set of artistic domains and the corresponding outputs  
 161 produced by our model for multiple input identities.

162

## References

- 163 [1] Sizhe An, Hongyi Xu, Yichun Shi, Guoxian Song, Umit Y.  
164 Ogras, and Linjie Luo. Panohead: Geometry-aware 3d  
165 full-head synthesis in 360 degrees. In *Proceedings of the*  
166 *IEEE/CVF Conference on Computer Vision and Pattern*  
167 *Recognition (CVPR)*, pages 20950–20959, 2023. 1
- 168 [2] Bahri Batuhan Bilecen, Ahmet Berke Gökmen, Furkan  
169 Güzelant, and Aysegül Dündar. Identity preserving 3d head  
170 stylization with multiview score distillation. In *Proceedings*  
171 *of the IEEE/CVF International Conference on Computer Vi-*  
172 *sion (ICCV)*, 2025. Poster. 1, 2, 4
- 173 [3] Tim Brooks, Aleksander Holynski, and Alexei A. Efros.  
174 Instructpix2pix: Learning to follow image editing instruc-  
175 tions. In *Proceedings of the IEEE/CVF Conference on Com-*  
176 *puter Vision and Pattern Recognition (CVPR)*, pages 18392–  
177 18402, 2023. 2, 4
- 178 [4] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos  
179 Zafeiriou. Arcface: Additive angular margin loss for deep  
180 face recognition. In *Proceedings of the IEEE/CVF Con-*  
181 *ference on Computer Vision and Pattern Recognition*, pages  
182 4690–4699, 2019. 1
- 183 [5] Rinon Gal, Or Patashnik, Haggai Maron, Amit H. Bermano,  
184 Gal Chechik, and Daniel Cohen-Or. Stylegan-nada: Clip-  
185 guided domain adaptation of image generators. *ACM Trans.*  
186 *Graph.*, 41(4), 2022. 1, 4
- 187 [6] Yuming Gu, Phong Tran, Yujian Zheng, Hongyi Xu, Heyuan  
188 Li, Adilbek Karmanov, and Hao Li. Diffportrait360: Consis-  
189 tent portrait diffusion for 360 view synthesis. In *Proceedings*  
190 *of the IEEE/CVF Conference on Computer Vision and Pat-*  
191 *tern Recognition (CVPR)*, pages 26263–26273, 2025. 1
- 192 [7] Biwen Lei, Kai Yu, Mengyang Feng, Miaomiao Cui, and Xu-  
193 ansong Xie. Diffusiongan3d: Boosting text-guided 3d gener-  
194 ation and domain adaptation by combining 3d gans and dif-  
195 fusion priors. In *Proceedings of the IEEE/CVF Conference*  
196 *on Computer Vision and Pattern Recognition (CVPR)*, pages  
197 10487–10497, 2024. 1, 4
- 198 [8] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or,  
199 and Dani Lischinski. Styleclip: Text-driven manipulation of  
200 stylegan imagery. In *Proceedings of the IEEE/CVF inter-*  
201 *national conference on computer vision*, pages 2085–2094,  
202 2021. 1, 4
- 203 [9] Kunpeng Song, Ligong Han, Bingchen Liu, Dimitris  
204 Metaxas, and Ahmed Elgammal. Diffusion guided do-  
205 main adaptation of image generators. In *Proceedings of the*  
206 *IEEE/CVF Winter Conference on Applications of Computer*  
207 *Vision (WACV)*, 2024. 1, 2, 4
- 208 [10] Qixun Wang, Xu Bai, Haofan Wang, Zekui Qin, Anthony  
209 Chen, Huaxia Li, Xu Tang, and Yao Hu. Instantid: Zero-shot  
210 identity-preserving generation in seconds. *arXiv preprint*  
211 *arXiv:2401.07519*, 2024. 2, 4
- 212 [11] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiao-  
213 gang Xu, Jiashi Feng, and Hengshuang Zhao. Depth any-  
214 thing v2. In *Advances in Neural Information Processing Sys-*  
215 *tems*, pages 21875–21911. Curran Associates, Inc., 2024. 1



Figure 7. Additional qualitative results across a wide range of identities and style domains.



Figure 8. Additional qualitative results across a wide range of identities and style domains.