

# V-PartSwap: Motion-Consistent Facial Part Transfer in Videos via Alignment-Aware Diffusion

Anonymous CVPR VGBE submission

Paper ID \*\*\*\*\*

## Abstract

001 *Facial part transfer in videos, such as swapping eyes, eye-*  
002 *brows, nose, or mouth between identities, remains underex-*  
003 *plored compared to full-face video swapping. While image-*  
004 *based methods exist, extending this task to videos is chal-*  
005 *lenging: the transferred part must follow the target’s mo-*  
006 *tion and expression while remaining visually consistent with*  
007 *surrounding regions. Directly injecting reference appear-*  
008 *ance often leads to pasted artifacts or temporal inconsis-*  
009 *tencies. We propose V-PartSwap, a training-free frame-*  
010 *work for reference-guided facial part transfer built upon*  
011 *Diffusion Autoencoders (DiffAE). Our approach performs*  
012 *region-restricted semantic feature blending in the DiffAE*  
013 *encoder to inject reference appearance while preserving*  
014 *non-edited regions. To ensure motion consistency, we align*  
015 *reference parts to the target video using landmark-driven*  
016 *thin-plate spline (TPS) warping and enforce temporal co-*  
017 *herence through region-guided flow-based attention. We*  
018 *also introduce a new benchmark for facial part transfer in*  
019 *videos. Experiments show that our method produces motion-*  
020 *consistent edits with competitive visual quality compared to*  
021 *existing editing methods.*

## 022 1. Introduction

023 Part-level facial editing is increasingly demanded in digital  
024 media production. In particular, facial part transfer in videos  
025 enables applications such as avatar customization and syn-  
026 thetic dataset creation, where users modify a specific region  
027 (e.g., eyes or mouth) while preserving the person’s overall  
028 identity, expression, and motion (Figure 1). Unlike full-face  
029 video swapping, this task requires introducing identity cues  
030 locally while maintaining compatibility with the surrounding  
031 facial appearance and motion.

032 However, facial part transfer in videos remains challeng-  
033 ing. Facial regions are strongly coupled with their surround-  
034 ings in both appearance and motion. Existing image-based  
035 methods, e.g. FuseAnyPart [20], mainly focus on spatial re-  
036 alism and treat frames independently, which often leads to

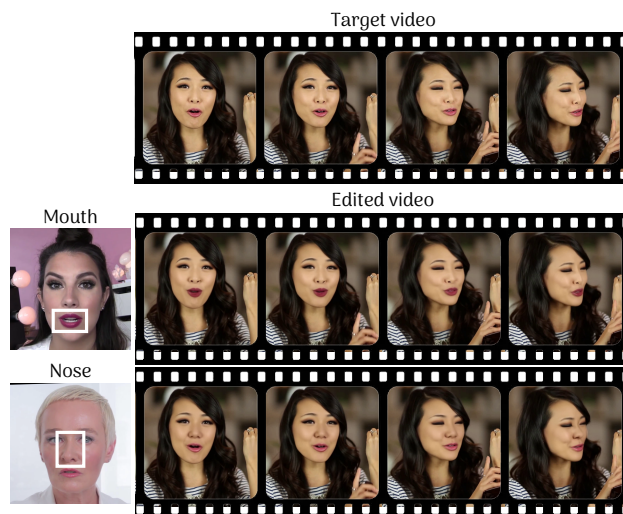


Figure 1. Given a target video (top) and reference images specifying the mouth or nose (left), our method generates edited videos (right) in which the transferred part follows the target’s motion while maintaining a realistic appearance and temporal consistency.

temporally inconsistent edits when applied to videos. Simply injecting reference appearance without considering geometric and motion compatibility can produce pasted-looking artifacts or unstable temporal behavior.

We propose V-PartSwap, an alignment-aware framework for facial part transfer in videos built upon a reconstruction-based diffusion autoencoder (DiffAE) [15]. Our goal is motion-consistent editing, where the transferred appearance follows the target video’s motion while remaining visually consistent with surrounding regions. To achieve this, we inject reference appearance through region-restricted semantic feature blending in a pretrained DiffAE encoder. To ensure geometric compatibility, we further apply landmark-driven thin-plate spline (TPS) warping to adapt the reference part to the target’s per-frame facial configuration. Finally, stochastic boundary blending restores fine details, and optical-flow-guided regional attention enforces temporal consistency.

To enable systematic evaluation, we construct a benchmark based on FaceForensics++ [16] with 100 target identi-

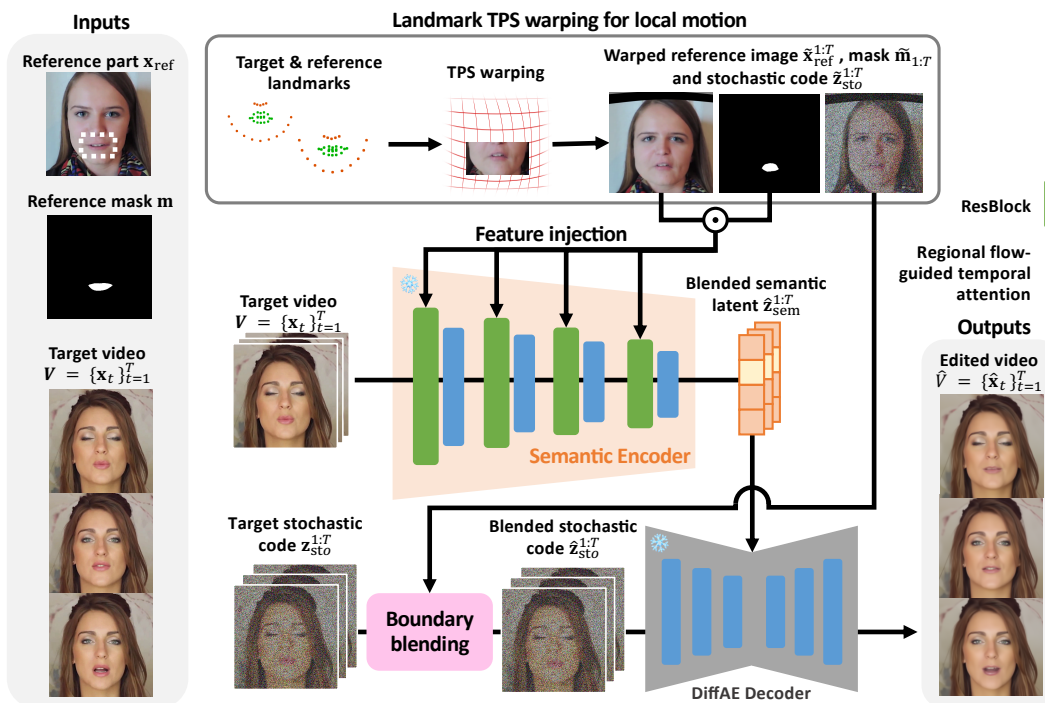


Figure 2. **Overview of V-PartSwap.** Our pipeline combines (i) semantic feature blending, (ii) landmark-guided TPS alignment, (iii) stochastic boundary blending, and (iv) regional flow-guided temporal attention for motion-consistent facial part transfer.

056 ties and reference parts for eyebrows, eyes, nose, and mouth.  
 057 Experiments show that our approach produces motion-  
 058 consistent edits while preserving overall facial realism.

## 059 2. Related Works

060 Facial part transfer aims to replace a specific facial compo-  
 061 nent of a target face using a reference image while preserv-  
 062 ing the remaining facial regions. Existing methods mainly ad-  
 063 dress this task in static images using GAN or diffusion-based  
 064 frameworks, including region-based style transfer, latent  
 065 editing, and part-aware fusion [7, 13, 20, 21, 23]. However,  
 066 these approaches operate on individual frames and do not  
 067 explicitly model temporal consistency.

068 Full-face swapping has been extensively studied in both  
 069 images and videos. Early works rely on GAN-based iden-  
 070 tity injection [2, 14], while recent diffusion-based meth-  
 071 ods improve realism and controllability through condi-  
 072 tional generation [9, 12]. Video extensions further ad-  
 073 dress identity–motion disentanglement and temporal coher-  
 074 ence [1, 8, 18]. In contrast, facial part transfer requires  
 075 localized identity injection while preserving global identity  
 076 and motion.

077 Our method builds on Diffusion Autoencoders (Dif-  
 078 fAE) [15], which enable controllable editing through a struc-  
 079 tured semantic–stochastic latent space.

## 080 3. Method: V-PartSwap

081 Given a target video  $V = \{\mathbf{x}_t\}_{t=1}^T$ , a reference image  $\mathbf{x}_{\text{ref}}$ ,  
 082 and a reference mask  $\mathbf{m}$ , our goal is to synthesize an edited

083 video  $\hat{V} = \{\hat{\mathbf{x}}_t\}_{t=1}^T$  in which the selected facial part matches  
 084 the reference appearance while preserving the target identity  
 085 and motion.

086 As shown in Fig. 2, our pipeline consists of four steps:  
 087 (i) region-aware semantic feature blending in DiffAE, (ii)  
 088 landmark-guided TPS warping to align the reference part  
 089 with the target motion, (iii) boundary-aware blending of the  
 090 stochastic latent code for detail recovery, and (iv) regional  
 091 flow-guided temporal attention for temporal consistency.

### 092 3.1. Semantic Feature Injection

093 **Diffusion autoencoder.** We build on DiffAE [15], which  
 094 decomposes an image into a semantic latent  $\mathbf{z}_{\text{sem}}$  and a  
 095 stochastic latent  $\mathbf{z}_{\text{sto}}$ . The semantic code captures high-  
 096 level facial structure, while the stochastic code preserves  
 097 residual details. For each frame  $\mathbf{x}_t$ , we obtain  $\mathbf{z}_{\text{sem}}^t =$   
 098  $\text{Enc}(\mathbf{x}_t)$ ,  $\mathbf{z}_{\text{sto}}^t = \text{DDIM\_inv}(\mathbf{x}_t, \mathbf{z}_{\text{sem}}^t)$ .

099 **Region-aware semantic feature blending.** To transfer a  
 100 specific facial part while preserving the remaining regions,  
 101 we inject reference semantics only inside the target part mask.  
 102 Let  $\mathbf{h}_t^i$  and  $\mathbf{h}_{\text{ref}}^i$  denote intermediate encoder features of the  
 103 target and reference at resolution  $i$ . Given a resized mask  
 104  $\mathbf{m}^i$ , we perform residual blending

$$105 \tilde{\mathbf{h}}_t^i = \mathbf{h}_t^i + \mathbf{m}^i \odot (\mathbf{h}_{\text{ref}}^i - \mathbf{h}_t^i), \quad (1)$$

106 which injects reference semantics only within the selected  
 107 region. Applying this operation across multiple resolutions  
 108 yields a semantic latent  $\tilde{\mathbf{z}}_{\text{sem}}$  encoding the reference part  
 109 identity.

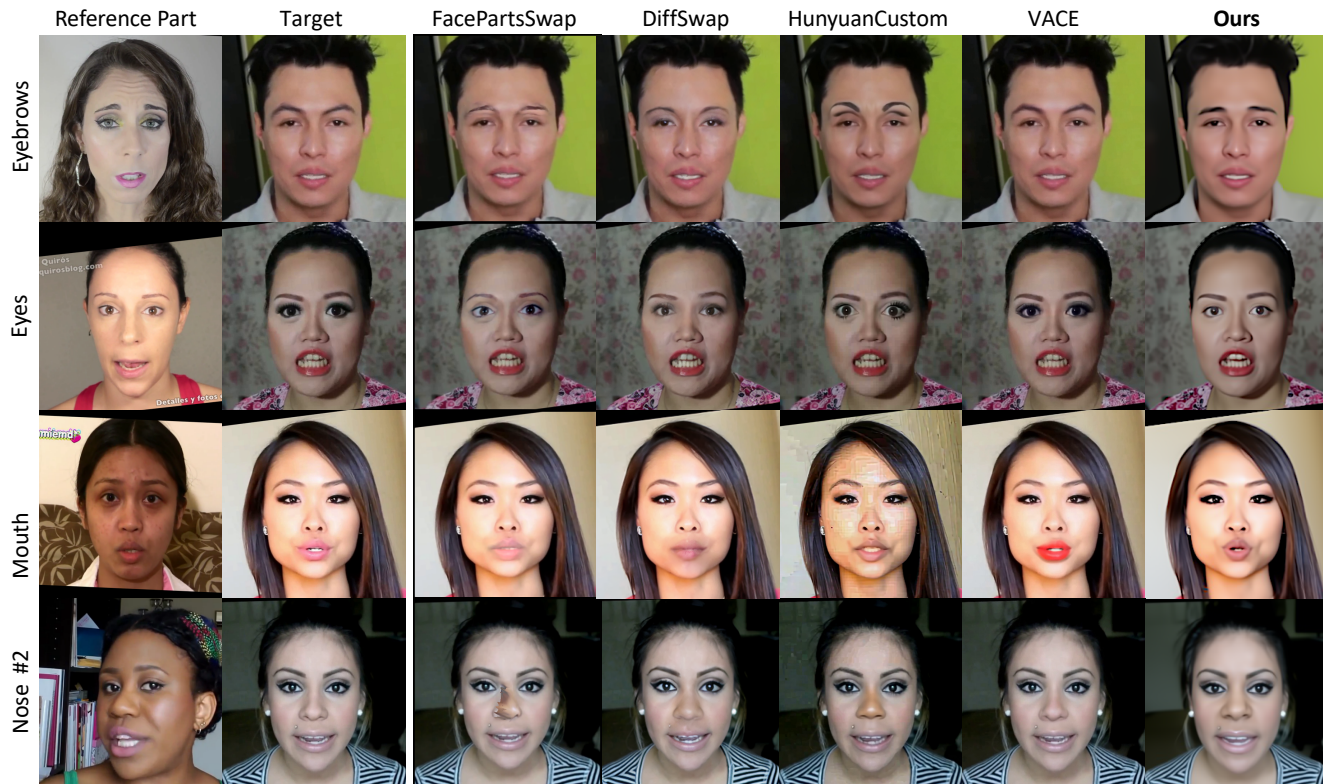


Figure 3. Qualitative comparison of facial-part transfer in videos. Each example shows the target video, reference image, and edited outputs compared with our baselines.

### 110 3.2. Landmark TPS Warping for Local Motion

111 Facial parts such as the mouth exhibit non-rigid motion that  
 112 cannot be captured by global head pose alone. To align  
 113 the reference appearance with the target motion, we warp  
 114 the reference image, mask, and stochastic latent code to  
 115 each frame using landmark-guided thin-plate spline (TPS)  
 116 warping. Let  $\{\mathbf{p}_j\}_{j \in \mathcal{I}}$  and  $\{\mathbf{q}_j^t\}_{j \in \mathcal{I}}$  denote the reference and  
 117 target landmark coordinates for a selected facial region  $\mathcal{I}$   
 118 at frame  $t$ . We estimate a TPS mapping  $\Phi_t(\mathbf{u})$  that aligns  $\mathbf{q}_j^t$   
 119 to  $\mathbf{p}_j$ , producing a sampling grid  $\mathbf{G}_t$ . We warp the reference  
 120 image  $\mathbf{x}_{\text{ref}}$ , mask  $\mathbf{m}$ , and stochastic latent  $\mathbf{z}_{\text{sto}}^{\text{ref}}$  by the grid  
 121  $\tilde{\mathbf{x}}_{\text{ref}}^t, \tilde{\mathbf{m}}_t, \mathbf{z}_{\text{sto}}^{\text{ref} \rightarrow t} = \mathcal{W}(\cdot, \mathbf{G}_t)$ . This alignment transfers the  
 122 target’s local deformation to the reference part, enabling  
 123 motion-consistent editing.

### 124 3.3. Boundary Stochastic Blending

125 Replacing the stochastic latent code inside the entire mask  
 126 can introduce artifacts. Instead, we blend stochastic latents  
 127 only within a narrow band around the warped mask  
 128 boundary. Given the TPS-warped mask, we construct a  
 129 soft boundary weight map  $\tilde{\mathbf{b}}_t$  using morphological opera-  
 130 tions followed by Gaussian smoothing. Let  $\mathbf{z}_{\text{sto},t}^{\text{tgt}}$  denote  
 131 the target stochastic latent and  $\mathbf{z}_{\text{sto}}^{\text{ref} \rightarrow t}$  the TPS-warped refer-  
 132 ence latent. We compute the blended latent as  $\hat{\mathbf{z}}_{\text{sto}}^t =$   
 133  $(1 - \tilde{\mathbf{b}}_t) \odot \mathbf{z}_{\text{sto},t}^{\text{tgt}} + \tilde{\mathbf{b}}_t \odot \mathbf{z}_{\text{sto}}^{\text{ref} \rightarrow t}$ , allowing fine details to be  
 134 restored while maintaining smooth transitions.

### 3.4. Regional Flow-guided Temporal Attention

To improve temporal stability, we apply a flow-guided tem-  
 poral attention module following FLATTEN [3]. The atten-  
 tion is restricted to a region of interest defined by the  
 warped mask, allowing features to be aggregated only from  
 motion-corresponding locations across frames without af-  
 fecting unedited regions.

## 4. Experiments

We describe our evaluation benchmark and then present the  
 experiments.

### 4.1. Experiment Setup

**Evaluation benchmark.** Since no benchmark exists for  
 video facial part transfer, we construct one based on Face-  
 Forensics++ [16]. We select clips containing at least 120  
 consecutive frames and randomly sample 100 videos as tar-  
 gets. For reference images, we curate a pool for four facial  
 parts (eyebrows, eyes, nose, and mouth). Candidate frames  
 are sampled from videos that do not overlap with the tar-  
 get identities. Each frame undergoes face alignment and  
 semantic parsing [24] to extract the corresponding part mask.  
 Quality filters are applied to remove blurred images, trun-  
 cated masks, and extreme expressions. The final reference  
 pool contains 320 aligned samples (80 per part), each con-  
 sisting of an RGB image, part mask, and facial landmarks.

Table 1. Quantitative comparison of the common subset of facial parts (Eye, Mouth, Nose), which are supported by all evaluated methods. ROI-based fidelity metrics (CLIP-I, DINO-I) are macro-averaged over the three parts. Identity and background metrics are computed outside the edited ROI. Temporal and global consistency metrics are computed on the full edited videos corresponding to the same subset. LPIPS<sub>BG</sub> denotes the background LPIPS.  $E_{warp}$  is reported in  $10^{-3}$ .  $\uparrow$  indicates higher is better and  $\downarrow$  indicates lower is better. Best results are **bold**, and second-best results are underlined.

Method	Quality		Motion		Editing Fidelity		Identity	
	$E_{warp} \downarrow$	FVD $\downarrow$	Pose $\downarrow$	Expr. $\downarrow$	CLIP-I $\uparrow$	DINO-I $\uparrow$	ArcFace $\uparrow$	LPIPS <sub>BG</sub> $\downarrow$
FacePartsSwap [5]	9.232	30.723	1.756	0.142	<b>0.579</b>	<b>0.384</b>	0.929	0.160
E4S [13]	9.439	39.249	1.802	0.141	0.568	<u>0.375</u>	0.632	0.215
DiffSwap [23]	9.452	60.043	1.814	<u>0.116</u>	0.563	0.364	0.944	<u>0.0320</u>
FuseAnyPart [20]	9.618	72.524	2.032	0.188	0.563	0.359	0.770	0.0737
VideoSwap [6]	9.351	212.41	2.215	0.195	0.554	0.344	0.482	0.158
VACE [11]	<u>9.232</u>	42.935	<u>1.728</u>	<b>0.0762</b>	0.561	0.357	<b>0.973</b>	<b>0.0261</b>
HunyuanCustom [10]	9.244	<b>23.379</b>	<b>1.727</b>	0.145	0.574	0.366	<u>0.959</u>	0.0441
MoCha [22]	9.626	119.97	2.965	0.289	<u>0.577</u>	0.365	0.0804	0.341
<b>Ours</b>	<b>9.143</b>	<u>27.848</u>	1.919	0.135	0.565	0.371	0.868	0.0693

159 **Implementation Details.** We implement our method on  
160 the DiffAE FFHQ256 model [15]. All videos and reference  
161 images are aligned using Dlib and resized to  $256^2$ . Optical  
162 flow is estimated using RAFT [19]. TPS warping is disabled  
163 for nose transfer, as nasal variation is mainly governed by  
164 global pose and exhibits minimal non-rigid deformation. All  
165 experiments are conducted on NVIDIA A100 GPUs. Editing  
166 a 120-frame video takes approximately 305 seconds. Our  
167 method operates in a fully training-free manner.

168 **Baselines.** We compare with representative facial part  
169 swapping methods, including FacePartsSwap [5], E4S [13],  
170 DiffSwap [23], and FuseAnyPart [20]. We also include  
171 general reference-based video editing systems, including  
172 VideoSwap [6], VACE [11], HunyuanCustom [10], and  
173 MoCha [22]. We use the official implementations and pre-  
174 trained models for all baselines. We do not include the re-  
175 sults of eyebrow transfer for FuseAnyPart [20], as it requires  
176 retraining of the method.

177 **Evaluation Metrics.** We evaluate editing quality along  
178 four aspects. *Temporal and video quality:* Optical-flow war-  
179 ping error ( $E_{warp}$ ) and Fréchet Video Distance (FVD). *Motion*  
180 *preservation:* Pose error (HopeNet [17]), and expression  
181 error (Deep3DFaceRecon [4]). *Editing fidelity:* CLIP-I and  
182 DINO-I similarity between edited parts and reference parts.  
183 *Identity preservation:* ArcFace similarity and background  
184 LPIPS computed outside the edited region.

## 185 4.2. Main Results

186 **Visual comparisons.** Figure 3 presents qualitative com-  
187 parisons with representative baselines. DiffSwap often fails  
188 to fully replace the target attributes, while FacePartsSwap  
189 frequently produces visible boundary artifacts. VACE tends  
190 to generate weak edits for fine-grained attributes, such as  
191 eyebrows and the nose, whereas HunyuanCustom sometimes  
192 produces appearance mismatches due to differences in skin  
193 tone or illumination. In contrast, our method consistently

194 transfers the desired facial parts while maintaining natural  
195 appearance and boundary coherence.

196 **Quantitative comparisons.** Table 1 reports quantitative re-  
197 sults on the common subset of facial parts (eyes, mouth, and  
198 nose). Our method achieves the lowest warping error  $E_{warp}$ ,  
199 indicating stronger motion consistency between the edited  
200 regions and the target video. It also obtains the second-best  
201 FVD score, suggesting improved perceptual video quality.  
202 Despite being training-free, our approach maintains appear-  
203 ance fidelity (CLIP-I and DINO-I) comparable to recent  
204 methods, while preserving identity and background consis-  
205 tency across frames.

206 Image-based part transfer methods such as FacePartsS-  
207 wap and E4S achieve strong appearance fidelity but often in-  
208 troduce boundary artifacts. Video editing approaches such as  
209 VideoSwap and MoCha exhibit larger identity drift, reflected  
210 by lower ArcFace similarity. Diffusion-based part-transfer  
211 methods achieve competitive fidelity but exhibit poorer over-  
212 all quality. Overall, our method achieves stronger motion  
213 alignment while maintaining competitive appearance fidelity  
214 and identity preservation.

## 215 5. Conclusions

216 This paper introduced V-PartSwap, a training-free frame-  
217 work for motion-consistent facial part transfer in videos.  
218 Built upon DiffAE, our method performs localized semantic  
219 blending in the latent space while preserving non-edited fa-  
220 cial regions. By incorporating landmark-guided alignment  
221 and region-guided temporal attention, V-PartSwap preserves  
222 reference appearance while maintaining compatibility with  
223 the target face and temporal consistency across frames.

224 The method does not explicitly model teeth, which may  
225 lead to artifacts when the mouth opens widely, and large ap-  
226 pearance mismatches between reference and target identities  
227 may produce overly smoothed results. Addressing these chal-  
228 lenges is left for future work.

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  
270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285

## References

- [1] Sanoojan Baliah, Yohan Abeysinghe, Rusiru Thushara, Khan Muhammad, Abhinav Dhall, Karthik Nandakumar, and Muhammad Haris Khan. Vface: A training-free approach for diffusion-based video face swapping. In *IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 4315–4324, 2026. 2
- [2] Renwang Chen, Xuanhong Chen, Bingbing Ni, and Yanhao Ge. Simswap: An efficient framework for high fidelity face swapping. In *ACM MM*, pages 2003–2011, 2020. 2
- [3] Yuren Cong, Mengmeng Xu, Christian Simon, Shoufa Chen, Jiawei Ren, Yanping Xie, Bodo Rosenhahn, Tao Xiang, and Sen He. FLATTEN: optical flow-guided attention for consistent text-to-video editing. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2024. 3
- [4] Yu Deng, Jiaolong Yang, Sicheng Xu, Dong Chen, Yunde Jia, and Xin Tong. Accurate 3d face reconstruction with weakly-supervised learning: From single image to image set. In *CVPRW*, 2019. 4
- [5] Claudio Ferrari, Matteo Serpentoni, Stefano Berretti, and Alberto Del Bimbo. What makes you, you? analyzing recognition by swapping face parts. In *Proceedings of IEEE International Conference on Pattern Recognition (ICPR)*, pages 945–951. IEEE, 2022. 4
- [6] Yuchao Gu, Yipin Zhou, Bichen Wu, Licheng Yu, Jia-Wei Liu, Rui Zhao, Jay Zhangjie Wu, David Junhao Zhang, Mike Zheng Shou, and Kevin Tang. Videoswap: Customized video subject swapping with interactive semantic point correspondence. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7621–7630, 2024. 4
- [7] Jingtuo Guo and Yi Liu. Facial parts swapping with generative adversarial networks. *Journal of Visual Communication and Image Representation*, 78:103152, 2021. 2
- [8] Xu Guo, Fulong Ye, Xinghui Li, Pengqi Tu, Pengze Zhang, Qichao Sun, Songtao Zhao, Xiangwang Hou, and Qian He. Dreamid-v: Bridging the image-to-video gap for high-fidelity face swapping via diffusion transformer. *arXiv preprint arXiv:2601.01425*, 2026. 2
- [9] Yue Han, Junwei Zhu, Keke He, Xu Chen, Yanhao Ge, Wei Li, Xiangtai Li, Jiangning Zhang, Chengjie Wang, and Yong Liu. Face-adapter for pre-trained diffusion models with fine-grained id and attribute control. In *Proceedings of European Conference on Computer Vision (ECCV)*, pages 20–36. Springer, 2024. 2
- [10] Teng Hu, Zhentao Yu, Zhengguang Zhou, Sen Liang, Yuan Zhou, Qin Lin, and Qinglin Lu. Hunyuancustom: A multimodal-driven architecture for customized video generation. *arXiv preprint arXiv:2505.04512*, 2025. 4
- [11] Zeyinzi Jiang, Zhen Han, Chaojie Mao, Jingfeng Zhang, Yulin Pan, and Yu Liu. VACE: All-in-one video creation and editing. In *Proceedings of IEEE/CVF International Conference on Computer Vision (ICCV)*, 2025. 4
- [12] Kihong Kim, Yunho Kim, Seokju Cho, Junyoung Seo, Jisu Nam, Kychul Lee, Seungryong Kim, and KwangHee Lee. Diffface: Diffusion-based face swapping with facial guidance. *PR*, 163:111451, 2025. 2
- [13] Zhian Liu, Maomao Li, Yong Zhang, Cairong Wang, Qi Zhang, Jue Wang, and Yongwei Nie. Fine-grained face swapping via regional gan inversion. In *Proceedings of IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8578–8587, 2023. 2, 4 286  
287  
288  
289  
290
- [14] Yuval Nirkin, Yosi Keller, and Tal Hassner. Fsgan: Subject agnostic face swapping and reenactment. In *Proceedings of IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 7184–7193, 2019. 2 291  
292  
293  
294
- [15] Konpat Preechakul, Nattanat Chatthee, Suttisak Widadwongsa, and Supasorn Suwajanakorn. Diffusion autoencoders: Toward a meaningful and decodable representation. In *Proceedings of IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10619–10629, 2022. 1, 2, 4 295  
296  
297  
298  
299  
300
- [16] Andreas Rossler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Nießner. Faceforensics++: Learning to detect manipulated facial images. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1–11, 2019. 1, 3 301  
302  
303  
304  
305
- [17] Nataniel Ruiz, Eunji Chong, and James M. Rehg. Fine-grained head pose estimation without keypoints. In *CVPRW*, 2018. 4 306  
307  
308
- [18] Hao Shao, Shulun Wang, Yang Zhou, Guanglu Song, Dailan He, Zhuofan Zong, Shuo Qin, Yu Liu, and Hongsheng Li. Vividface: A robust and high-fidelity video face swapping framework. In *NeurIPS*, 2025. 2 309  
310  
311  
312
- [19] Zachary Teed and Jia Deng. RAFT: recurrent all-pairs field transforms for optical flow. In *Proceedings of European Conference on Computer Vision (ECCV)*, pages 402–419, 2020. 4 313  
314  
315  
316
- [20] Yaohua Wang, Siying Cui, Aixi Zhang, Wei-Long Zheng, Senzhang Wang, et al. Fuseanypart: Diffusion-driven facial parts swapping via multiple reference images. In *NeurIPS*, 2024. 1, 2, 4 317  
318  
319  
320
- [21] Zhiliang Xu, Xiyu Yu, Zhibin Hong, Zhen Zhu, Junyu Han, Jingtuo Liu, Errui Ding, and Xiang Bai. Facecontroller: Controllable attribute editing for face in the wild. In *Proceedings of AAAI Conference on Artificial Intelligence (AAAI)*, pages 3083–3091, 2021. 2 321  
322  
323  
324  
325
- [22] Zhengbo Xu, Jie Ma, Ziheng Wang, Zhan Peng, Jun Liang, and Jing Li. End-to-end video character replacement without structural guidance. *arXiv preprint arXiv:2601.08587*, 2026. 4 326  
327  
328  
329
- [23] Wenliang Zhao, Yongming Rao, Weikang Shi, Zuyan Liu, Jie Zhou, and Jiwen Lu. Diffswap: High-fidelity and controllable face swapping via 3d-aware masked diffusion. In *Proceedings of IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8568–8577, 2023. 2, 4 330  
331  
332  
333  
334
- [24] Yinglin Zheng, Hao Yang, Ting Zhang, Jianmin Bao, Dongdong Chen, Yangyu Huang, Lu Yuan, Dong Chen, Ming Zeng, and Fang Wen. General facial representation learning in a visual-linguistic manner. In *Proceedings of IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18697–18709, 2022. 3 335  
336  
337  
338  
339  
340