

# VIDEO FACE RE-AGING: TOWARD TEMPORALLY CONSISTENT FACE RE-AGING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Video face re-aging deals with altering the apparent age of a person to the target age in videos. This problem is challenging due to the lack of paired video datasets maintaining temporal consistency in identity and age. Most re-aging methods process each image individually without considering the temporal consistency of videos. While some existing works address the issue of temporal coherence through video facial attribute manipulation in latent space, they often fail to deliver satisfactory performance in age transformation. To tackle the issues, we propose (1) a novel synthetic video dataset that features subjects across a diverse range of age groups; (2) a baseline architecture designed to validate the effectiveness of our proposed dataset, and (3) the development of novel metrics tailored explicitly for evaluating the temporal consistency of video re-aging techniques. Our comprehensive experiments on public datasets, including VFHQ and CelebV-HQ, show that our method outperforms existing approaches in age transformation accuracy and temporal consistency. Notably, in user studies, our method was preferred for temporal consistency by 48.1% of participants for the older direction and by 39.3% for the younger direction.

## 1 INTRODUCTION

Video face re-aging or video-based face re-aging aims to transform the apparent age in facial videos while ensuring the temporal consistency in both age and identity. This field holds significance relevance across diverse domains, including computer graphics, forensics, entertainment, and advertising. Despite the extensive research conducted in this domain, the challenge remains largely unexplored when it comes to videos. One of the remaining challenges is that existing image-based methods yield inconsistent identities when applied to videos or consecutive frames featuring varying expressions, viewpoints, and lighting conditions.

Seminal studies Alaluf et al. (2022); Tzaban et al. (2022) have leveraged StyleGAN-based Karras et al. (2019) frameworks to develop techniques for manipulating facial attributes in videos, aiming for greater attribute consistency. Zoss et al. (2022) have utilized the synthetic images labeled with existing re-aging techniques Alaluf et al. (2021) to curate a paired dataset for re-aging. Their study shows that supervised training on synthetic dataset yields favorable outcomes for still images. Preechakul et al. (2022); Chen & Lathuilière (2023); Li et al. (2023); Wahid et al. (2024) introduce diffusion models to tackle this problem. Lin et al. (2024) employs a reward-based approach by establishing a consensus between the aging process and the aging personalization agents to generate robust faces. However, these approaches are trained on static images and significantly suffer from temporal inconsistencies.

Thus, training face re-aging methods on videos is beneficial for addressing the temporal consistency problem in video re-aging tasks. Therefore, we propose a pipeline to generate a synthetic video dataset. This dataset comprises paired data for supervised training, consisting of various ages, poses, and expressions. The creation of this dataset involves three major steps.

Firstly, we utilize StyleGAN to synthesize a face image. Then we apply existing re-aging method SAM Alaluf et al. (2021) to generate images of the same person with varying ages (Sec. 3.1.1). Next, we build key frames consisting of various poses and expressions of that individual (Sec. 3.1.2). Lastly, we introduce natural motion with frame interpolation (Sec. 3.1.3).

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

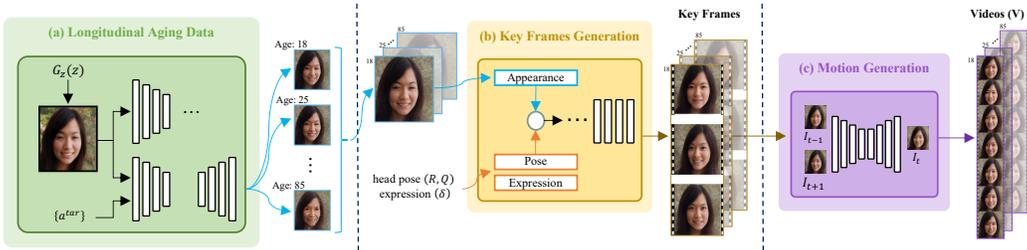


Figure 1: Our proposed pipeline to construct the video dataset for re-aging. Firstly, high-resolution synthetic facial images are created using StyleGAN Karras et al. (2019). Subsequently, images of individuals at different target ages are generated using SAM Alaluf et al. (2021) for age transformation. Next, key frames are produced by employing OSFV, which alters the pose and expression of these synthetic images. This is achieved without relying on driving images, instead using random values for rotation, translation, and expression keypoints. Finally, motion is added to these key frames using FILM Reda et al. (2022), creating smooth and high-fidelity motion videos of subjects at different ages.

In addition to the dataset, we introduce a baseline architecture designed to utilize the temporal coherence inherent in our proposed video dataset. This architecture primarily comprises recurrent blocks, employing a fusion-based approach that leverages concatenated inputs to exploit temporal consistency. Drawing inspiration from seminal works in video generation Clark et al. (2019); Saito et al. (2020); Tulyakov et al. (2018), we incorporate a video discriminator equipped with 3D convolutional layers to ensure both realism and natural motion in the generated videos.

Recognizing that existing aging metrics are not well-suited for video-based methods, we address this gap by developing novel metrics to assess temporal continuity in video-based re-aging. Through extensive experiments, we demonstrate that our video-based architecture produces remarkable results and outperforms existing state-of-the-art methods across various public datasets.

Given the challenges and limitations of current methods, our work introduces several contributions in video re-aging as follows.

1. We introduce a pipeline designed to generate a synthetic video dataset specifically for video re-aging. This dataset features age-paired videos of individuals across various ages, poses, and expressions, allowing model to be trained with supervised learning.
2. We present a baseline network architecture custom-designed for our synthesized video dataset. Our generator is built upon a combination of recurrent blocks consisting of U-Net architecture, utilizing both 2D image-based and 3D video-based conditional discriminators.
3. We propose novel metrics to evaluate the temporal consistency of video re-aging methods. These include Temporal Regional Wrinkle Consistency (TRWC) and Temporally Age Preservation. These metrics provide a robust framework for evaluating the quality of age transformations over time.

## 2 RELATED WORKS

### 2.1 IMAGE-BASED FACE RE-AGING

The study in Antipov et al. (2017b) pioneered the use of a conditional GAN for face aging. Subsequently, several influential works such as Antipov et al. (2017a); Wang et al. (2018); Or-EI et al. (2020); Li et al. (2021); Yao et al. (2021b); Kim et al. (2024); Guo et al. (2024) emerged, expanding on this concept. Similar to Li et al. (2021), SAM Alaluf et al. (2021) emphasizes continuous age progression. SAM is a StyleGAN-based model Karras et al. (2019) capable of generating high-resolution images. In contrast to other methods, it does not use an age classifier to estimate the input age. FRAN Zoss et al. (2022) trains a simple encoder-decoder network in a supervised manner with a generated synthetic dataset. Rather than adopting age classifiers or embedding, it extends the input to a 5-channel image, including two binary masks for input and target ages. CUSP Gomez-Trenado

et al. (2022) introduces style and content encoders to disentangle the style and content of the input. This approach also incorporates the GB algorithm Springenberg et al. (2014) within the CUSP module, ensuring that only age-relevant features are processed. AgeTransGAN Hsu et al. (2022) disentangles the encoded image into identity and age components with two independent modules. PADA Li et al. (2023) and FADING Chen & Lathuilière (2023) adopt a text-driven approach and integrates the pre-trained CLIP Radford et al. (2021) and diffusion models. However, these works only focus on images, and directly applying these methods on individual frames does not consider temporal consistency, affecting the quality of re-aging quality over time.

## 2.2 VIDEO-BASED FACE RE-AGING

Most video face re-aging methods transform the face by manipulating the age in latent space, except Duong et al. (2019) that proposes a reinforcement learning method for the sequence of video frames. Yao et al. (2021a) proposes editing the vectors in the StyleGAN latent space with various pre-processing steps that are independently applied to every single frame of an input video. Tzaban et al. (2022) identify the inconsistencies in PTI Roich et al. (2022) and suggest e4e Tov et al. (2021) encoder with it for finding the pivots. Video editing methods often crop faces as a preprocessing step. To overcome this problem, Yang et al. (2023) addresses this issue by proposing changes in the initial StyleGAN layers to overcome the cropping problem. Kim et al. (2023) proposes a diffusion-based editing approach by using Preechakul et al. (2022) that disentangles the video into time-dependent features (such as motion) that are applied to each frame and time-independent features (such as identity) which are shared across all the frames. While improving temporal consistency over image-based methods, they still struggle to accurately transform faces to the target age.

t

## 2.3 VIDEO FACE RE-AGING DATASET

Video face re-aging presents significant challenges, primarily due to the lack of dedicated video datasets. Existing image-based face re-aging datasets are typically labeled either automatically using an age classifier or manually via crowdsourcing Zhang et al. (2022); Liu et al. (2021); Karras et al. (2019); Rothe et al. (2018). However, these datasets do not provide paired ages for supervised training. Zheng et al. (2017) proposed a technique to generate a synthetic dataset that is labelled through SAM Alaluf et al. (2021). Duong et al. (2019) also labelled a video dataset, which remains private.

# 3 VIDEO FACE RE-AGING FRAMEWORK

We first describe the pipeline of synthesizing the proposed video dataset. This is followed by introducing our baseline architecture and loss functions. Lastly, we present novel metrics specifically designed for video re-aging methods to quantify their re-aging performance over time.

## 3.1 RE-AGING VIDEO DATASET

### 3.1.1 IMAGE-BASED FACE RE-AGING DATASET

Creating a high-quality image-based re-aging dataset is a crucial in our pipeline because it directly influences the quality of our subsequent video re-aging dataset. However, obtaining these paired image datasets is a challenging task. To overcome this limitation, we turn to insights from Zoss et al. (2022), which show that training on synthetic datasets can yield realistic results on real images. For instance, the neural network can learn how wrinkles change from the synthetic images and apply this knowledge to real images. We refer to these learned changes as delta images  $D_t$ , as illustrated in Fig. 2.

Firstly, we leverage StyleGAN Karras et al. (2020) which takes a random noise as input and generate high-resolution synthetic image. Then we utilize SAM Alaluf et al. (2021) that manipulates the latent vector with StyleGAN for age transformation. The intuition behind choosing SAM over other existing methods is its re-aging performance in terms of age error, that is also evident through our experiments. Using this approach, we construct a synthetic facial image dataset:

$$I^{tar} = SAM(G_z(z); a^{tar}) \quad (1)$$

Given a random noise  $z$ , StyleGAN  $G_z(\cdot)$  produces a sample image, which is then used as input for SAM Alaluf et al. (2021) along with a target age  $a^{tar}$ . As a result we get an image  $I^{tar}$  with apparent age  $a^{tar}$  as shown in Fig. 1 (a). Readers can refer to Zoss et al. (2022) for more details about image-level re-aging process.

### 3.1.2 KEY FRAMES GENERATION

The subsequent step in the outlined pipeline involves key frames generation. This key frames capture specific snapshots or moments in videos to enable seamless transitions between sequences. The challenge lies in acquiring diverse facial images to adequately encapsulate the video’s dynamics, including various poses and expressions. To address this, we utilize a recent face reenactment technique Wang et al. (2021a) that modifies the pose and expression of a source image based on a driving image. Incorporating this method into our pipeline allows us to generate multiple images of an individual with varying poses and expressions. These resulting images serve as key frames.

In this study, we employ the off-the-shelf model OSFV<sup>1</sup> Wang et al. (2021a) to produce synthetic key frames using our image dataset. Our approach diverges from the original methodology in that we rely solely on source images from our dataset. Instead of utilizing driving images, which can be challenging to gather, we employ random values for the rotation matrix  $R$ , translation matrix  $Q$ , and expression keypoints  $\delta$  to generate various poses and expressions as shown in Fig. 1 (b). We can write this process as follows:

$$K = G_{kp}(I^{tar}, R, Q, \delta), \quad (2)$$

where  $K$  represents generated key frames and  $G_{kp}$  denotes key frames generator. We generate eight different key frames through this method. While it is also possible to repeat this procedure to create a motion video, but we have observed a degradation of the quality in the resultant videos. Therefore, we come up with an alternative approach to address this problem in next section.

### 3.1.3 MOTION GENERATION

The final step of our pipeline is motion generation. We leverage the recent work in frame interpolation methods to ensure smooth and high-fidelity motion. Specifically, we employ the method presented in Reda et al. (2022) to the eight key frames generated in the previous step by recursively generating intermediate frames between them. This iterative process is executed for two consecutive frames and is repeated for every subject across all ages as follows:

$$I_t = FI(I_{t-1}, I_{t+1}), \quad (3)$$

where  $FI$  is motion generation network Reda et al. (2022) and  $I_t$  is the frame at  $t$  time-step. As a result, we obtain smooth and high-fidelity paired videos for every subject of different ages as shown in Fig. 1 (c).

## 3.2 NETWORK ARCHITECTURE

### 3.2.1 GENERATOR

To fully utilize the advantage of our novel video re-aging dataset, we carefully design the generation scheme to accurately transform faces over time. Following the approach in Zoss et al. (2022), we adopt an alternative method for incorporating input and target ages into our model structure. Recent methods Gomez-Trenado et al. (2022); Hsu et al. (2022) tend to use pre-trained age classifiers to estimate the input age. We concatenate the input frame  $I_t$  at  $t$  time-step with two spatial masks (one for input age and one for target age) over channel dimensions. These two masks contain constant

<sup>1</sup>Note that we trained OSFV on VFHQ Xie et al. (2022) dataset for a resolution of  $512 \times 512$ . The training takes 30 days on 8 A100 GPUs. Additional details are provided in the supplementary materials.

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229

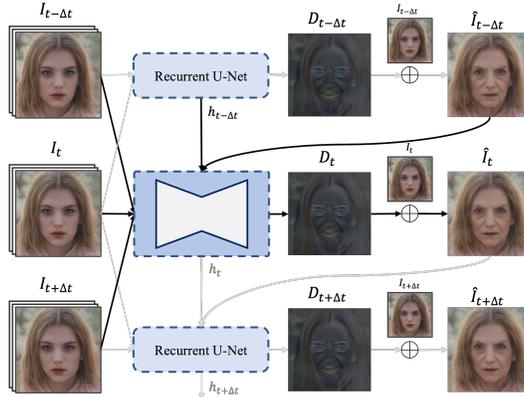


Figure 2: Overview of our generator for video re-aging.

230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241

values representing input and output ages, ranging from 0 to 100 normalized between 0 and 1. This results in a 5-channel masked frame  $I_t^{mask}$ . Mathematically,

$$I_t^{mask} = [I_t, M^{inp}, M^{tar}], \quad (4)$$

where  $M^{inp}$  and  $M^{tar}$  are spatial masks for input and target ages.  $[\cdot]$  denotes the channel-wise concatenation operation. If  $V^{mask} = \{I_1^{mask}, I_2^{mask}, \dots, I_N^{mask}\}$  refers to any input masked sequence  $V^{mask}$  for  $N$  total number of frames, our generator  $G$  produces the output video  $\hat{V}$ .

242  
243

$$\hat{V} = G(V^{mask}), \quad (5)$$

244  
245  
246  
247  
248  
249  
250  
251

The proposed  $G$  adopts a recursive scheme to consider the temporal information of videos, inspired by Wu et al. (2022); Tian et al. (2021); Zhu et al. (2023), namely recurrent block function  $RB(\cdot)$ . We employ the U-Net architecture within the  $RB$  function. Please refer to Sec. A for the details of U-Net network architecture. This  $RB$  stacks the multiple arguments. First, we concatenate the consecutive input frames  $[I_{t-\Delta t}^{mask}, I_t^{mask}, I_{t+\Delta t}^{mask}]$ . Here,  $I_{t-\Delta t}^{mask}$  and  $I_{t+\Delta t}^{mask}$  refers to the two adjacent frames of  $I_t^{mask}$  with interval step  $\Delta t$ . We further concatenate the resulting stacked frames with previous hidden state  $h_{t-\Delta t}$  and previous output frame  $\hat{I}_{t-\Delta t}$ . Mathematically,

252  
253

$$h_t, D_t = RB([I_{t-\Delta t}^{mask}, I_t^{mask}, I_{t+\Delta t}^{mask}, \hat{I}_{t-\Delta t}, h_{t-\Delta t}]) \quad (6)$$

254  
255  
256

As a result, we obtain hidden state  $h_t$  and delta image  $D_t$ . Once the delta image is attained, we can easily obtain output frame  $\hat{I}_t$  through element-wise summation between input  $I_t$  and  $D_t$ :

257  
258

$$\hat{I}_t = D_t + I_t. \quad (7)$$

259  
260  
261  
262  
263

We have now processed three frames to get the re-aged output of the middle frame. For the remaining consecutive frames, we can pass the hidden state  $h_t$  and output frame  $\hat{I}_t$  along with the next consecutive frames in Eq. 6 for the next iteration. This process is repeated for all  $N$  frames of the input video  $V^{mask}$  to generate the output video  $\hat{V}$ , as illustrated in Fig. 2.

264  
265

### 3.2.2 DISCRIMINATOR

266  
267  
268  
269

In addition to the image discriminator, we introduce a video discriminator to assess the consistency of age-related high frequency details across consecutive frames. We observed that using only image discriminator produces inconsistent output, leading to flickering. The comparison is presented in our ablation studies. The image discriminator employs a PatchGAN architecture Isola et al. (2017) to differentiate realistic images from synthetic ones, while for the video discriminator, we adopt a

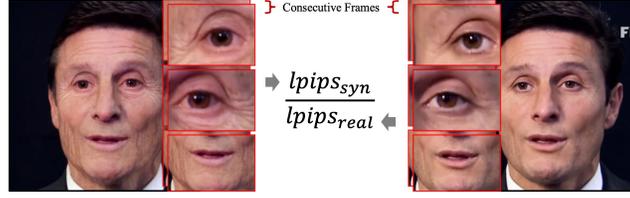


Figure 3: Conceptual overview of proposed TRWC.

spatio-temporal convolutional network that utilizes 3D convolution layers, similar to the approach used in previous works such as Vondrick et al. (2016) and Tulyakov et al. (2018). Specifically, we exploit the spatio-temporal information from consecutive generated frames in the form of delta images  $D_t$  as inputs to important details over time without depending on facial images.

### 3.2.3 LOSS FUNCTIONS

We use  $L1$  loss, LPIPS loss Zhang et al. (2018), and adversarial loss functions Lim & Ye (2017). We also use adversarial loss functions for image and video discriminators  $\mathcal{L}_{adv,I}$  and  $\mathcal{L}_{adv,V}$ . Our total objective is obtained as:

$$\mathcal{L} = \lambda_{L1} \mathcal{L}_{L1}(\hat{V}, V_{gt}) + \lambda_{adv,V} \mathcal{L}_{adv,V}(\hat{V}, M^{tar}) + \lambda_{adv,I} \mathcal{L}_{adv,I}(\hat{V}, M^{tar}) + \lambda_p \mathcal{L}_{LPIPS}(\hat{V}, V_{gt}),$$

where  $V_{gt}$  is the ground-truth video corresponding to the target age obtained from dataset building pipeline. Here, we set the lambda values as  $\lambda_{L1} = 1.0$ ,  $\lambda_{adv,I} = 0.025$ ,  $\lambda_{adv,V} = 0.025$ ,  $\lambda_p = 1.0$ .

## 3.3 PROPOSED METRICS

Evaluating age transformation performance in a video context is both essential and challenging. According to the best of our knowledge, there is no specific metric to assess the temporal consistency of re-aging methods. Most existing evaluation metrics are designed for image-based re-aging methods as they measure the age transformation performance on individual images Zoss et al. (2022) or use metrics designed to evaluate the continuity of identity drift Alaluf et al. (2022). Regarding this, we propose two metrics, TRWC and T-Age, capable of effectively assessing temporal continuity in terms of age in video-based re-aging task.

### 3.3.1 TEMPORAL REGIONAL WRINKLE CONSISTENCY (TRWC)

To evaluate the consistency of video re-aging methods, we need to consider the aging-related details in facial areas such as wrinkles where it matters. This is because these regions are most vulnerable and simultaneously important in terms of human perception. Therefore, we focus on aging-related facial areas, especially Crow’s feet and Nasolabial folds, motivated from Wang et al. (2021b); Li et al. (2022). These areas are more likely to show signs of aging because they are used a lot for facial expressions and speech. We extract the region-of-interest (ROI) around these areas to observe changes over time. However, when analyzing perceptual differences between frames, factors such as facial angles and expressions can be sensitive. To address this, we calculate the LPIPS of the synthesized image and normalize it using the LPIPS of the real image, ensuring we consider differences in each image over time. Fig. 3 explains the conceptual schematic of TRWC. Mathematically, we define TRWC as follows.

$$\text{TRWC}_{\Delta t} = \frac{1}{3(N - \Delta t)} \sum_{r=1}^3 \sum_{t=1}^{N-\Delta t} \frac{\text{lpiPs}(f_{roi}^r(\hat{I}_t), f_{roi}^r(\hat{I}_{t+\Delta t}))}{\text{lpiPs}(f_{roi}^r(I_t), f_{roi}^r(I_{t+\Delta t}))}, \quad (8)$$

Table 1: Quantitative comparison on CelebV &amp; VFHQ datasets. The best results are highlighted in bold.

Dataset	Models	Young $\rightarrow$ Old		Old $\rightarrow$ Young	
		TRWC <sub>1</sub> ↓	T-Age↓	TRWC <sub>1</sub> ↓	T-Age↓
CelebV-HQ	AgeTransGAN	4.25	1.26	3.12	2.29
	CUSP	-	-	1.90	2.19
	FADING	4.21	3.20	2.23	3.14
	<b>OURS</b>	<b>1.38</b>	<b>0.84</b>	<b>0.70</b>	<b>1.03</b>
VFHQ	AgeTransGAN	4.26	1.61	2.92	2.04
	CUSP	-	-	1.76	2.14
	FADING	3.79	2.77	2.04	2.47
	<b>OURS</b>	<b>1.35</b>	<b>1.16</b>	<b>0.69</b>	<b>1.34</b>

where  $N$  is the number of frames and  $\Delta t$  is time interval between consecutive frames. The function  $f_{roi}^r(\cdot)$  is ROI function that are obtained by facial landmark detector Deng et al. (2018). Here, we consider only three ROI regions, left-eye, right-eye, and mouth, indexed by  $r$ .

### 3.3.2 TEMPORAL-AGE (T-AGE) PRESERVATION

Drawing inspiration from the TL-ID metric Tzaban et al. (2022), which proposes metrics for identity consistency in videos, we introduce T-Age for video re-aging. T-Age measures the age difference between two adjacent frames using cosine similarity, utilizing an off-the-shelf age classifier Rothe et al. (2015). A lower T-Age value indicates a more consistent age representation across the frames.

## 4 EXPERIMENTS

### 4.1 EXPERIMENT SETUP

In this section, we show the superiority of our method by comparing with all existing state-of-the-art re-aging methods. This includes HRFAE Yao et al. (2021b), SAM Alaluf et al. (2021), FRAN Zoss et al. (2022), CUSP Gomez-Trenado et al. (2022), AgeTransGAN Hsu et al. (2022), Diffusion AE Preechakul et al. (2022), STIT Tzaban et al. (2022), StyleGANEX Yang et al. (2023), and Diffusion VAE Kim et al. (2023), FADING Chen & Lathuilière (2023) totaling 10 methodologies.

We have trained our methods on the proposed synthetic videos generated through our pipeline. We have shared the sample images of our generated dataset in Fig. 7. For the test set, we choose the CelebV-HQ Zhu et al. (2022) and VFHQ dataset Xie et al. (2022) as video test set. We consider three target age groups, (18, 25, 35) with input age as 85 for *Old*  $\rightarrow$  *Young* task and three age groups, (65, 75, 85) with input age as 18 for *Young*  $\rightarrow$  *Old* task. We provide the additional details in the supplementary materials (Sec. A).

### 4.2 METRICS

We evaluate the performance of re-aging models based on their ability to transform to the target age while ensuring temporal consistency in the re-aged image. Given that we do not have access to ground-truth for real videos, we employ metrics that do not rely on ground-truth. Specifically, we use five metrics to evaluate the results. We calculate mean absolute error (MAE) between the estimate ages, computed with the pre-trained age classifier to quantify the age transformation quality whereas temporal consistency is measured by TRWC and T-Age. Despite knowing that many existing methods employ DEX that might bias the test results, we opted for DEX due to its accuracy and stability. However, it’s important to note that we did not use any age estimation network in our training process. Lastly, we use TGID Tzaban et al. (2022) to evaluate the global identity similarity.

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

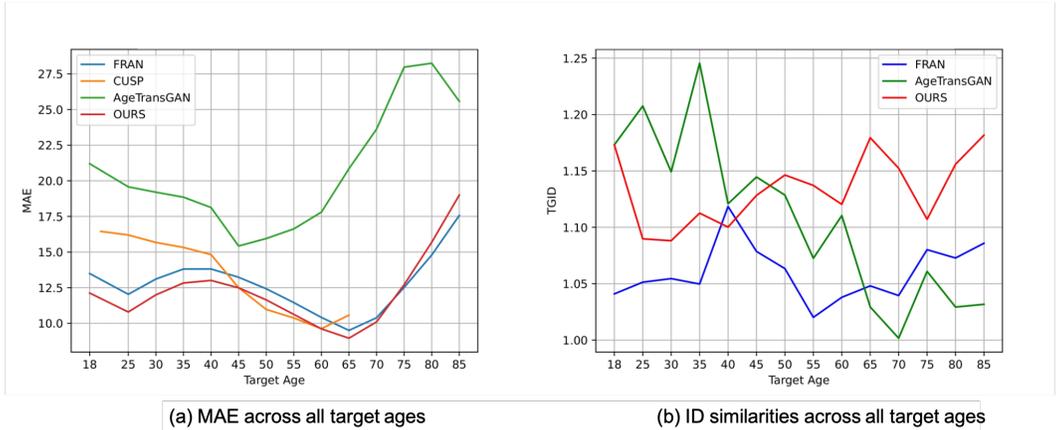


Figure 4: Comparison of (a) Mean Absolute Error (MAE) and (b) ID Similarity (TGID) across a range of target ages for different methods. Lower MAE and higher TGID values indicate better performance.

Table 2: User Study Results. Participants evaluated the methods according to four criteria: Age Accuracy (AA), Identity Preservation (IP), Temporal Consistency (TC), and Overall Naturalness (ON). The highest-scoring results are highlighted in bold. Scores are in percentages (%).

Models	Young → Old				Old → Young			
	AA↑	IP↑	TC↑	ON↑	AA↑	IP↑	TC↑	ON↑
AgeTransGAN	10.56	11.80	11.14	12.38	22.39	31.83	27.06	25.23
CUSP	21.81	<b>46.07</b>	21.02	29.02	12.56	14.31	11.82	12.85
<b>OURS</b>	<b>67.64</b>	42.13	<b>67.84</b>	<b>58.60</b>	<b>65.05</b>	<b>53.86</b>	<b>61.12</b>	<b>61.92</b>

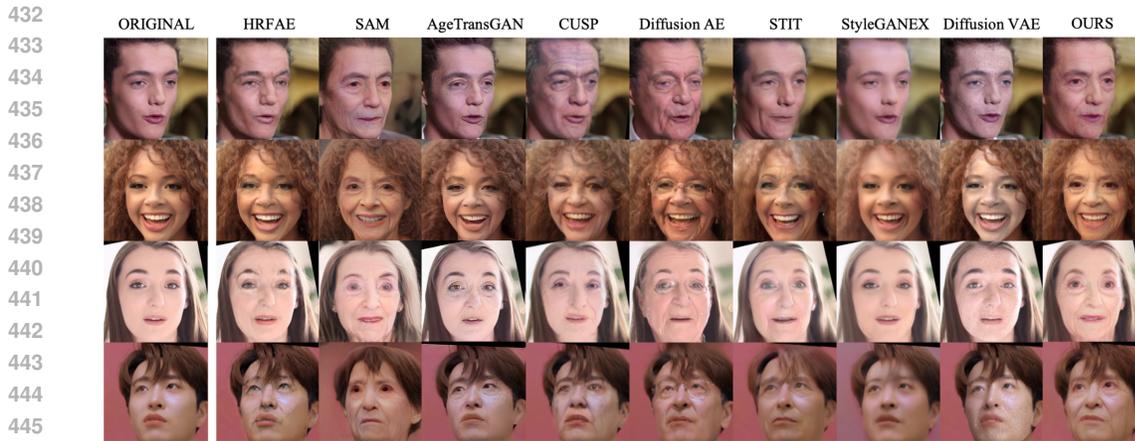
## 5 COMPARISON RESULTS

### 5.1 QUANTITATIVE RESULTS

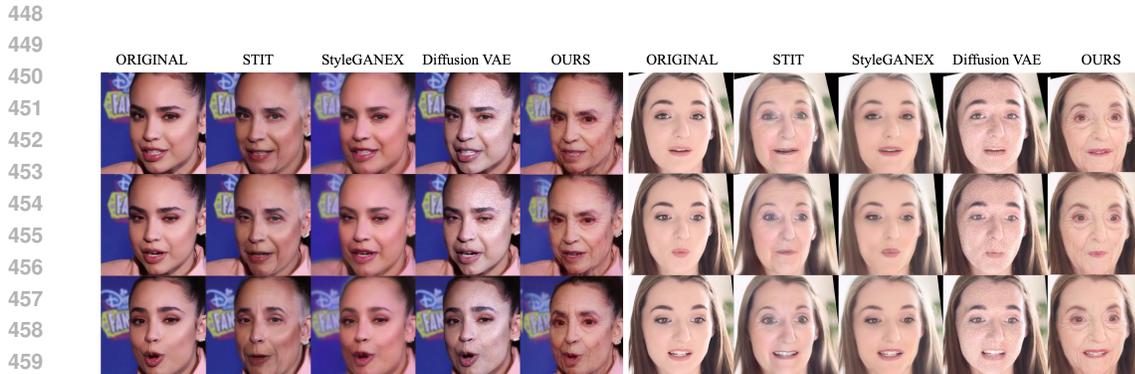
We quantitatively compare our method with state-of-the-art methods in Table 1. Note that we omit the CUSP result for *Young* → *Old* due to its maximum target age limitation (65) Gomez-Trenado et al. (2022). Our results indicate that model trained on videos show higher TRWC, suggesting that it maintains greater temporal stability and exhibits lower perceptual differences between the consistent frames. On other side, When evaluating age transformation continuity using T-AGE, our method outperforms the other methods. Furthermore, it shows that diffusion models, i.e., FADING Chen & Lathuilière (2023), are also temporally inconsistent. These results suggest that these metrics’ performances are strongly correlated with the target ages. Therefore, we show the MAE for all the target ages in Fig. 4 (a). Furthermore, conserving the identity is important when changing the age. Therefore, by following Tzaban et al. (2022), we compare the identity similarity in Fig. 4 (b). The trends show that our methods preserve the identity when changing the ages, even for older ages, as other datasets are skewed towards younger ages. We can observe that the age transformation performance of our method is consistently better than other state-of-the-art methods, despite being trained on synthetic videos. Additionally, our method also improves overall temporal consistency. It is worth mentioning that all the results we presented are in line with user studies (Table. 2), which further affirm our newly proposed metrics.

### 5.2 USER STUDY

We conducted a subjective evaluation via a user study, presenting 15 sequences to a total of 43 anonymous participants for the older direction and 36 participants for the younger direction. This study compares our method with state-of-the-art techniques such as CUSP Gomez-Trenado et al. (2022), and AgeTransGAN Hsu et al. (2022) on *Young* → *Old* and *Old* → *Young* tasks. As illus-



447 Figure 5: Qualitative comparison with existing state-of-the-art methods. The target age is set to 85.



461 Figure 6: Comparison with video editing methods.

462  
463  
464  
465  
466  
467  
468  
469  
470

trated in Table. 2, a clear majority of participants found our method superior, particularly in terms of age accuracy, temporal consistency, and naturalness. These findings highlight the effectiveness of our video-based approach in enhancing both temporal consistency and age transformation capabilities, thereby demonstrating that proposed metrics are reliable indicators of cognitive temporal consistency. For example, the significance of a 0.18 ( $\pm 1$ ) difference in TRWC in Table. 1 is evident by the users' choices.

### 471 5.3 QUALITATIVE RESULTS

472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

We present our qualitative comparison with the existing state-of-the-art in Fig. 5. It is evident that SAM is effective in age transformation, while STIT and StyleGANEX struggle in this aspect. However, SAM fails to preserve attributes such as identity, pose, and expression because it strongly depends on age classifier loss, which is proficient at altering age but not at maintaining other key attributes. On the other hand, STIT and StyleGANEX, which utilize average vectors related to age changes, fail to properly consider the individual samples. As a result, they encounter difficulties in age transformation, particularly when dealing with wild samples. In such cases, their methods become more susceptible due to their strong reliance on trained data for constructing the latent space. Diffusion AE achieves successful age transformation but introduces artifacts such as glasses, leading to identity drift. While the Diffusion VAE works for small changes, it fails entirely in age transformation when dealing with a significant age gap for all subjects. Image-based methodologies, such as HRFAE and AgeTransGAN, inherently exhibit lower age transformation capabilities in wild scenarios. In contrast, while CUSP achieves successful age transformation in certain samples, it often tends to produce artifacts. In contrast, our method successfully performs age transformation while maintaining stable results in terms of identity, expression, and pose. While FRAN also achieves sim-

Table 3: Ablation studies: We ablate the (Left) key frame generation and (Right) different frame interpolation methods to select the highest performance method based on CPBD and warping error.

Method	MAE↓		ID↑		Methods	DQBC	EMA-VFI	FILM
	18	85	18	85				
StyleHEAT	8.06	4.41	0.64	0.50	CPBD↑	0.5115	0.4213	<b>0.6061</b>
StyleMask	7.23	9.98	0.61	0.37	Warp Error↓	0.00118	0.00119	<b>0.00104</b>
OSFV	<b>1.62</b>	<b>0.71</b>	<b>0.90</b>	<b>0.90</b>				

ilar quality, our results possess more natural wrinkle details, with the differences being particularly pronounced around the eyes.

Furthermore, Fig. 6 presents a performance comparison of different frames to demonstrate the efficacy of our approach against other video editing methods. It can be observed that, despite producing temporally consistent results, their ability to transform age is quite limited. Overall, these comparison results indicate that leveraging a video dataset enables our method to outperform the existing methods, yielding natural age progression and high-quality results, even under extreme conditions. We share more results in Sec. B and discuss the limitations of our work in D.

## 6 ABLATION STUDY

In this section, we ablate the selection of the modules of proposed pipeline except the choice of SAM Alaluf et al. (2021) which is already validated in Zoss et al. (2022).

### 6.0.1 KEY FRAME GENERATION

In Table 3 (a), we compare the performance of three reenactment methods: OSFV Wang et al. (2021a), StyleMask Bounareli et al. (2023), and StyleHeat Yin et al. (2022), for key frame generation. The best method is selected for our re-aging task based on lower MAE and higher identity similarity (ID). For ID, we calculate the cosine similarity between the vectors of the face image, extracted using ArcFace Deng et al. (2019). Given different pose and expressions, it is crucial to maintain the age of a person without losing its identity. The results show that OSFV Wang et al. (2021a) successfully preserves both the identity and age of the person, while other methods fail to maintain the person’s identity and age. We further investigate the training configurations of Wang et al. (2021a) in Sec. C.

### 6.0.2 MOTION GENERATION

In this ablation study, we experiment with state-of-the-art frame interpolation models Reda et al. (2022); Zhang et al. (2023b); Zhou et al. (2023) to determine the most effective motion generation method. We evaluate image quality and motion consistency using non-reference-based metrics such as cumulative probability of blur detection (CPBD) Narvekar & Karam (2009; 2011) for image quality and warping error Zhang et al. (2023a) for motion consistency. Table 3 (b) suggests that all methods exhibits lower error in terms of temporal consistency. However, FILM Reda et al. (2022) outperforms other models by producing sharp results, while the outputs of the remaining models tend to be blurry.

## 7 CONCLUSION

In this paper, we introduced a paired video dataset for the video re-aging. This dataset encompasses subjects from a wide range of age groups. Shifting from the conventional model-centric focus, we adopted a data-centric approach. Inspired by recent advancements in face reenactment and frame interpolation, we encompasses various facial poses and expressions in the proposed data. Consequently, we proposed a baseline network architecture to evaluate the proposed dataset, emphasizing both temporal consistency and the quality of age transformations. Furthermore, we formulated two novel metrics to evaluate temporal consistency in video re-aging, which consider age-relevant features such as facial wrinkles over time. We validated our method with comprehensive experiments.

## REFERENCES

- 540  
541  
542 Yuval Alaluf, Or Patashnik, and Daniel Cohen-Or. Only a matter of style: Age transformation using  
543 a style-based regression model. *ACM Transactions on Graphics (TOG)*, 40(4):1–12, 2021.
- 544 Yuval Alaluf, Or Patashnik, Zongze Wu, Asif Zamir, Eli Shechtman, Dani Lischinski, and Daniel  
545 Cohen-Or. Third time’s the charm? image and video editing with stylegan3. In *European Con-  
546 ference on Computer Vision*, pp. 204–220. Springer, 2022.
- 547 Grigory Antipov, Moez Baccouche, and Jean-Luc Dugelay. Boosting cross-age face verification  
548 via generative age normalization. In *2017 IEEE International Joint Conference on Biometrics  
549 (IJCB)*, pp. 191–199. IEEE, 2017a.
- 550 Grigory Antipov, Moez Baccouche, and Jean-Luc Dugelay. Face aging with conditional generative  
551 adversarial networks. In *2017 IEEE international conference on image processing (ICIP)*, pp.  
552 2089–2093. IEEE, 2017b.
- 553 Stella Bounareli, Christos Tzelepis, Vasileios Argyriou, Ioannis Patras, and Georgios Tzimiropou-  
554 los. Stylemask: Disentangling the style space of stylegan2 for neural face reenactment. *IEEE  
555 Conference on Automatic Face and Gesture Recognition*, 2023.
- 556 Xiangyi Chen and Stéphane Lathuilière. Face aging via diffusion-based editing. *arXiv preprint  
557 arXiv:2309.11321*, 2023.
- 558 Aidan Clark, Jeff Donahue, and Karen Simonyan. Adversarial video generation on complex datasets.  
559 *arXiv preprint arXiv:1907.06571*, 2019.
- 560 Jiankang Deng, Anastasios Roussos, Grigorios Chrysos, Evangelos Ververas, Irene Kotsia, Jie Shen,  
561 and Stefanos Zafeiriou. The menpo benchmark for multi-pose 2d and 3d facial landmark local-  
562 ization and tracking. *IJCV*, 2018.
- 563 Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin  
564 loss for deep face recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision  
565 and Pattern Recognition*, pp. 4690–4699, 2019.
- 566 Chi Nhan Duong, Khoa Luu, Kha Gia Quach, Nghia Nguyen, Eric Patterson, Tien D Bui, and  
567 Ngan Le. Automatic face aging in videos via deep reinforcement learning. In *Proceedings of the  
568 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10013–10022, 2019.
- 569 Guillermo Gomez-Trenado, Stéphane Lathuilière, Pablo Mesejo, and Óscar Cordón. Custom struc-  
570 ture preservation in face aging. In *European Conference on Computer Vision*, pp. 565–580.  
571 Springer, 2022.
- 572 Yingchun Guo, Mingyuan Li, and Gang Yan. Identity-preserving face aging with multi-attribution  
573 fusion and multi-scale attention. *IEEE Signal Processing Letters*, 2024.
- 574 Gee-Sern Hsu, Rui-Cang Xie, Zhi-Ting Chen, and Yu-Hong Lin. Agetransgan for facial age trans-  
575 formation with rectified performance metrics. In *European Conference on Computer Vision*, pp.  
576 580–595. Springer, 2022.
- 577 Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with  
578 conditional adversarial networks. *Proceedings of the IEEE/CVF Conference on Computer Vision  
579 and Pattern Recognition*, 2017.
- 580 Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative  
581 adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and  
582 Pattern Recognition*, pp. 4401–4410, 2019.
- 583 Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyz-  
584 ing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF Conference on  
585 Computer Vision and Pattern Recognition*, pp. 8110–8119, 2020.
- 586 Bumsoo Kim, Abdul Muqet, Kyuchul Lee, and Sanghyun Seo. Toonaging: Face re-aging upon  
587 artistic portrait style transfer. *arXiv preprint arXiv:2402.02733*, 2024.
- 588  
589  
590  
591  
592  
593

- 594 Gyeongman Kim, Hajin Shim, Hyunsu Kim, Yunjey Choi, Junho Kim, and Eunho Yang. Diffu-  
595 sion video autoencoders: Toward temporally consistent face video editing via disentangled video  
596 encoding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recog-  
597 nition*, pp. 6091–6100, 2023.
- 598 Peipei Li, Rui Wang, Huaibo Huang, Ran He, and Zhaofeng He. Pluralistic aging diffusion autoen-  
599 coder. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*,  
600 pp. 22613–22623, October 2023.
- 601 Xiaoming Li, Shiguang Zhang, Shangchen Zhou, Lei Zhang, and Wangmeng Zuo. Learning dual  
602 memory dictionaries for blind face restoration. *IEEE Transactions on Pattern Analysis and Ma-  
603 chine Intelligence*, 2022.
- 604 Zeqi Li, Ruowei Jiang, and Parham Aarabi. Continuous face aging via self-estimated residual age  
605 embedding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recog-  
606 nition*, pp. 15008–15017, 2021.
- 607 Jae Hyun Lim and Jong Chul Ye. Geometric gan. *arXiv preprint arXiv:1705.02894*, 2017.
- 608 Ling Lin, Hao Liu, Jinqiao Liang, Zhendong Li, Jiao Feng, and Hu Han. Consensus-agent deep  
609 reinforcement learning for face aging. *IEEE Transactions on Image Processing*, 2024.
- 610 Lu Liu, Haibo Yu, Shenghui Wang, Lili Wan, and Shanshan Han. Learning shape and texture  
611 progression for young child face aging. *Signal Processing: Image Communication*, 93:116127,  
612 2021.
- 613 Arsha Nagrani, Joon Son Chung, and Andrew Zisserman. Voxceleb: a large-scale speaker identifi-  
614 cation dataset. *arXiv preprint arXiv:1706.08612*, 2017.
- 615 Niranjan D Narvekar and Lina J Karam. A no-reference perceptual image sharpness metric based  
616 on a cumulative probability of blur detection. In *2009 International Workshop on Quality of  
617 Multimedia Experience*, pp. 87–91. IEEE, 2009.
- 618 Niranjan D Narvekar and Lina J Karam. A no-reference image blur metric based on the cumulative  
619 probability of blur detection (cpbd). *IEEE Transactions on Image Processing*, 20(9):2678–2683,  
620 2011.
- 621 Roy Or-El, Soumyadip Sengupta, Ohad Fried, Eli Shechtman, and Ira Kemelmacher-Shlizerman.  
622 Lifespan age transformation synthesis. In *Computer Vision–ECCV 2020: 16th European Confer-  
623 ence, Glasgow, UK, August 23–28, 2020, Proceedings, Part VI 16*, pp. 739–755. Springer, 2020.
- 624 Konpat Preechakul, Nattanat Chatthee, Suttisak Wizadwongsa, and Supasorn Suwajanakorn. Dif-  
625 fusion autoencoders: Toward a meaningful and decodable representation. In *Proceedings of the  
626 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10619–10629, 2022.
- 627 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
628 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
629 models from natural language supervision. In *International conference on machine learning*, pp.  
630 8748–8763. PMLR, 2021.
- 631 Fitsum Reda, Janne Kontkanen, Eric Tabellion, Deqing Sun, Caroline Pantofaru, and Brian Curless.  
632 Film: Frame interpolation for large motion. In *European Conference on Computer Vision*, pp.  
633 250–266. Springer, 2022.
- 634 Daniel Roich, Ron Mokady, Amit H Bermano, and Daniel Cohen-Or. Pivotal tuning for latent-based  
635 editing of real images. *ACM Transactions on graphics (TOG)*, 42(1):1–13, 2022.
- 636 Rasmus Rothe, Radu Timofte, and Luc Van Gool. Dex: Deep expectation of apparent age from a  
637 single image. In *Proceedings of the IEEE international conference on computer vision workshops*,  
638 pp. 10–15, 2015.
- 639 Rasmus Rothe, Radu Timofte, and Luc Van Gool. Deep expectation of real and apparent age from  
640 a single image without facial landmarks. *International Journal of Computer Vision*, 126(2-4):  
641 144–157, 2018.

- 648 Masaki Saito, Shunta Saito, Masanori Koyama, and Sosuke Kobayashi. Train sparsely, generate  
649 densely: Memory-efficient unsupervised training of high-resolution temporal gan. *International*  
650 *Journal of Computer Vision*, 128(10-11):2586–2606, 2020.
- 651
- 652 Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for  
653 simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*, 2014.
- 654
- 655 Yu Tian, Jian Ren, Menglei Chai, Kyle Olszewski, Xi Peng, Dimitris N Metaxas, and Sergey  
656 Tulyakov. A good image generator is what you need for high-resolution video synthesis. *arXiv*  
657 *preprint arXiv:2104.15069*, 2021.
- 658
- 659 Omer Tov, Yuval Alaluf, Yotam Nitzan, Or Patashnik, and Daniel Cohen-Or. Designing an encoder  
660 for stylegan image manipulation. *ACM Transactions on Graphics (TOG)*, 40(4):1–14, 2021.
- 661
- 662 Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. Mocogan: Decomposing motion  
663 and content for video generation. In *Proceedings of the IEEE/CVF Conference on Computer*  
664 *Vision and Pattern Recognition*, pp. 1526–1535, 2018.
- 665
- 666 Rotem Tzaban, Ron Mokady, Rinon Gal, Amit Bermano, and Daniel Cohen-Or. Stitch it in time:  
667 Gan-based facial editing of real videos. In *SIGGRAPH Asia 2022 Conference Papers*, pp. 1–9,  
668 2022.
- 669
- 670 Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating videos with scene dynamics.  
671 *Advances in neural information processing systems*, 29, 2016.
- 672
- 673 Junaid Wahid, Fangneng Zhan, Pramod Rao, and Christian Theobalt. Diffage3d: Diffusion-based  
674 3d-aware face aging. *arXiv preprint arXiv:2408.15922*, 2024.
- 675
- 676 Ting-Chun Wang, Arun Mallya, and Ming-Yu Liu. One-shot free-view neural talking-head synthesis  
677 for video conferencing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
678 *Pattern Recognition*, pp. 10039–10049, 2021a.
- 679
- 680 Xintao Wang, Yu Li, Honglun Zhang, and Ying Shan. Towards real-world blind face restoration  
681 with generative facial prior. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
682 *and Pattern Recognition*, 2021b.
- 683
- 684 Zongwei Wang, Xu Tang, Weixin Luo, and Shenghua Gao. Face aging with identity-preserved  
685 conditional generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on*  
686 *Computer Vision and Pattern Recognition*, pp. 7939–7947, 2018.
- 687
- 688 Yanze Wu, Xintao Wang, Gen Li, and Ying Shan. Animesr: learning real-world super-resolution  
689 models for animation videos. *Advances in Neural Information Processing Systems*, 35:11241–  
690 11252, 2022.
- 691
- 692 Liangbin Xie, Xintao Wang, Honglun Zhang, Chao Dong, and Ying Shan. Vfhq: A high-quality  
693 dataset and benchmark for video face super-resolution. In *The IEEE Conference on Computer*  
694 *Vision and Pattern Recognition Workshops (CVPRW)*, 2022.
- 695
- 696 Shuai Yang, Liming Jiang, Ziwei Liu, and Chen Change Loy. Styleganex: Stylegan-based manipu-  
697 lation beyond cropped aligned faces. *arXiv preprint arXiv:2303.06146*, 2023.
- 698
- 699 Xu Yao, Alasdair Newson, Yann Gousseau, and Pierre Hellier. A latent transformer for disentangled  
700 face editing in images and videos. In *Proceedings of the IEEE/CVF international conference on*  
701 *computer vision*, pp. 13789–13798, 2021a.
- 702
- 703 Xu Yao, Gilles Puy, Alasdair Newson, Yann Gousseau, and Pierre Hellier. High resolution face age  
704 editing. In *2020 25th International conference on pattern recognition (ICPR)*, pp. 8624–8631.  
705 IEEE, 2021b.
- 706
- 707 Fei Yin, Yong Zhang, Xiaodong Cun, Mingdeng Cao, Yanbo Fan, Xuan Wang, Qingyan Bai,  
708 Baoyuan Wu, Jue Wang, and Yujiu Yang. Styleheat: One-shot high-resolution editable talking  
709 face generation via pre-trained stylegan. In *European conference on computer vision*, pp. 85–101.  
710 Springer, 2022.

- 702 Bowen Zhang, Chenyang Qi, Pan Zhang, Bo Zhang, HsiangTao Wu, Dong Chen, Qifeng Chen,  
703 Yong Wang, and Fang Wen. Metaportrait: Identity-preserving talking head generation with fast  
704 personalized adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
705 *Pattern Recognition*, pp. 22096–22105, 2023a.
- 706 Guozhen Zhang, Yuhan Zhu, Haonan Wang, Youxin Chen, Gangshan Wu, and Limin Wang. Ex-  
707 tracting motion and appearance via inter-frame attention for efficient video frame interpolation.  
708 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.  
709 5682–5692, 2023b.
- 710 Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable  
711 effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE/CVF Conference*  
712 *on Computer Vision and Pattern Recognition*, pp. 586–595, 2018.
- 713 Yuxuan Zhang, Xin Wang, M Saad Shakeel, Hao Wan, and Wenxiong Kang. Learning upper patch  
714 attention using dual-branch training strategy for masked face recognition. *Pattern Recognition*,  
715 126:108522, 2022.
- 716 Tianyue Zheng, Weihong Deng, and Jiani Hu. Cross-age lfw: A database for studying cross-age  
717 face recognition in unconstrained environments. *arXiv preprint arXiv:1708.08197*, 2017.
- 718 Chang Zhou, Jie Liu, Jie Tang, and Gangshan Wu. Video frame interpolation with densely queried  
719 bilateral correlation. *arXiv preprint arXiv:2304.13596*, 2023.
- 720 Hao Zhu, Wayne Wu, Wentao Zhu, Liming Jiang, Siwei Tang, Li Zhang, Ziwei Liu, and  
721 Chen Change Loy. CelebV-HQ: A large-scale video facial attributes dataset. In *ECCV*, 2022.
- 722 Jingyuan Zhu, Huimin Ma, Jiansheng Chen, and Jian Yuan. Motionvideogan: A novel video gen-  
723 erator based on the motion space learned from image pairs. *IEEE Transactions on Multimedia*,  
724 2023.
- 725 Gaspard Zoss, Prashanth Chandran, Eftychios Sifakis, Markus Gross, Paulo Gotardo, and Derek  
726 Bradley. Production-ready face re-aging for visual effects. *ACM Transactions on Graphics (TOG)*,  
727 41(6):1–12, 2022.
- 728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755