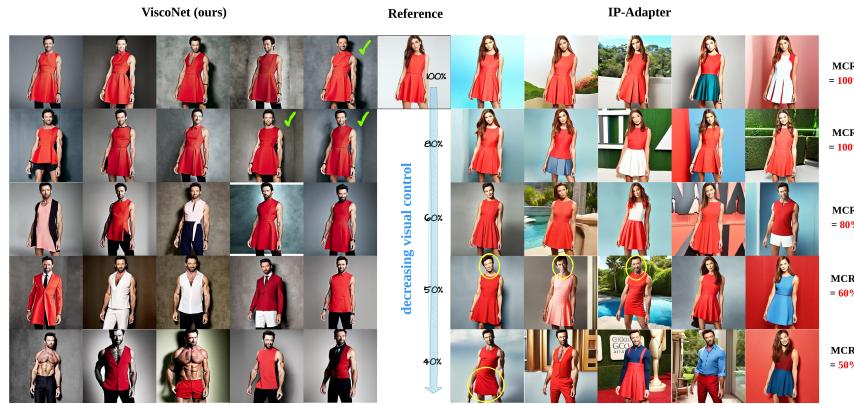


Appendix:



(a) IP-Adapter suffer 100% mode collapse at control strength over 80% and unable to generate the target person *Hugh Jackman*. Its visual conditioning power is much weaker when it finally escapes mode collapse at lower control strength, and **unable to generate the short pant**, and correct clothing style and color. In contrast, we are robust against mode collapse and avoid much of the problems above suffered by IP-Adapter, and able to generate desired results ✓ at 100% control strength, preserving faithfulness of both the person identity and clothing appearance.



(b) The conflict between the feminine reference image and Hugh Jackman's masculine image creates more conflict and hence mode collapse as suffered by IP-Adapter. IP-Adapter struggles to generate correct faces and pleated dress patterns (circled in yellow) at weaker control strength. This does not affect our method.

Fig. 15: Comparing the effect of control strength on re-identification task. IP-Adapter suffers much more severe mode collapse and struggles to create perfect image balancing the reference image and text prompt of *Hugh Jackman*.

control strengths from both our and IP-Adapter. With 100% strength, although IP-Adapter can reconstruct the reference image, it suffers 100%

Overall, our method is effectively mode collapse free at 60% while IP-Adapter still has 67% MCR. As shown in Figure 15, although high control strength introduces some mode collapse to our method. However, we can still generate high-quality images, preserving visual conditioning and a person’s identity.

621 A.2 Further Quantitative Comparison

We further explore qualitative results in this section. Unlike Figure 7-10, where we slide along the control strength on the same random seed to demonstrate latent space discontinuity, we extend Section A to present the best samples across all control strengths from both methods for direct comparison, as shown in Figure 16-18.

615
616
617
618
619
620617
618
619
620621
622
623
624
625
626



(a) Visual reference taken from the unseen test dataset.



(b) Unlike other movie stars with more diverse costumes, Prince Charles' limited clothing range presents the toughest challenge. (Top) IP-Adatper cannot produce any image of Prince Charles wearing the reference clothing. (Bottom) despite the extreme data gap, our method can produce reasonable images.

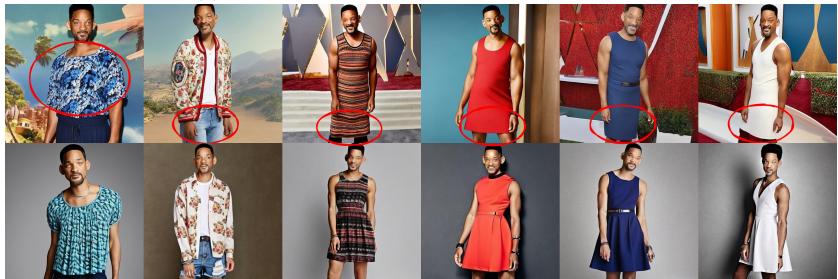
Fig. 16: Most challenging example in re-identification task - Prince Charles.

627 Among all the celebrities mentioned in the text prompt, *Prince Charles*¹ -
 628 known for having a limited wardrobe of formal attires in public images - presents
 629 the greatest challenge to the generalization capability of the models. IP-Adapter
 630 encounters difficulties and fails to generate any image of Prince Charles in casual
 631 or feminine clothing, as depicted in the reference image (Figure 16). In contrast,
 632 our method achieves reasonable success despite the monumental challenge. Figure
 633 17 - 18 shows samples from the rest of the text prompts used in the experiment.
 634 Overall, IP-Adapter needs to have much-lowered control strength to escape mode
 635 collapse, resulting in loss of fidelity in clothing to the reference images, including
 636 the incorrect length of pants or dress, wrong color and pattern, i.e., loss of the
 637 pleated dress pattern, it previously able to generate (Figure 15b).

¹ Stable Diffusion was trained on dated data before Prince Charles ascended to be king, so we adhere to his old title in the experiment.



(a) Visual reference taken from the unseen test dataset.



(b) Will Smith: (top) IP-Adaptor showing incorrect clothing color, length, or style (no pleated dress pattern). (bottom) Ours



(c) Dwayne Johnson: (top) IP-Adaptor (bottom) Ours



(d) Hugh Jackman: (top) IP-Adaptor (bottom) Ours.

Fig. 17: Re-identification comparison with IP-Adapter.



(a) Visual reference taken from the unseen test dataset.



(b) Keanu Reeves: (top) IP-Adaptor (bottom) Ours.



(c) Robert Downey Jr.: (top) IP-Adaptor (bottom) Ours



(d) Tom Cruise: (top) IP-Adaptor (bottom) Ours

Fig. 18: Re-identification comparison with IP-Adapter.

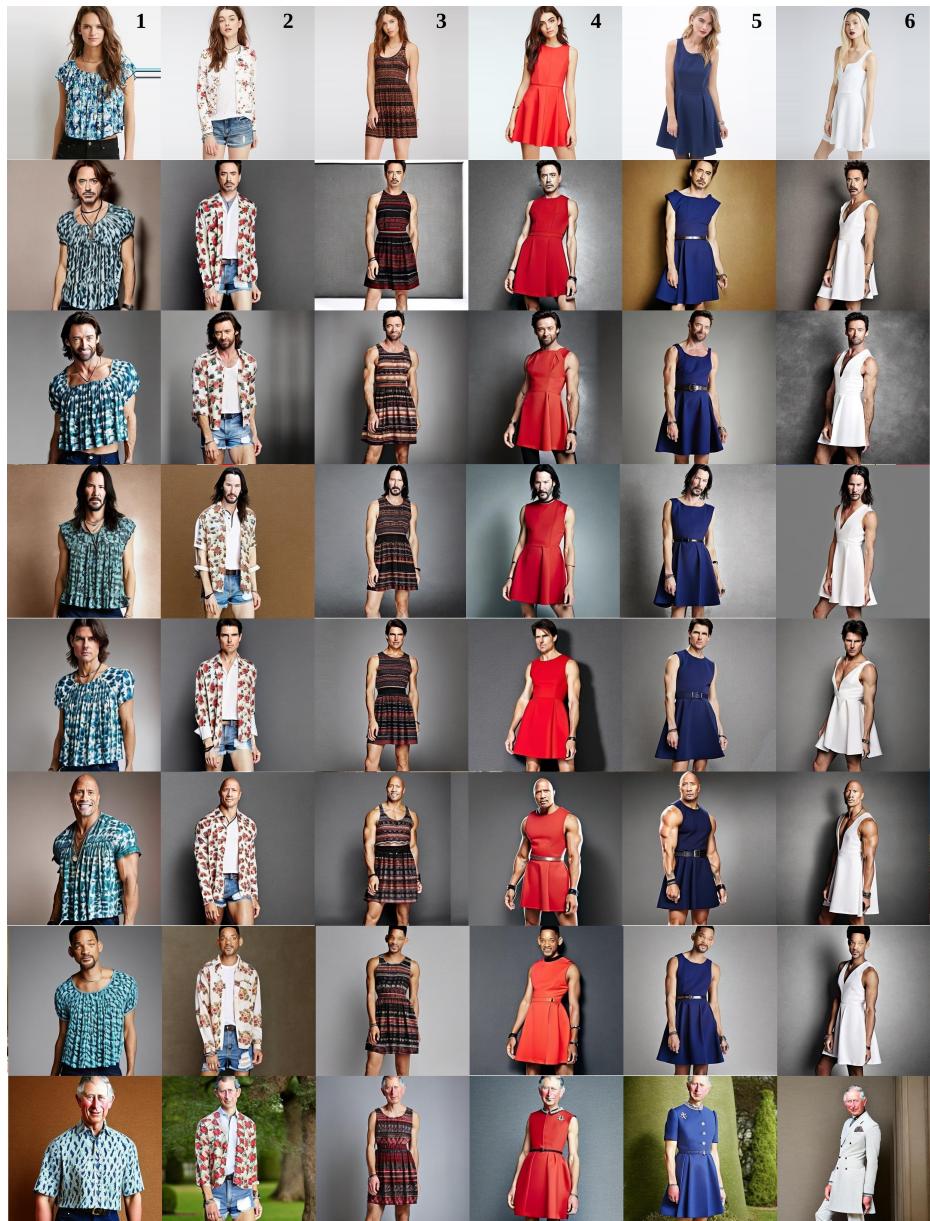


Fig. 19: Putting our images together shows the consistency of our method in delivering celebrity re-identification.

638

B Versatile Human Image Generation Task

638

639

B.1 Re-identification (visual prompt)

639

640

Figure 20 shows by conditioning on face and hair images, our method generates realistic people with diverse skin tones and body shapes correctly matching the faces despite the DeepFashion dataset consisting of more than 90% of female images, predominately fair-skinned women.

640

641

642

643



Fig. 20: Re-identification with a visual prompt.

644

B.2 Stylization

644

645

646

647

648

649

650

651

Figure 21 and Figure 22 show that our visual conditioning is effective across many image domains in creating a desired person's appearance, including various painting styles and also 3D objects such as statues, sculptures, toys, and 3D graphics. Some image domains have distinctive characteristics with considerable divergence from real photos, such as cartoons with disproportionate bigger heads, which can lead to a higher mode collapse rate. We circumvent this by removing the face mask to create results such as in Figure 21 and Figure 22.

645

646

647

648

649

650

651

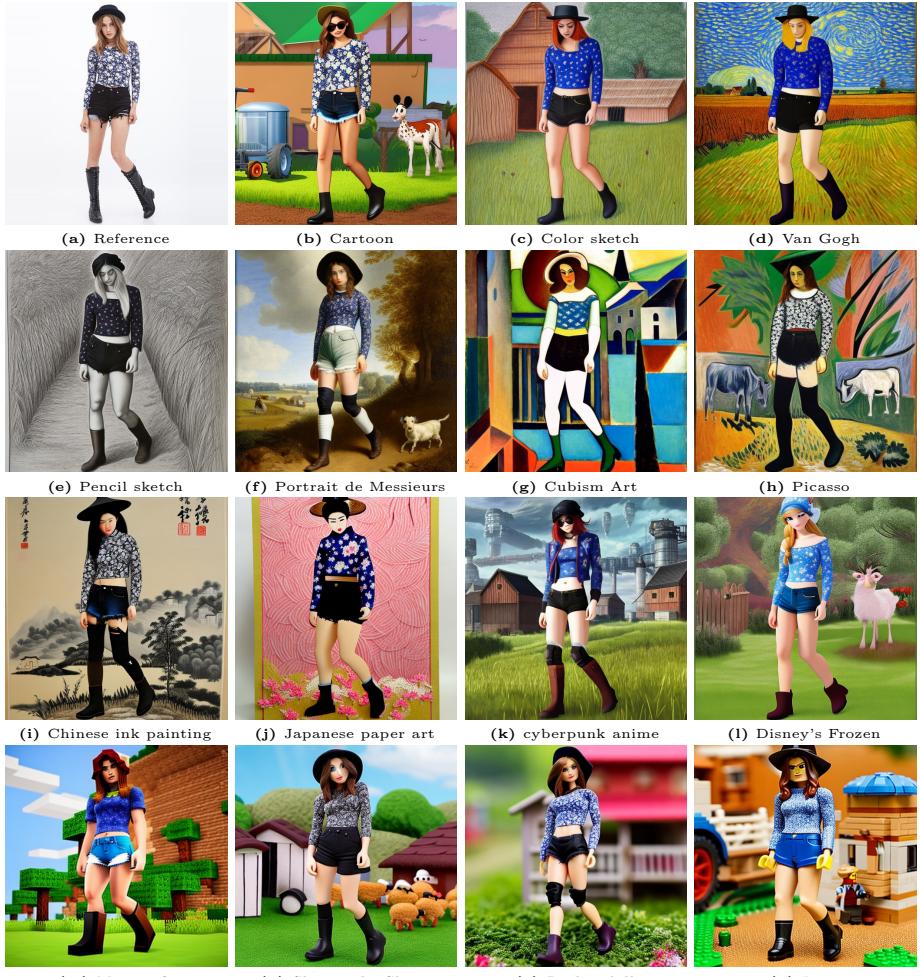


Fig. 21: Stylization. Text prompt: “*a woman, in farm.*”

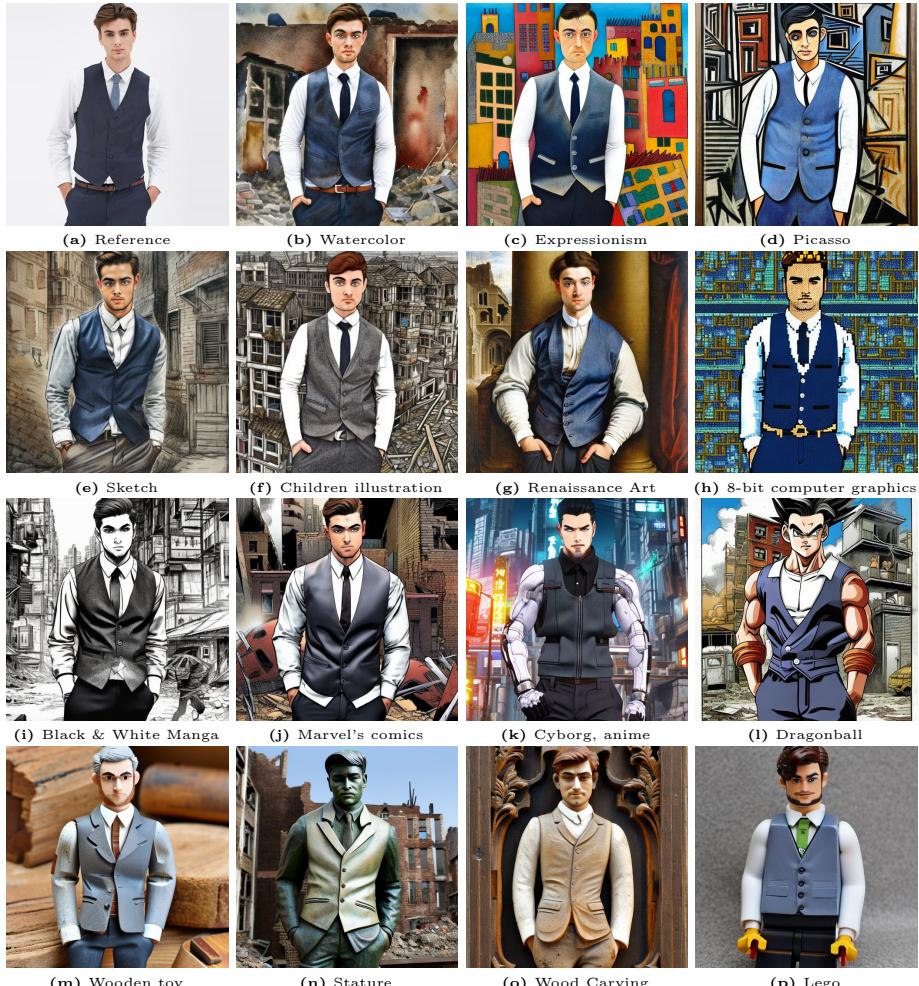


Fig. 22: Stylization. Text prompt: “*a man, in a derelict city.*”

B.3 Pose Re-target

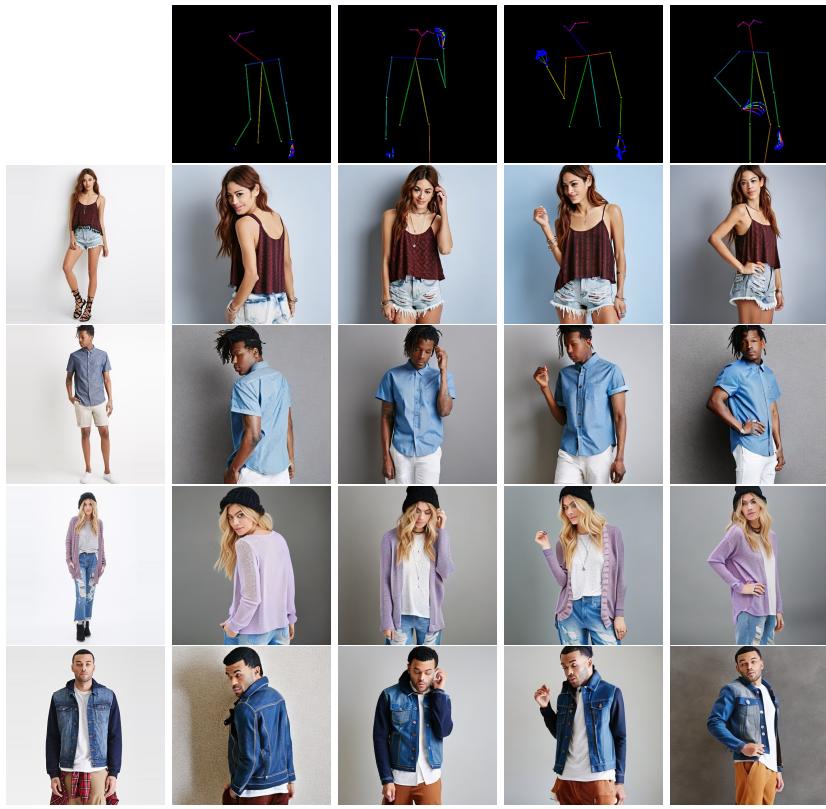


Fig. 23: Pose Transfer from (a) reference person to new poses in (b)-(e)

B.4 Virtual Try-on

Figure 24 demonstrates how we perform fashion virtual try-on using visual and text prompts. Figure 25 illustrates the culmination of our methods, showcasing the seamless integration of re-identification, virtual try-on, and pose re-target.



Fig. 24: High-resolution virtual try-on with real-world background. (Top) reference fashion for visual conditioning. (Bottom): virtual try-on results.

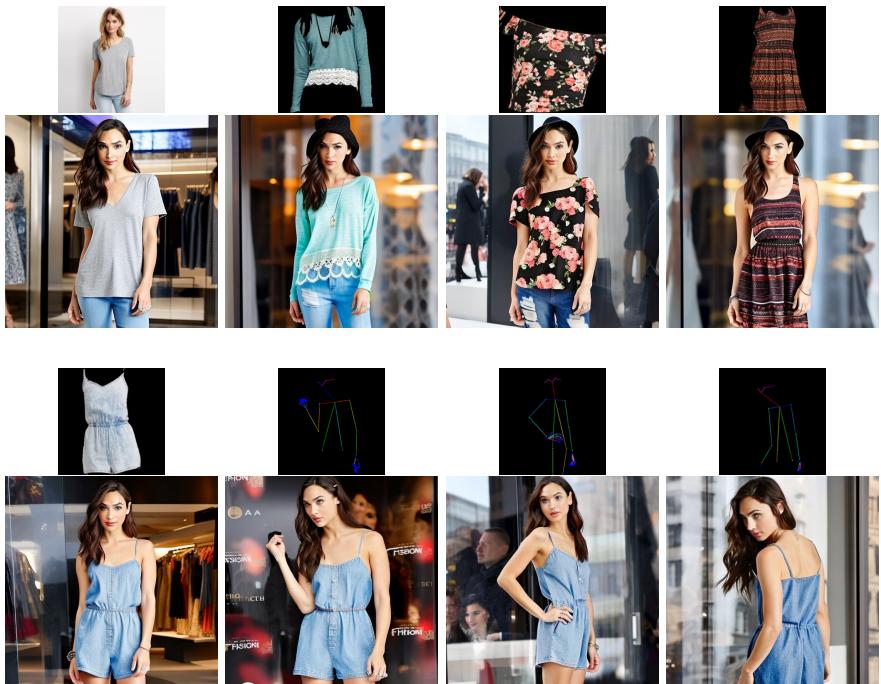


Fig. 25: Combining re-identification, virtual try-on, and pose re-target, we showcase examples of posing fashion with celebrity avatars.

C Quantitative Result

C.1 Section 4.1: Mode Collapse Quantitative Result

Table 2 shows the quantitative results corresponding to Figure 6 in Section 4.1 - Mode Collapse and Control Strength. Our method produces a higher CLIP score than the baseline at various control strengths, indicating less mode collapse. This is more evident in CLIP accuracy; at control strength 0.5, we achieve 100% (or 0% MCR) while baselines have only 46% and 63% ControlNet and IP-Adapter, respectively.

Strength	0.0	0.2	0.3	0.4	0.5	0.6	0.8	1.0
<u>CLIP score</u>								
ControlNet	0.2720	0.2620	0.2440	0.2340	0.2300	0.2300	0.2240	0.2260
IP-Adapter	0.2900	0.2920	0.2780	0.2620	0.2360	0.2120	0.1780	0.1900
ViscoNet(Ours)	0.2860	0.2940	0.2920	0.2900	0.2800	0.2660	0.2420	0.2220
<u>CLIP accuracy</u>								
ControlNet	0.8660	0.7720	0.7000	0.6180	0.4620	0.5500	0.5020	0.5760
IP-Adapter	0.9800	0.9700	0.8800	0.7500	0.6300	0.4100	0.1500	0.2100
ViscoNet(Ours)	1.0000	1.0000	1.0000	1.0000	1.0000	0.9000	0.7000	0.5760
<u>Pose accuracy (OKS)</u>								
ControlNet	0.0880	0.4139	0.6610	0.8305	0.8596	0.8852	0.9223	0.9348
IP-Adapter	0.5379	0.5412	0.6060	0.6813	0.7546	0.8074	0.9010	0.9298
ViscoNet(Ours)	0.0446	0.1654	0.3869	0.6580	0.7845	0.8253	0.8824	0.9102

Table 2: Reduced control strength results in higher CLIP scores and accuracy, translating to less mode collapse.

Figure 26 shows the breakdown of CLIP accuracy across the image styles in Table 2. Based on the same Stable Diffusion model, all models have shown the highest mode collapse rate in Van Gogh’s painting style, while Ukiyoe is the least affected.

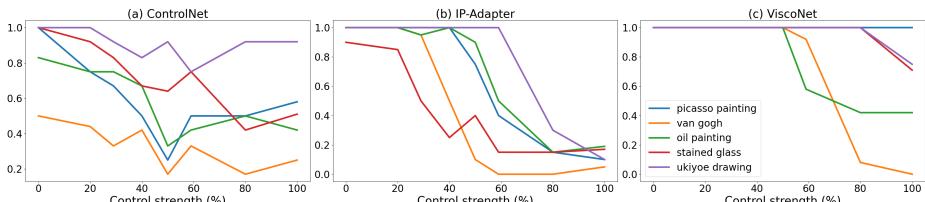


Fig. 26: CLIP accuracy - comparing different image styles.

669

C.2 Section 4.3: Human Evaluation Result

670
671
672

We further conducted a more extensive scale user study on Amazon Mechanical Turk (AMT) to measure the real-life preferences between our model and the HIG baseline approaches. We perform a 4-way comparison, asking workers to select their best preference from randomly shuffled samples, as shown in Figure 27.



Fig. 27: Screenshot of user study presented to users for evaluating the quality of the stylization against the three baselines.

673

670

671

672

Image Styles	Human Evaluation				
	HumanSD	ControlNet	T2I-Adapter	ViscoNet (Ours)	Ours (%)
Ukiyoe	27	32	4	37	37%
Cyberpunk anime	23	13	21	41	41%
Stained glass	0	32	23	45	45%
Van Gogh	2	13	9	76	76%
Picasso	0	13	42	45	45%
Oil Painting	9	11	7	73	73%
Disney	5	23	5	67	67%
Total	77	139	111	384	
Average	9.43%	19.9%	15.9%	54.9%	

Table 3: Our method scores the highest in human evaluation, proving its ability to generate good-quality, diverse image styles.