Now You See Me, Now You Don't: A Unified Framework for Expression Consistent Anonymization in Talking Head Videos

Anil Egin Inria Sophia Antipolis Valbonne/France

Andrea Tangherloni Universita Bocconi Milan/Italy

anil.eqin@studbocconi.it

andrea.tangherloni@unibocconi.it

Antitza Dantcheva Inria Sophia Antipolis Valbonne/France

antitza.dantcheva@inria.fr

Abstract

Face video anonymization is aimed at privacy preservation while allowing for the analysis of videos in a number of computer vision downstream tasks such as expression recognition, people tracking, and action recognition. We propose here a novel unified framework referred to as Anon-NET, streamlined to de-identify facial videos, while preserving age, gender, race, pose, and expression of the original video. Specifically, we inpaint faces by a diffusion-based generative model guided by high-level attribute recognition and motion-aware expression transfer. We then animate deidentified faces by video-driven animation, which accepts the de-identified face and the original video as input. Extensive experiments on the datasets VoxCeleb2, CelebV-HQ, and HDTF, which include diverse facial dynamics, demonstrate the effectiveness of AnonNET in obfuscating identity while retaining visual realism and temporal consistency. The code of AnonNet will be publicly released.

1. Introduction

Video anonymization is aimed at effectively obscuring identity-information, such as faces or voices, without compromising the integrity or usability of the content. Such anonymization has been fueled by ethical, legal or practically necessity - increasingly essential in a world, where visual data has become omnipresent. For instance, medical therapy sessions recorded for research require video anonymization to protect patient identities, particularly facial features, while preserving related expressions and emotions, which are pertinent for research.

In addition, legal frameworks such as the General Data



Figure 1. Qualitative comparison between original and anonymized face pairs. Left: CelebA-HQ. Right: LFW. Each row shows original/anonymized pairs.

Protection Regulation (*GDPR*)¹ impose strict constraints on collection, processing, and dissemination of personal data, including biometric identifiers such as images and videos of the human face. More recently, the European Union's Artificial Intelligence Act (*AI Act*))², adopted in 2024, introduced a tiered risk-based framework that places heightened scrutiny on AI systems handling biometric and

¹https://gdpr-info.eu

²https://artificialintelligenceact.eu

identity-sensitive data. These evolving regulations reinforce the demand for anonymization methods that ensure privacy protection, while preserving the utility of data for downstream computer vision tasks. Traditional approaches including pixelation, blurring, and masking often degrade video quality and compromise associated applicability in such tasks [28].

Modern deep generative models, especially Generative Adversarial Networks (GANs), have significantly advanced visual realism. *Image anonymization* inpainting-based approaches, such as DeepPrivacy [16] and DeepPrivacy2 [15] replace only sensitive facial regions, thereby preserving the surrounding content of the facial area. Conversely, fully synthetic pipelines such as FALCO [2] generate artificial facial images, while preserving high-level attributes including age, gender and race, however may introduce inconsistencies in expression or struggle with robustness under diverse pose and lighting conditions due to reliance on GAN-inversion and matching in a synthetic latent space.

W.r.t. *video anonymization*, face-swapping methods [47] can inadvertently preserve identity-specific features, compromising unlinkability [32]. Finally, landmark-based motion transfer techniques risk motion artifacts, in cases when tracking is imperfect [29].

Motivated by the above, in this work we propose *Anon-NET*, a multi-stage framework that (a) *synthesizes new identities*, while preserving facial attributes. Our approach employs diffusion-based inpainting guided by structural priors for comprehensive identity obfuscation, avoiding the limitations of reference-based methods. Further, (b) a landmark-free motion transfer module ensures realistic expressions without relying on explicit keypoint tracking, thereby mitigating alignment fail cases. By restricting modifications to the face region, AnonNET retains original scene context, offering high-quality, privacy-preserving video anonymization.

The main contributions of our work include the following.

- A *novel multi-stage framework* for video anonymization that synthesizes new facial identities while preserving age, gender, race, and expressions.
- A new dataset of anonymized videos pertained to Vox-Celeb, CelebV, and HDTF datasets, providing a valuable resource for future research in areas such as deepfake detection.
- A comprehensive evaluation of our pipeline against stateof-the-art models on image level, as well as with regard to the re-identification, identity consistency, and expressionaware downstream utility on video level.

2. Related Work

2.1. Image Anonymization

Traditional Techniques include pixelation, blurring, as well as masking obscure facial features, in order to hinder identity recognition, preserving general image context. However, such techniques can significantly degrade images, impeding tasks such as expression analysis or attribute prediction [3, 18]. In addition, prior work has demonstrated that such anonymization techniques can be partially reversible [7], raising concerns about related robustness in adversarial settings and the potential for identity leakage, rendering such methods insufficient for scenarios that demand both, privacy protection and downstream utility.

Adversarial Techniques employ adversarial training, balancing image anonymization and utility preservation within a min-max framework. Nasr et al. [24] used adversarial regularization to defend against membership inference, while Wu et al. [39, 40] leveraged GANs to deidentify faces without compromising action recognition. Nonetheless, directly training GANs in the image space is challenging w.r.t. fine-grained preservation of expressions and background details due to high pixel-space complexity.

To address these challenges, some methods manipulate latent representations [21], disentangling identity-specific traits from other image-features. Nevertheless, such approaches may inadvertently preserve cues correlated with identity or require accurate facial landmarks or segmentation, which can be error-prone [13, 22, 31]. DeepPrivacy2 [15] extended a guided GAN framework, in order to anonymize full-body images, relying on precise facial segmentation and landmark detection. Failure cases in the latter can compromise both, anonymity and image quality.

Diffusion-Based Approaches iteratively refine noisy inputs, employing U-Net-like architectures to anonymize faces [19, 20, 25]. Such methods often apply constrained transforms for obfuscation, yielding convincing results. Related limitations include the lack of control of diverse attributes such as age, gender, race, and expressions. Our approach similarly uses diffusion-based inpainting, however conditions it on user-specified attributes, granting finegrained control over face appearance, while preserving nonidentity aspects of the scene. Therefore, AnonNET enables flexible anonymization that balances privacy with visual fidelity and enables tasks such as emotion recognition or pose estimation.

2.2. Motion Transfer

Motion transfer focuses on generating videos, with an appearance stemming from a reference image, guided by video-motion, *e.g.*, head pose, expressions [12, 42]. Early works used explicit keypoints or 3D models [35], however experiencing tracking errors in complex motions. More re-

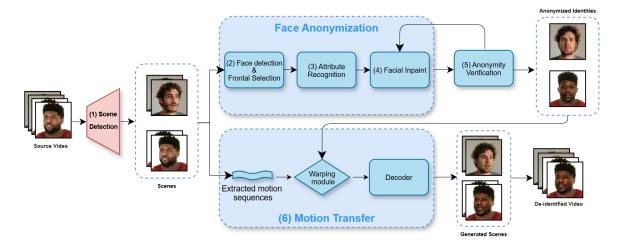


Figure 2. Overview of our multi-stage anonymization AnonNET-pipeline. (1) Scene changes are detected, and identities are tracked. (2) Faces are detected and single frontal frame is selected per scene identity. (3) Facial attributes are recognized. (4) A diffusion-based model inpaints the masked face. (5) Current anonymity is evaluated. (6) Landmark-free motion transfer reintroduces natural head movement. (7) Frames are reassembled for a coherent output video.

cent approaches learn to warp or synthesize motion in the latent space, reducing reliance on landmark detection [10].

Non-Diffusion techniques like FOMM [30], MRAA [46], SAFA [33], or Face vid2vid [34] detect implicit keypoints to warp source images, often faster, however lacking background stitching or fine eye/lip retargeting. They also are challenged under extreme head rotations or multiple faces.

Landmark-free solutions synthesize animation in the latent space, circumventing explicit structural constraints. For instance, the *Latent Image Animator* LIA [36] builds a latent motion dictionary, whereas LivePortrait [10] employs stitching and retargeting for flexible portrait animation. Such methods capture subtle expressions without the necessity of structural information. We note extreme motion or occlusion remain challenging.

Additional Models such as FADM [46], Face Adapter [11], AniPortrait [37], X-Portrait [41], and MegActor [44] achieve high-fidelity facial reenactment, often at the cost of increased computational complexity. These approaches frequently rely on 3D priors or explicit landmark extraction, rendering them less suitable for large-scale or real-time applications.

In contrast, AnonNET prioritizes computational efficiency, avoiding the necessity for 3D reconstruction or facial landmarks. This enables efficient anonymization, allowing for processing of longer video sequences.

2.3. Video Anonymization

RID-TWIN [23] used BLIP for face captioning and stable diffusion jointly with MediaPipe-based segment extraction. However, it focused on automatic de-identification rather

than preserving specific attributes like age, gender, race, or expressions. SAFA was leveraged [33] for head motion transfer, which is prone to landmark-based errors. SAFA used self-supervised landmark-like keypoints, however associated low resolution leads to motion and appearance artifacts.

AI Stylization [43] constitutes a perceptual approach for anonymization, replacing facial realism with *artistic* abstraction. After an initial facial feature randomization stage, the method applied cubist and painterly stylizations to anonymized faces, aiming to preserve emotional salience and enhance viewer empathy. However, the approach sacrificed photorealism entirely, as renderings tend to be stylized and visually inconsistent across frames. Structural fidelity is not maintained, and lack of temporal coherence and realism limits its use cases.

In contrast, our AnonNET-pipeline conditions diffusionbased inpainting on user-defined attribute priors, enabling consistent preservation of essential characteristics. Specifically, we adopt the *landmark-free* motion transfer framework LIA and LivePortrait, in order to transfer expression and head pose. Additionally, we systematically detect *scene changes* for improved temporal consistency and seamless anonymization across longer video sequences, therefore accommodating complex multi-scene videos effectively.

3. Method

3.1. Overview of the Proposed Framework

We propose a *multi-stage* pipeline (see Figure 2), streamlined to obfuscate identity, while preserving attributes such as age, gender, race, and at the same time ensuring tempo-

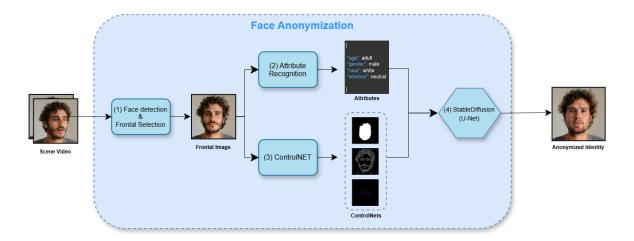


Figure 3. Overview of the expression-consistent face anonymization module in AnonNET. (1) A face detection and frontal selection stage extracts a frame containing a frontal face from the input scene video. Then this frontal image is processed in parallel by two branches: (2) An attribute recognition module that infers semantic attributes such as age, gender, race, and expression; (3) a ControlNet module, which extracts structural guidance (*e.g.*, face mask, lineart, pose) for conditioning the generative model. (4) Stable Diffusion based on U-Net synthesizes an anonymized face conditioned on both, extracted attributes and ControlNet features.

ral consistency represented by expression and head poses. In particular, our AnonNET includes following stages.

- (1) Scene Detection & Identity Clustering. We segment the input video via FFmpeg-based scene change detection, and then cluster faces across scenes using VGG-Face2 [5] embeddings and cosine-distance thresholds. Scenes containing the same individual share a consistent anonymized identity throughout the video.
- (2) Face Detection & Frontal Selection. RetinaFace [9] localizes the face region. A single representative frame, associated to a frontal pose, is selected per scene. This reduces flickering and computational burden by focusing anonymization on a single frame per segment.
- (3) Attribute Recognition. We estimate age, gender, race, and emotion via DeepFace [27], retaining high-level features that do not reveal identity.
- (4) **Diffusion-Based Inpainting.** We adopt Realistic Vision V5.0 to inpaint the masked face guided by **Control-Nets** for segmentation mask, lineart, and openpose maintain structural fidelity along with **Attribute-Conditioned Prompt** that we provide key attributes while discarding identity-specific details.
- (5) Anonymity Verification. To ensure identity obfuscation and avoid leakage at the end of the process, we include a verification module that verifies the cosine similarity between VGG-Face2 [5] embeddings of the original and anonymized images. In case that the similarity exceeds a threshold, the inpainting is re-triggered with higher stochasticity to enforce stronger anonymization.
- **(6) Landmark-Free Motion Transfer.** We select the frameworks **LIA** and **LivePortrait** to warp the anonymized

face per frame, replicating natural head movements from the original video, without explicit landmark tracking. This approach mitigates flickering and alignment errors.

(7) Video Reassembly. Processed frames are merged back into the original scene structure, retaining audio and background context.

We note that resulting videos *preserve facial attributes*, *exhibit strong identity obfuscation, and temporal coherence*, see Supplementary Material for videos. Our modular design allows for each stage, namely scene detection, face detection, inpainting, motion transfer to be independently improved or replaced. We proceed to elaborate on each stage.

3.2. Video Preprocessing and Scene Detection

- (1) Scene Change Detection. We detect coarse scene boundaries using FFmpeg's scene change filter (select='gt(scene,X)'), which flags transitions based on frame-wise histogram differences. To avoid over-segmentation, we refine these segments by computing mean RGB differences and merging visually similar intervals under a shared scene_id. This simple yet effective two-stage strategy yields stable scene partitions and reduces redundancy, enabling consistent identity tracking and efficient anonymization across temporally coherent regions.
- (2) Frontal Selection For each scene, we select a *single* representative frame to serve as the anchor for anonymization and motion transfer. To ensure full facial coverage and minimize downstream hallucination, we prioritize frames

with a frontal head pose, where both geometric structure and semantic attributes are fully visible.

We estimate head pose using the face_alignment [4] library in 2D landmark mode, extracting 68 facial landmarks per frame. A subset of six key points (nose tip, chin, eye corners, mouth corners) is selected and matched to a predefined 3D face model. We then solve the perspective-n-point (PnP) problem via OpenCV's solvePnP, computing the 3D rotation vector of the head. Frames with fewer than 80% of landmarks falling within image boundaries are discarded. Among valid candidates, the frame with the smallest absolute pitch and yaw is selected as the scene frontal frame.

This selection strategy ensures that identity obfuscation operates on a complete and unobstructed face. Since motion transfer is applied post-anonymization, any missing or occluded facial regions in the frontal frame would otherwise be synthesized without constraint—potentially leading to artifacts or semantic drift.

3.3. Face Detection and Attribute-Guided Prompt Generation

(3) Attribute Recognition. The localized face is then passed to the DeepFace library [27], in order to extract coarse demographic and affective attributes, including age, gender, race, and emotion. These attributes are used to guide the anonymization process in a non-identifying manner.

3.4. Diffusion-Based Inpainting

(4) **Diffusion-Based Inpainting.** Towards conditioning the diffusion-based inpainting, we construct a descriptive prompt, encoding the extracted attributes. Additionally, we apply a negative prompt to suppress undesired artifacts such as distortions, unrealistic textures, or cartoon-like features. The final prompt directs the generation toward a photorealistic, high-fidelity identity with preserved semantic traits. An example prompt is:

```
A photorealistic portrait of a middle-aged Asian female, with a neutral expression.
```

We perform identity obfuscation via latent diffusion-based inpainting using Realistic Vision V5.0, a publicly available checkpoint based on Stable Diffusion v1.5 [26].

ControlNet applies structural guidance from the head mask to confine synthesis to the face region while preserving surrounding content.

Lineart & OpenPose ControlNets provide edge and pose priors to enforce geometric and expression fidelity during generation.

We use the DPMSolver++ scheduler for efficient denoising, typically over 20--70 steps with a guidance scale between 8--20, depending on dataset characteristics. A VAE

with perceptual reconstruction loss is used to map images between pixel and latent space, supporting visually coherent and detailed inpainting.

3.5. Anonymity Verification and Motion Transfer

- (5) Anonymity Verification. To ensure successful identity obfuscation, we capture the cosine similarity between the VGG-Face2 embeddings of original and anonymized images. In case that the cosine distance score is below a threshold of 0.3, indicating potential identity leakage, we re-trigger the inpainting process. In this second pass, we introduce greater stochasticity by increasing the prompt guidance scale and reducing ControlNet conditioning strength. Additionally, we extend the number of denoising steps by 5, in order to allow the diffusion process to deviate further from the original identity while still maintaining attribute consistency.
- **(6) Motion Transfer.** We combine landmark-free frameworks (**LIA**, **LivePortrait**) with scene stitching and eye/lip retargeting to animate anonymized faces across frames:
- 1. **Encoding:** The original (non-anonymized) source frame is encoded to obtain a latent motion representation $z_{s \to r}$, capturing pose and expression dynamics.
- 2. **Framewise Motion Codes:** For each target frame, motion offsets $\mathbf{w}_{r \to d}$ are predicted based on pose and expression changes.
- 3. Flow Field Synthesis: The anonymized source frame image is encoded, and the learned flow map (from $z_{s \to r} + \mathbf{w}_{r \to d}$) is used to warp it, transferring motion to the anonymized face, while preserving appearance.
- 4. **Refinement:** LivePortrait optionally enhances eye and lip dynamics (*e.g.*, blinks, speech) for improved realism. By decoupling motion from identity and operating in latent and flow-guided spaces, these approaches avoid explicit landmark detection, reduce artifacts, and increase robustness to rapid motion and partial occlusion, while maintaining temporal consistency.
- (5) Video Reassembly. Finally, we reassemble processed scenes into the final video. Specifically, each anonymized segment is integrated using original timestamps, ensuring proper alignment. We note that the audio is unaltered and resynchronized to preserve speech and background sounds.

Scenes with no faces remain untouched; multi-person scenarios are deferred for future work due to the complexities of simultaneous multi-face motion transfer.

4. Results

4.1. Experimental Setup

We proceed to comprehensively evaluate AnonNET's performance and compare related results to state-of-the-art anonymization methods, providing quantitative analysis on identity obfuscation and attribute retention. Further, an ablation study illustrates the impact of core components.

4.1.1. Datasets

We evaluate our framework on following two widely used *image* datasets.

- CelebA-HQ [17] contains 30,000 high-resolution face images annotated with 40 facial attributes, including age, gender, race, and appearance traits. The diversity in pose and lighting allow for our evaluation on attributepreserving anonymization.
- LFW [14] includes 13,233 images of 5,749 identities captured in unconstrained conditions. We focus on identity obfuscation and generalization under varying image quality and occlusions.

For *video*-based anonymization, we additionally use following datasets.

- CelebV-HQ [49] constitutes a curated high-resolution video face dataset.
- **VoxCeleb2** [8] represents a subset of 50,000 clips featuring over thousands of identities in varied conditions.
- **HDTF** [48] comprises expressive head motion and fine-grained lip sync.

4.1.2. Comparative Methods

We compare AnonNET against following *image* anonymization frameworks. DeepPrivacy2 [15] is prominent for removing identity cues, while preserving contextual and structural consistency. At the same time CIAGAN [22] represents a competitive former approach that modifies latent identity features while retaining key facial attributes.

W.r.t. video-anonymization, we compare AnonNET to RID-TWIN [23] that constitutes a recent end-to-end video anonymization pipeline with temporal coherence.

4.1.3. Implementation Details

AnonNET integrates a multi-stage pipeline with tailored configurations per component:

Anonymization: Based on Realistic Vision V5.0 with:

- Denoising steps: 20–70 (dataset-specific tuning)
- Guidance scale: 8–20 (for prompt expressiveness)
- ControlNets: Segmentation, LineArt, and OpenPose
- DPMSolver Scheduler for accelerated sampling

Motion Transfer: Landmark-free video animation via LivePortrait [10] and LIA [35], enabling smooth expression preservation across frames.

Computational Time. Anonymizing the 50,000-clip VoxCeleb2 subset takes approximately 160 hours with LIA and 185 hours with LivePortrait, using a single A100 GPU.

4.2. Evaluation Metrics

We adopt an evaluation framework, aimed at assessing identity obfuscation, attribute retention, and visual quality.

Re-identification Rate (Re@1). To evaluate identity leakage, we measure the rank-1 re-identification accuracy using face embeddings extracted from VGGFace2 [5] and CASIA-WebFace [45] models. For each anonymized image, we compute its embedding and retrieve the closest match from the original image set based on cosine similarity. A sample is considered successfully re-identified if its nearest neighbor corresponds to the same identity as the original. The final Re@1 score is computed as the ratio of correctly re-identified samples over the total number of anonymized images. Lower scores indicate stronger identity obfuscation.

Image Quality and Aesthetics. We use the **Q-Align/One-Align** [38] metric to estimate perceptual quality (Qual) and visual appeal (Aes). These scores are averaged across all images and videos.

Pose and Gaze Preservation. For each image, facial landmarks are first localized with MTCNN. Pose angles (pitch/yaw) are extracted via Dlib's face pose estimator. Gaze direction is then evaluated using **L2CS-Net** [1], computing the mean absolute error (MAE) between original and anonymized outputs.

Expression Preservation. Expression labels are predicted pre- and post-anonymization using the **DeepFace** library [27]. Accuracy is defined as the fraction of samples, where the predominant expression label is retained.

Temporal Identity Consistency (Video). For each anonymized video, we compute the average cosine distance of DINO [6] embeddings between consecutive frames, comparing with the same metric on original videos. This as-

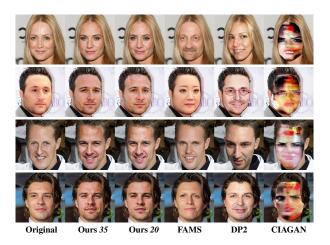


Figure 4. Qualitative face anonymization results pertained to the CelebA-HQ dataset. Each row corresponds to an input image (left column), and columns show outputs from various image-anonymization methods.

sesses intra-video consistency post-anonymization.

Encoding	CelebA-HQ		LFW	
	VGG↓	CASIA↓	VGG↓	CASIA↓
DeepPrivacy2 [15]	0.008	0.008	0.023	0.017
DeepPrivacy [16]	0.011	0.036	0.015	0.027
CIAGAN [22]	0.004	0.022	0.009	0.002
FALCO [2]	0.017	0.028	0.016	0.021
CAMOUFLaGE [25]	0.096	0.100	0.102	0.116
AnonNET(steps = 20)	0.073	0.031	0.056	0.039
AnonNET($steps = 35$)	0.041	0.017	0.042	0.027

Table 1. Re-identification rate employing VGGFace2 and CASIA pertaining to the CelebA-HQ [17] and LFW [14] datasets. Lower scores denote a lower similarity between anonymized and original images and are therefore better.

4.3. Quantitative Evaluation

Re-identification Performance. Table 1 reports rank-1 re-identification accuracy (Re@1) employing VGGFace2 and CASIA-WebFace embeddings on CelebA-HQ and LFW. While CIAGAN and DeepPrivacy2 achieve the lowest scores overall, AnonNET (35 steps) remains competitive, with Re@1 values of 0.041 (VGG) on CelebA-HQ and 0.042 on LFW. Compared to CIAGAN and DeepPrivacy, AnonNET provides a consistent drop in re-identification, while retaining attribute fidelity and performs favorably relative to recent diffusion-based anonymization baselines such as FALCO and CAMOUFLaGE.

Perceptual Quality and Aesthetic Appeal. As shown in Table 2, AnonNET outperforms all baselines on Q-Align quality and aesthetic scores across both datasets. W.r.t. CelebA-HQ, it achieves the highest quality (4.164) and aesthetics (3.332) scores, exceeding both, DeepPrivacy2 and ground truth. W.r.t. LFW, AnonNET similarly leads with 2.887 (Qual) and 1.939 (Aes), indicating strong generalization to unconstrained, low-resolution data. In contrast, CIA-GAN yields lower perceptual scores than both, AnonNET and DeepPrivacy2, consistent with degradation as noted in prior results. Beyond preserving visual realism, AnonNET systematically reduces artifacts and low-quality regions, owing to its attribute-guided diffusion design. These results are coherent and photorealistic outputs enhance compatibility with downstream tasks such as expression recognition or affective computing.

Trade-off. The above results confirm that AnonNET offers a favorable privacy—utility trade-off: while not achieving the absolute lowest Re@1, it delivers superior image quality and aesthetics, which are essential for downstream usability.



Figure 5. Qualitative comparison of original and anonymized frames pertained to the VoxCeleb2 dataset using AnonNET. Each column shows original/anonymized pairs.

Encoding	CelebA-HQ		LFW	
	Qual↑	Aes↑	Qual↑	Aes↑
GT	4.035	2.932	2.047	1.191
DeepPrivacy2	3.551	1.875	2.025	1.099
CIAGAN	1.011	1.361	1.006	1.466
AnonNET(steps = 20)	4.074	3.055	2.914	1.904
AnonNET($steps = 35$)	4.164	3.332	2.887	1.939

Table 2. Quality and aesthetics scores for anonymized images of CelebA-HQ and LFW.

Video-level Evaluation. Table 3 compares identity preservation, perceptual quality, and aesthetics between ground truth videos and those anonymized by AnonNET. Across all three datasets, AnonNET improves both, quality and aesthetics scores over the original videos while maintaining comparable levels of identity suppression. On CelebV-HQ and HDTF, our method achieves higher quality (*e.g.*, 4.153 versus 4.045 on HDTF) and aesthetics (*e.g.*, 3.021 versus 2.938), with only marginal differences in identity preservation.

W.r.t. VoxCeleb2, which encompasses low resolution and challenging visual settings in related videos, Anon-NET produces clean and coherently anonymized faces, raising quality from 2.493 to 2.859 and aesthetics from 1.606 to 1.949. These results highlight the robustness of our pipeline, even in unconstrained settings. Indeed, the new synthesized face improves structural integrity and overall visual consistency, rendering anonymized videos amenable to downstream tasks such as tracking or expression analysis.

Dataset		GT		Aı	nonNET	
	id_pres↓	qual ↑	aes ↑	id_pres↓	qual ↑	aes ↑
CelebV-HQ	0.011	3.800	2.718	0.013	3.907	2.886
VoxCeleb2	0.021	2.493	1.606	0.022	2.859	1.949
HDTF	0.008	4.045	2.938	0.007	4.153	3.021

Table 3. Comparison of identity preservation, quality, and aesthetics in videos for ground truth and AnonNET across datasets. Lower is better for identity preservation (id_pres); higher is better for quality (qual) and aesthetics (aes).

As shown in Figure 5, AnonNET reconstructs sharper







(a) Gender mismatch

(b) Expression mismatch

(c) Gaze inconsistency

Figure 6. Limitations of the proposed AnonNET framework. Each image shows a comparison between the original and anonymized version.

facial features and preserves expressions more consistently than the original VoxCeleb2 frames, which suffer from heavy blur and compression. This visual improvement, especially in motion-rich regions, renders previously unusable videos viable for downstream tasks such as expression analysis or video reenactment.

Table 4 shows that AnonNET, guided by OpenPose, achieves superior pose preservation and competitive gaze alignment. Even with fewer denoising steps, it outperforms all baselines, highlighting its efficiency and motion consistency.

Dataset	CelebA-HQ		
	Pose ↓	Gaze↓	
DeepPrivacy2	0.140	0.244	
FALCO	0.088	0.258	
FAMS[20]	0.048	0.161	
AnonNET(steps = 20)	0.014	0.187	
AnonNET(steps = 35)	0.015	0.172	

Table 4. Pose and gaze preservation (lower is better) on CelebA-HQ.

Tables 5 and 6 summarize AnonNET's performance on CelebA-HQ and LFW. our proposed method achieves nearperfect anonymization rates with no detection failures. Attribute preservation remains high across both datasets, particularly for gender and race, while expression accuracy is lower on LFW due to its greater variability and resolution constraints. These results confirm AnonNET's robustness across datasets with differing visual and demographic characteristics.

Anonymization Results				
Total images	30,000			
Successfully anonymized	29,997			
Anonymization failures	3			
Face detection failures	0			
Attribute Preservation				
Race (%)	79.5			
Gender (%)	99.4			
Age (mean \pm std)	(1.87,4.23)			
Expression (%)	74.7			

Table 5. Anonymization statistics and attribute preservation accuracy on the CelebA-HQ dataset.

Anonymization Results				
Total images	13,233			
Successfully anonymized	12,912			
Anonymization failures	321			
Face detection failures	0			
Attribute Preservation				
Race (%)	87.1			
Gender (%)	99.3			
Age (mean \pm std)	(2.69, 6.29)			
Expression (%)	52.9			

Table 6. Anonymization statistics and attribute preservation accuracy on the LFW dataset.

Limitations. Figure 6 highlights failure cases of Anon-NET. Since attribute guidance relies on pretrained recognition networks, errors in gender, expression, or gaze estimation can propagate to the anonymized output. These limitations suggest the need for more robust or fine-tuned attribute predictors, especially for handling edge cases and underrepresented demographics.



Figure 7. Qualitative comparison of motion transfer models. Each row corresponds to following frames of a video, columns correspond Original, Live Portrait, and LIA, respectively.

5. Conclusions

In this work, we introduced AnonNET, a unified multistage framework for anonymizing talking head videos, placing emphasis on preserving key facial attributes. We presented extensive evaluations, demonstrating the ability of AnonNET to obfuscate identity obfuscation, while allowing for further analysis. As opposed to the state of the art, AnonNET is more robust to diverse poses, lighting conditions, and motion dynamics, rendering it suitable for realworld applications such as journalism, therapy, and human-computer interaction.

Future work will explore extending AnonNET to full-body video anonymization, incorporating audio-driven synchronization for improved lip consistency, and enhancing expression preservation in more dynamic conversational settings. These directions aim to further broaden the applicability of anonymized video content in sensitive or privacy-critical domains.

References

- [1] Ahmed A. Abdelrahman, Thorsten Hempel, Aly Khalifa, and Ayoub Al-Hamadi. L2cs-net: Fine-grained gaze estimation in unconstrained environments. 2022. 6
- [2] Simone Barattin, Christos Tzelepis, Ioannis Patras, and Nicu Sebe. Attribute-preserving face dataset anonymization via latent code optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8001–8010, 2023. 2, 7
- [3] M. Barrett et al. Face anonymization: A comparative study of traditional methods. *Journal of Privacy and Security*, 4: 123–145, 2017.
- [4] Adrian Bulat and Georgios Tzimiropoulos. How far are we from solving the 2d & 3d face alignment problem? (and a dataset of 230,000 3d facial landmarks). In *International Conference on Computer Vision*, 2017. 5
- [5] Qiong Cao, Li Shen, Weidi Xie, Omkar M Parkhi, and Andrew Zisserman. Vggface2: A dataset for recognising faces across pose and age. In 2018 13th IEEE international conference on automatic face & gesture recognition (FG 2018), pages 67–74. IEEE, 2018. 4, 6
- [6] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In Proceedings of the IEEE/CVF international conference on computer vision, pages 9650–9660, 2021. 6
- [7] S. Choi et al. Blurring and masking for face obfuscation: A comprehensive review. Security and Privacy Research, 8: 233–250, 2021. 2
- [8] Joon Son Chung, Arsha Nagrani, and Andrew Zisserman. Voxceleb2: Deep speaker recognition. *arXiv preprint arXiv:1806.05622*, 2018. 6
- [9] Jiankang Deng, Jia Guo, Evangelos Ververas, Irene Kotsia, and Stefanos Zafeiriou. Retinaface: Single-shot multilevel face localisation in the wild. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5203–5212, 2020. 4
- [10] Jianzhu Guo, Dingyun Zhang, Xiaoqiang Liu, Zhizhou Zhong, Yuan Zhang, Pengfei Wan, and Di Zhang. Liveportrait: Efficient portrait animation with stitching and retargeting control. arXiv preprint arXiv:2407.03168, 2024. 3, 6
- [11] 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 *European Conference on Computer Vision*, pages 20–36. Springer, 2024. 3

- [12] Li Hu. Animate anyone: Consistent and controllable imageto-video synthesis for character animation. In *Proceedings of* the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 8153–8163, 2024. 2
- [13] Shengshan Hu, Xiaogeng Liu, Yechao Zhang, Minghui Li, Leo Yu Zhang, Hai Jin, and Libing Wu. Protecting facial privacy: Generating adversarial identity masks via style-robust makeup transfer. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 15014–15023, 2022. 2
- [14] Gary B Huang, Marwan Mattar, Tamara Berg, and Eric Learned-Miller. Labeled faces in the wild: A database forstudying face recognition in unconstrained environments. In Workshop on faces in 'Real-Life' Images: detection, alignment, and recognition, 2008. 6, 7
- [15] Håkon Hukkelås and Frank Lindseth. Deepprivacy2: Towards realistic full-body anonymization. In *Proceedings of* the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV), pages 1329–1338, 2023. 2, 6, 7
- [16] Håkon Hukkelås, Rudolf Mester, and Frank Lindseth. Deepprivacy: A generative adversarial network for face anonymization. In *International Symposium on Visual Computing (ISVC)*, pages 565–578. Springer, 2019. 2, 7
- [17] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. arXiv preprint arXiv:1710.10196, 2017. 6, 7
- [18] L. Kramer et al. Pixelation as a privacy-preserving technique in facial recognition systems. *International Journal of Pri*vacy Computing, 12:34–46, 2019. 2
- [19] Han-Wei Kung, Tuomas Varanka, Sanjay Saha, Terence Sim, and Nicu Sebe. Face anonymization made simple. *arXiv* preprint arXiv:2411.00762, 2024. 2
- [20] Han-Wei Kung, Tuomas Varanka, Sanjay Saha, Terence Sim, and Nicu Sebe. Face anonymization made simple. In *Proceedings of the Winter Conference on Applications of Computer Vision (WACV)*, pages 1040–1050, 2025. 2, 8
- [21] Minh-Ha Le and Niklas Carlsson. Styleid: Identity disentanglement for anonymizing faces. arXiv preprint arXiv:2212.13791, 2022. 2
- [22] Maxim Maximov, Ismail Elezi, and Laura Leal-Taixé. Ciagan: Conditional identity anonymization generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5447–5456, 2020. 2, 6, 7
- [23] Anirban Mukherjee, Monjoy Narayan Choudhury, and Dinesh Babu Jayagopi. Rid-twin: An end-to-end pipeline for automatic face de-identification in videos. arXiv preprint arXiv:2403.10058, 2024. 3, 6
- [24] Milad Nasr, Reza Shokri, and Amir Houmansadr. Machine learning with membership privacy using adversarial regularization. In *Proceedings of the 2018 ACM SIGSAC conference* on computer and communications security, pages 634–646, 2018. 2
- [25] Luca Piano, Pietro Basci, Fabrizio Lamberti, and Lia Morra. Latent diffusion models for attribute-preserving image anonymization. *arXiv preprint arXiv:2403.14790*, 2024. 2, 7

- [26] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695, 2022. 5
- [27] Sefik Serengil and Alper Ozpinar. A benchmark of facial recognition pipelines and co-usability performances of modules. *Journal of Information Technologies*, 17(2):95–107, 2024. 4, 5, 6
- [28] Yan Shoshitaishvili, Christopher Kruegel, and Giovanni Vigna. Portrait of a privacy invasion. *Proceedings on Privacy Enhancing Technologies*, 2015(1):41–60, 2015. 2
- [29] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. First order motion model for image animation. In *Neural Information Processing Systems* (NeurIPS), 2019. 2
- [30] Aliaksandr Siarohin, Oliver J Woodford, Jian Ren, Menglei Chai, and Sergey Tulyakov. Motion representations for articulated animation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13653–13662, 2021. 3
- [31] Qianru Sun, Liqian Ma, Seong Joon Oh, Luc Van Gool, Bernt Schiele, and Mario Fritz. Natural and effective obfuscation by head inpainting. In *Proceedings of the IEEE con*ference on computer vision and pattern recognition, pages 5050–5059, 2018. 2
- [32] Justus Thies, Michael Zollhofer, Marc Stamminger, Christian Theobalt, and Matthias Nießner. Neural texture synthesis and expression transfer for face video reenactment. In *IEEE Transactions on Visualization and Computer Graphics*, 2019. 2
- [33] Qiulin Wang, Lu Zhang, and Bo Li. Safa: Structure aware face animation. In 2021 International Conference on 3D Vision (3DV), pages 679–688. IEEE, 2021. 3
- [34] Ting-Chun Wang, Arun Mallya, and Ming-Yu Liu. One-shot free-view neural talking-head synthesis for video conferencing. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10039–10049, 2021. 3
- [35] Yaohui Wang, Di Yang, Francois Bremond, and Antitza Dantcheva. Latent image animator: Learning to animate images via latent space navigation. *arXiv preprint arXiv:2203.09043*, 2022. 2, 6
- [36] Yaohui Wang, Di Yang, Francois Bremond, and Antitza Dantcheva. Lia: Latent image animator. *IEEE Transac*tions on Pattern Analysis and Machine Intelligence, 46(12): 10829–10844, 2024. 3
- [37] Huawei Wei, Zejun Yang, and Zhisheng Wang. Aniportrait: Audio-driven synthesis of photorealistic portrait animation. arXiv preprint arXiv:2403.17694, 2024. 3
- [38] Haoning Wu, Zicheng Zhang, Weixia Zhang, Chaofeng Chen, Liang Liao, Chunyi Li, Yixuan Gao, Annan Wang, Erli Zhang, Wenxiu Sun, et al. Q-align: Teaching lmms for visual scoring via discrete text-defined levels. arXiv preprint arXiv:2312.17090, 2023. 6
- [39] Yifan Wu, Fan Yang, Yong Xu, and Haibin Ling. Privacyprotective-gan for privacy preserving face de-identification.

- Journal of Computer Science and Technology, 34:47–60, 2019. 2
- [40] Zhenyu Wu, Zhangyang Wang, Zhaowen Wang, and Hailin Jin. Towards privacy-preserving visual recognition via adversarial training: A pilot study. In *Proceedings of the Eu*ropean conference on computer vision (ECCV), pages 606– 624, 2018. 2
- [41] You Xie, Hongyi Xu, Guoxian Song, Chao Wang, Yichun Shi, and Linjie Luo. X-portrait: Expressive portrait animation with hierarchical motion attention. In ACM SIGGRAPH 2024 Conference Papers, pages 1–11, 2024. 3
- [42] Zhongcong Xu, Jianfeng Zhang, Jun Hao Liew, Hanshu Yan, Jia-Wei Liu, Chenxu Zhang, Jiashi Feng, and Mike Zheng Shou. Magicanimate: Temporally consistent human image animation using diffusion model. In *Proceedings of* the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 1481–1490, 2024. 2
- [43] Özge Nilay Yalçın, Vanessa Utz, and Steve DiPaola. Empathy through aesthetics: Using ai stylization for visual anonymization of interview videos. In *Proceedings of the 3rd Empathy-Centric Design Workshop: Scrutinizing Empathy Beyond the Individual*, pages 63–68, 2024. 3
- [44] Shurong Yang, Huadong Li, Juhao Wu, Minhao Jing, Linze Li, Renhe Ji, Jiajun Liang, and Haoqiang Fan. Megactor: Harness the power of raw video for vivid portrait animation. arXiv preprint arXiv:2405.20851, 2024. 3
- [45] Dong Yi, Zhen Lei, Shengcai Liao, and Stan Z Li. Learning face representation from scratch. arXiv preprint arXiv:1411.7923, 2014. 6
- [46] Bohan Zeng, Xuhui Liu, Sicheng Gao, Boyu Liu, Hong Li, Jianzhuang Liu, and Baochang Zhang. Face animation with an attribute-guided diffusion model. In *Proceedings of* the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 628–637, 2023. 3
- [47] W. Zhang, S. Shan, and X. Chen. Faceshifter: Towards high fidelity and occlusion aware face swapping. In *Proceed*ings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2020. 2
- [48] Zhimeng Zhang, Lincheng Li, Yu Ding, and Changjie Fan. Flow-guided one-shot talking face generation with a high-resolution audio-visual dataset. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3661–3670, 2021. 6
- [49] Hao Zhu, Wayne Wu, Wentao Zhu, Liming Jiang, Siwei Tang, Li Zhang, Ziwei Liu, and Chen Change Loy. Celebvhq: A large-scale video facial attributes dataset. In *European conference on computer vision*, pages 650–667. Springer, 2022. 6