

Identity-Preserving Image-to-Video Generation via Reward-Guided Optimization

Supplementary Material

To see the dynamic effect of our method and visual comparisons, please refer to our supplementary video. This document includes the following contents:

- Visualizations on Wan 2.2 5B
- Visualizations of multi-person scenarios
- Hacking metric
- Details on truncation gradient step K
- Ablation of loss weights
- Local temporal constraints
- Different face recognition models
- Face score variation with training steps
- Social impacts.

A. Visualizations on Wan 2.2 5B

In the main paper, we provide quantitative comparisons of Wan 2.2 5B with and without our reward optimization. Here we present qualitative visual comparisons, as shown in Fig. 1. Our method achieves better identity consistency preservation compared to the Wan 2.2 5B baseline.



Figure 1. **Visualizations on Wan 2.2 5B.** Our method preserves identity consistency better than Wan 2.2 5B baseline.

B. Visualizations of multi-person scenarios

Compared to injecting facial features directly into the base model, our RL-based approach not only requires no architectural modifications or additional modules, but also naturally generalizes to multi-person scenarios. Concat-ID[†] [4] injects reference face tokens to the video token sequence

and employs 3D self-attention for feature fusion. However, as the number of characters in a scene increases, this method struggles to maintain consistent identity representations across frames. In contrast, as shown in Fig. 2, our method achieves superior identity consistency in complex multi-person settings.

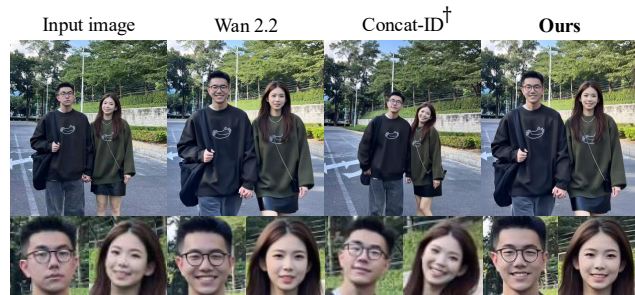


Figure 2. **Visual comparison on multi-person scenes.** Our method preserves identity consistency across frames better than other methods.

C. Hacking metric

We used a vision-language model (VLM) [1] as an auxiliary metric to detect “reward hacking” in generated videos—i.e., over-adherence to the first frame that yields unnaturally rigid faces, suppressed expressions, and poor responsiveness to prompts. The vision-language (VLM) was given the following prompt:

Analyze the following video and evaluate whether it exhibits signs of ‘reward hacking’ in the context of facial consistency optimization. Specifically, determine if the face remains too rigidly consistent with the first frame throughout the video, resulting in unnatural or overly static facial appearance, lack of natural expression changes, minimal facial dynamics, or absence of expected motion (e.g., subtle shifts due to speech, emotion, or camera movement). Consider the following aspects:

- Facial Motion: Is there realistic and natural variation in facial expressions or muscle movements across frames?
- Consistency vs. Stiffness: While the identity remains consistent, does the face appear unnaturally frozen or overly stabilized?

- 055 • Contextual Appropriateness: Given the con- 089
- 056 tent (e.g., talking, emotional expression, head 090
- 057 movement), is the level of facial motion ap- 091
- 058 propriate? 092
- 059 • Visual Artifacts: Are there any signs of 093
- 060 blending, warping, or smoothing artifacts 094
- 061 that suggest aggressive enforcement of sim- 095
- 062 ilarity to the first frame? 096
- 063 Is the video over-optimized (i.e., reward hack- 097
- 064 ing)? Only answer yes or no. 098

065 D. Details on truncation gradient step K

066 We found that as the truncated gradient step size increased, 089

067 the improvement in identity consistency began to level off. 090

Taking computational cost into account, we selected $K=4$. 091

K	1	2	3	4	5
FaceSim \uparrow	0.6593	0.6712	0.6835	0.6942	0.6966

068

069 E. Ablation of loss weights

070 We supervise the model with a combination of face reward 089

071 loss and KL-divergence regularization loss. In the main 090

072 paper, we discuss how KL-divergence regularization pre- 091

073 vents large-scale updates and distribution shift. Here we 092

074 analyze the impact of different loss weights for these two 093

075 components, as shown in Table 1. When the weight ratio is 094

076 $\lambda_2/\lambda_1 = 0$ or $\lambda_2/\lambda_1 = 1$, KL regularization is either negli- 095

077 gible or insufficient, resulting in a high rate of reward hack- 096

078 ing. Conversely, when the weight ratio reaches $\lambda_2/\lambda_1 =$ 097

079 100, KL regularization becomes overly restrictive, prevent- 098

080 ing effective reward-driven model optimization. Based on 099

081 comprehensive evaluation, we adopt $\lambda_1 = 0.1, \lambda_2 = 1$ dur- 100

082 ing training, corresponding to a weight ratio of $\lambda_2/\lambda_1 = 10$.

083

$$\mathcal{L} = \lambda_1 \mathcal{L}_{Reward} + \lambda_2 \mathcal{L}_{KL}, \quad (1)$$

Table 1. Ablation study on loss weights.

λ_2/λ_1	0	1	10	100
FaceSim \uparrow	0.7544	0.7215	0.6942	0.6371
Hacking \downarrow	58%	52%	10%	9%

084 F. Local temporal constraints

085 We have conducted experiments to incorporate local tem- 089

086 poral rewards. However, the results indicate that such con- 090

087 straints such constraints provide limited improvement in our 091

088 task for several reasons. Firstly, the primary challenge is not

089 pixel smoothness but identity fidelity. While VGG-based 090

091 perceptual loss or noise consistency can improve low-level 092

093 pixel smoothness and reduce jitter, they lack the discrimi- 094

095 native power to maintain fine-grained facial identity. Sec- 096

097 ondly, local rewards inherently tend to propagate errors. If 098

099 the identity in frame $t - 1$ starts to degrade, a local tem- 100

100 poral reward will supervise frame t to match that degraded 101

101 identity. In contrast, our face pool acts as a global anchor, 102

102 providing a stable and diverse reference set and covering a 103

103 wide range of angles and expressions. Lastly, current I2V 104

104 foundation models (e.g., Wan 2.2) already possess robust 105

105 temporal attention mechanisms that handle local frame-to- 106

106 frame consistency. 107

	Wan 2.2	LTC	Ours	Ours + LTC
FaceSim \uparrow	0.5780	0.6128	0.6942	0.6961
Time Flickering \uparrow	0.9676	0.9692	0.9690	0.9699

101

102 G. Different face recognition models

103 We employ the ArcFace model [2] as our facial embedding 104

104 extractor and compute the similarity between embeddings 105

105 of generated and ground-truth videos. To further validate 106

106 the robustness of our reward framework, we additionally 107

107 evaluate FaceSim using another face recognition model, 108

108 CurriculumFace [3]. As shown in Table 2, our method 109

109 consistently performs well across different face recognition 110

110 models. 111

Table 2. Different face recognition models.

	Wan 2.2 A14B	W/ reward model
FaceSim-Arc \uparrow	0.5780	0.6942
FaceSim-Cur \uparrow	0.5805	0.6989

110

111 H. Face score variation with training steps

112 As shown in Fig. 3, the face score on the training set 113

113 increases steadily throughout training, indicating that the 114

114 identity reward effectively enhances facial identity preser- 115

115 vation. 116

116 I. Social impacts

117 Our RL-based image-to-video diffusion method, which im- 118

118 proves human identity preservation, offers clear benefits for 119

119 creative production, accessibility, and scientific research. 120

120 However, it may also amplify risks of non-consensual deep- 121

121 fakes and misinformation. To mitigate these risks, we will 122

122 work continuously with legal and ethics experts and im- 123

123 pacted communities to iteratively advance safety measures 124

124 while maintaining the technology’s benefits. 125

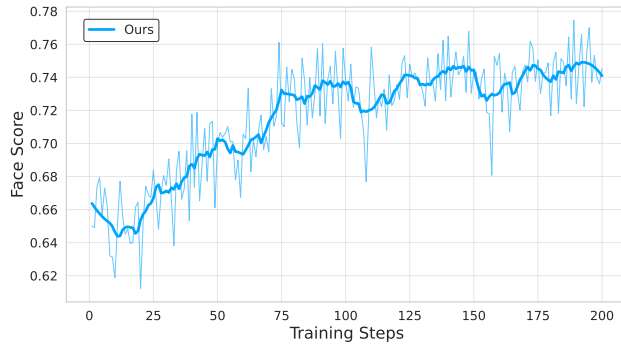


Figure 3. **Face score evolution during training in training set.** Scores increase steadily with training steps, indicating that the face reward effectively improves identity consistency.

125

References

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

- [1] Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasapat, Naveen Sachdeva, Inderjit Dhillon, Marcel Blisstein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025. 1
- [2] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *CVPR*, 2019. 2
- [3] Yuge Huang, Yuhan Wang, Ying Tai, Xiaoming Liu, Pengcheng Shen, Shaoxin Li, Jilin Li, and Feiyue Huang. Curricularface: adaptive curriculum learning loss for deep face recognition. In *CVPR*, 2020. 2
- [4] Yong Zhong, Zhuoyi Yang, Jiayan Teng, Xiaotao Gu, and Chongxuan Li. Concat-id: Towards universal identity-preserving video synthesis. *arXiv preprint arXiv:2503.14151*, 2025. 1